MC504 Sistemas Operacionais

Armazenamento

Prof. Dr. Eng. Isaías Bittencourt Felzmann

isaias@ic.unicamp.br

Campinas, 2s/2025

Copyright Note

The following set of slides are copyright Silberschatz, Galvin and Gagne, 2018. Modifications were made for their use in conjunction with MC504. The original material is available at os-book.com.

Os direitos autorais do conjunto de slides a seguir pertence a Silberschatz, Galvin and Gagne, 2018. Foram feitas modificações para seu uso em MC504. O material original está disponível em <u>os-book.com</u>.



Chapter 11: Mass-Storage Systems

- Overview of Mass Storage Structure
- HDD Scheduling
- NVM Scheduling
- Error Detection and Correction
- Storage Device Management
- Swap-Space Management
- Storage Attachment
- RAID Structure





Objectives

- Describe the physical structure of secondary storage devices and the effect of a device's structure on its uses
- Explain the performance characteristics of mass-storage devices
- Evaluate I/O scheduling algorithms
- Discuss operating-system services provided for mass storage, including RAID



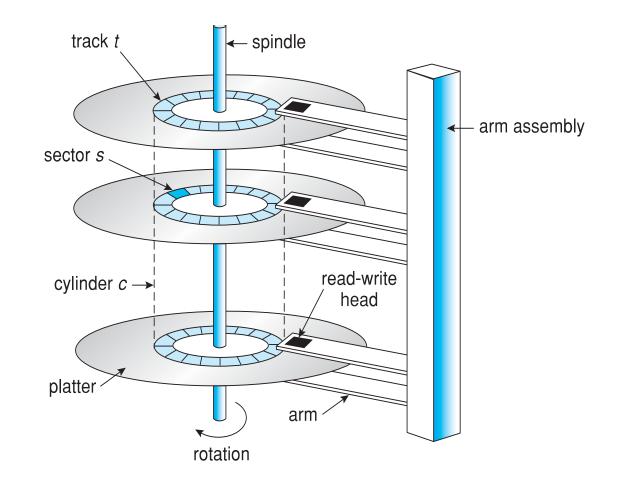
Overview of Mass Storage Structure

- Bulk of secondary storage for modern computers is hard disk drives
 (HDDs) and nonvolatile memory (NVM) devices
- HDDs spin platters of magnetically-coated material under moving read-write heads
 - Drives rotate at 60 to 250 times per second
 - Transfer rate is rate at which data flow between drive and computer
 - Positioning time (random-access time) is time to move disk arm to desired cylinder (seek time) and time for desired sector to rotate under the disk head (rotational latency)
 - Head crash results from disk head making contact with the disk surface -- That's bad
- Disks can be removable





Moving-head Disk Mechanism





Hard Disk Performance

- Access Latency = Average access time = average seek time + average latency
 - For fastest disk 3ms + 2ms = 5ms
 - For slow disk 9ms + 5.56ms = 14.56ms
- Average I/O time = average access time + (amount to transfer / transfer rate) + controller overhead
- For example to transfer a 4KB block on a 7200 RPM disk with a 5ms average seek time, 1Gb/sec transfer rate with a .1ms controller overhead =
 - 5ms + 4.17ms + 0.1ms + transfer time =
 - Transfer time = 4KB / 1Gb/s * 8Gb / GB * 1GB / 1024²KB = 32 / (1024²) = 0.031 ms
 - Average I/O time for 4KB block = 9.27ms + .031ms = 9.301ms





The First Commercial Disk Drive



1956
IBM RAMDAC computer included the IBM Model 350 disk storage system

5M (7 bit) characters 50 x 24" platters Access time = < 1 second





Nonvolatile Memory Devices

- If disk-drive like, then called solid-state disks (SSDs)
- Other forms include USB drives (thumb drive, flash drive), DRAM disk replacements, surface-mounted on motherboards, and main storage in devices like smartphones
- Can be more reliable than HDDs
- More expensive per MB
- Maybe have shorter life span need careful management
- Less capacity
- But much faster
- Busses can be too slow -> connect directly to PCI for example
- No moving parts, so no seek time or rotational latency





Nonvolatile Memory Devices

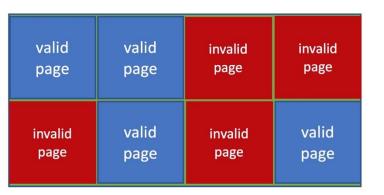
- Have characteristics that present challenges
- Read and written in "page" increments (think sector) but can't overwrite in place
 - Must first be erased, and erases happen in larger "block" increments
 - Can only be erased a limited number of times before worn out ~ 100,000
 - Life span measured in drive writes per day (DWPD)
 - A 1TB NAND drive with rating of 5DWPD is expected to have 5TB per day written within warrantee period without failing





NAND Flash Controller Algorithms

- With no overwrite, pages end up with mix of valid and invalid data
- To track which logical blocks are valid, controller maintains flash translation layer (FTL) table
- Also implements garbage collection to free invalid page space
- Allocates overprovisioning to provide working space for GC
- Each cell has lifespan, so wear leveling needed to write equally to all cells



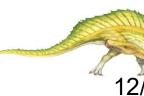
NAND block with valid and invalid pages





Volatile Memory

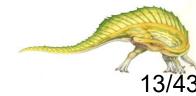
- DRAM frequently used as mass-storage device
 - Not technically secondary storage because volatile, but can have file systems, be used like very fast secondary storage
- RAM drives (with many names, including RAM disks) present as raw block devices, commonly file system formatted
- Computers have buffering, caching via RAM, so why RAM drives?
 - Caches / buffers allocated / managed by programmer, operating system, hardware
 - RAM drives under user control
 - Found in all major operating systems
 - Linux /dev/ram, macOS diskutil to create them, Linux
 /tmp of file system type tmpfs
- Used as high speed temporary storage
 - Programs could share bulk date, quickly, by reading/writing to RAM drive





Address Mapping

- Disk drives are addressed as large 1-dimensional arrays of logical blocks, where the logical block is the smallest unit of transfer
 - Low-level formatting creates logical blocks on physical media
- The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially
 - Sector 0 is the first sector of the first track on the outermost cylinder
 - Mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the cylinders from outermost to innermost
 - Logical to physical address should be easy
 - Except for bad sectors
 - Non-constant # of sectors per track via constant angular velocity





HDD Scheduling

- The operating system is responsible for using hardware efficiently —
 for the disk drives, this means having a fast access time and disk
 bandwidth
- Minimize seek time
- Seek time ≈ seek distance
- Disk bandwidth is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer

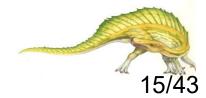


Disk Scheduling (Cont.)

- Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying "depth")
- Several algorithms exist to schedule the servicing of disk I/O requests
- The analysis is true for one or many platters
- We illustrate scheduling algorithms with a request queue (0-199)

98, 183, 37, 122, 14, 124, 65, 67

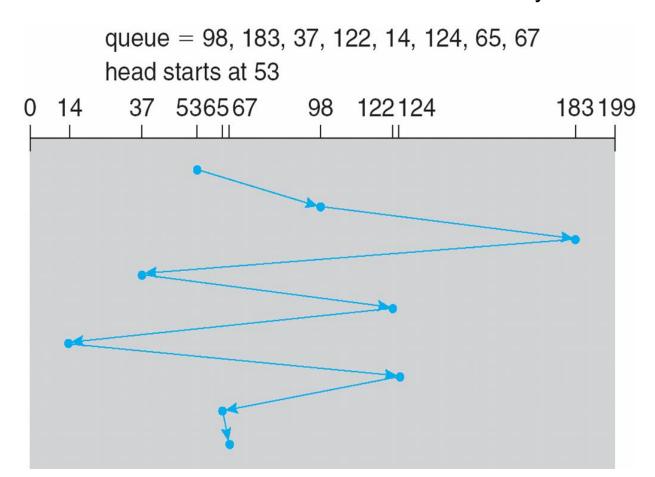
Head pointer 53





FCFS

Illustration shows total head movement of 640 cylinders



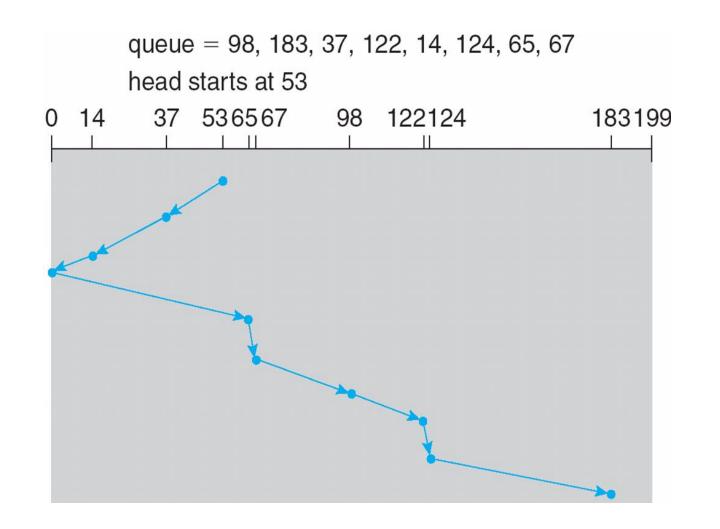


SCAN

- The disk arm starts at one end of the disk, and moves toward the other end, servicing requests until it gets to the other end of the disk, where the head movement is reversed and servicing continues.
- SCAN algorithm Sometimes called the elevator algorithm
- Illustration shows total head movement of 208 cylinders
- But note that if requests are uniformly dense, largest density at other end of disk and those wait the longest



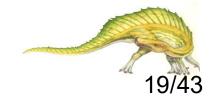
SCAN (Cont.)





C-SCAN

- Provides a more uniform wait time than SCAN
- The head moves from one end of the disk to the other, servicing requests as it goes
 - When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip
- Treats the cylinders as a circular list that wraps around from the last cylinder to the first one
- Total number of cylinders?

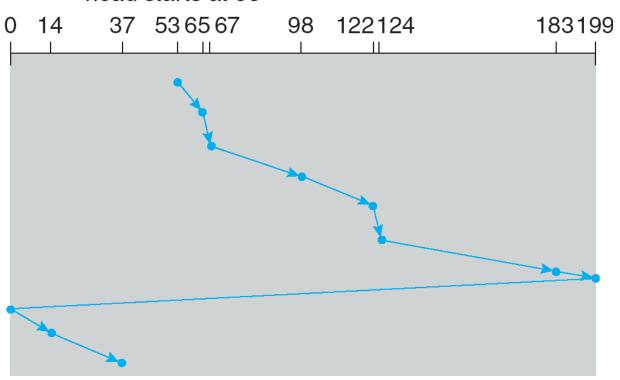




C-SCAN (Cont.)

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53





NVM Scheduling

- No disk heads or rotational latency but still room for optimization
- In RHEL 7 NOOP (no scheduling) is used but adjacent LBA requests are combined
 - NVM best at random I/O, HDD at sequential
 - Throughput can be similar
 - Input/Output operations per second (IOPS) much higher with NVM (hundreds of thousands vs hundreds)
 - But write amplification (one write, causing garbage collection and many read/writes) can decrease the performance advantage



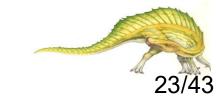
Storage Device Management

- Low-level formatting, or physical formatting Dividing a disk into sectors that the disk controller can read and write
 - Each sector can hold header information, plus data, plus error correction code (ECC)
 - Usually 512 bytes of data but can be selectable
- To use a disk to hold files, the operating system still needs to record its own data structures on the disk
 - **Partition** the disk into one or more groups of cylinders, each treated as a logical disk
 - Logical formatting or "making a file system"
 - To increase efficiency most file systems group blocks into clusters
 - Disk I/O done in blocks
 - File I/O done in clusters



RAID Structure

- RAID redundant array of inexpensive disks
 - multiple disk drives provides reliability via redundancy
- Increases the mean time to failure
- Mean time to repair exposure time when another failure could cause data loss
- Mean time to data loss based on above factors
- If mirrored disks fail independently, consider disk with 100,000 mean time to failure and 10 hour mean time to repair
 - Mean time to data loss is $100,000^2 / (2 * 10) = 500 * 10^6$ hours, or 57,000 years!
- Frequently combined with NVRAM to improve write performance
- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively





RAID (Cont.)

- Disk striping uses a group of disks as one storage unit
- RAID is arranged into six different levels
- RAID schemes improve performance and improve the reliability of the storage system by storing redundant data
 - Mirroring or shadowing (RAID 1) keeps duplicate of each disk
 - Striped mirrors (RAID 1+0) or mirrored stripes (RAID 0+1) provides high performance and high reliability
 - Block interleaved parity (RAID 4, 5, 6) uses much less redundancy
- RAID within a storage array can still fail if the array fails, so automatic replication of the data between arrays is common
- Frequently, a small number of hot-spare disks are left unallocated, automatically replacing a failed disk and having data rebuilt onto them



RAID Levels





(b) RAID 1: mirrored disks.



(c) RAID 4: block-interleaved parity.



(d) RAID 5: block-interleaved distributed parity.



(e) RAID 6: P + Q redundancy.





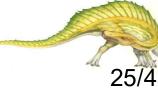








(f) Multidimensional RAID 6.

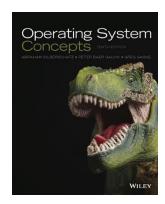




Other Features

- Regardless of where RAID implemented, other useful features can be added
- Snapshot is a view of file system before a set of changes take place (i.e. at a point in time)
 - More in Ch 12
- Replication is automatic duplication of writes between separate sites
 - For redundancy and disaster recovery
 - Can be synchronous or asynchronous
- Hot spare disk is unused, automatically used by RAID production if a disk fails to replace the failed disk and rebuild the RAID set if possible
 - Decreases mean time to repair

Bibliografia



Capítulo 11.



Capítulo –.