

Chapter 11: Mass-Storage Systems





Chapter 11: Mass-Storage Systems

- Overview of Mass Storage Structure
- HDD Scheduling
- Storage Device Management
- Swap-Space Management
- Storage Attachment
- RAID Structure





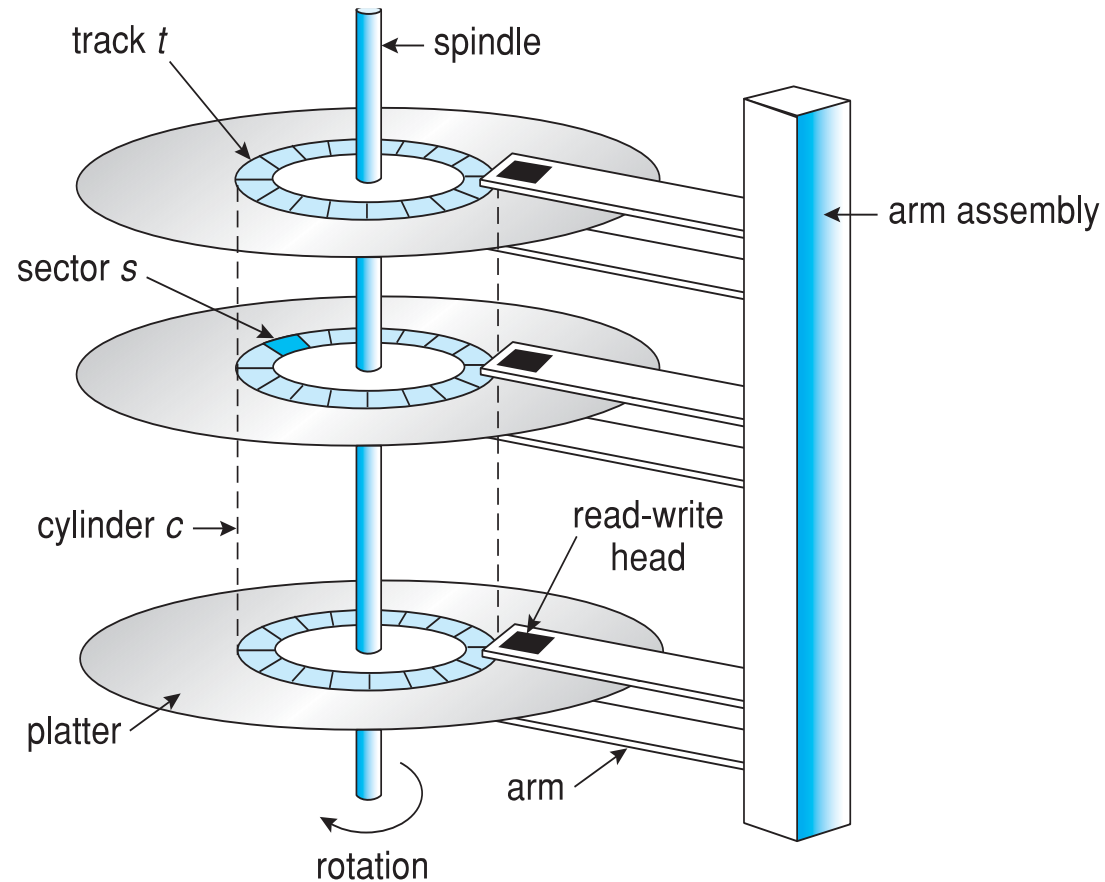
Objectives

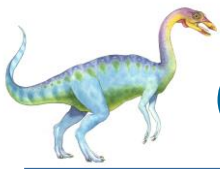
- Describe the **physical structure of secondary storage devices** and the **effect of a device's structure on its uses**
- Explain the performance characteristics of mass-storage devices
- Evaluate I/O scheduling algorithms
- Discuss **operating-system services provided for mass storage**, including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





Moving-head Disk Mechanism





Overview of Mass Storage Structure

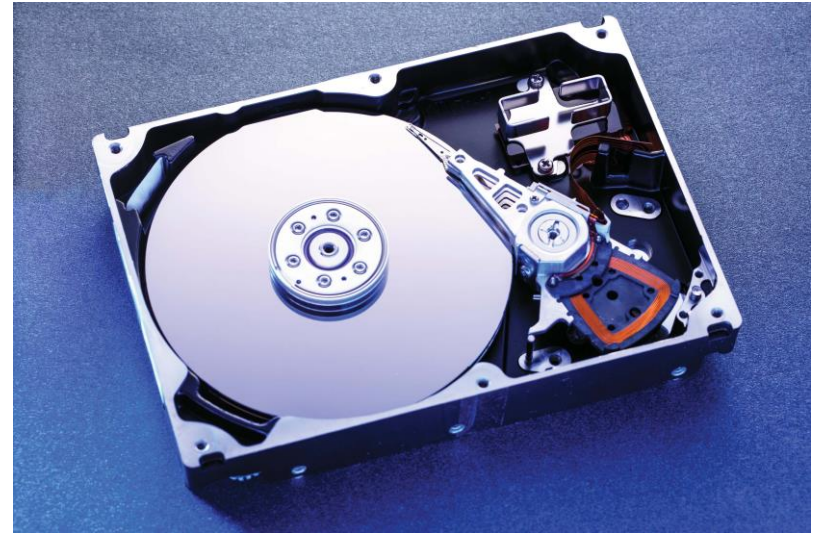
- Bulk of secondary storage for modern computers is **hard disk drives (HDDs)** and **nonvolatile memory (NVM)** devices
- **HDDs** spin **platters** of magnetically-coated material under moving **read-write heads**
 - Drives rotate at 60 to 250 times per second
 - **Transfer rate** is rate at which **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





Hard Disk Drives

- Common platter diameters range from **1.8 to 3.5 inches**
- **The sector size** 512 bytes until around 2010, post 4KB.
- The storage capacity: **GB and TB scale**
- A disk drive motor spins: 5,400, 7,200, 10,000, and 15,000 RPM.
- Typical disks can transfer tens to hundreds of **megabytes of data per 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
- **No moving parts**, so no seek time or rotational latency





Disk Attachment

- Host-attached storage accessed through I/O ports talking to **I/O busses**
- Several busses available, including **advanced technology attachment (ATA)**, **serial ATA (SATA)**, **eSATA**, **serial attached SCSI (SAS)**, **universal serial bus (USB)**, and **fibre channel (FC)**.
- Most common is SATA
- Because NVM much faster than HDD, new fast interface for NVM called **NVM express (NVMe)**, connecting directly to **Peripheral Component Interconnect (PCI)** bus





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**





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 \approx 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.)

- There are many sources of disk I/O request
 - OS
 - System processes
 - Users processes
- I/O request **includes** input or output mode, disk address, memory address, number of sectors to transfer
- OS maintains **queue of requests**, per disk or device
- Idle disk can immediately work on I/O request, busy disk means work must queue
 - Optimization algorithms only make sense when a queue exists
 - Both bandwidth and access time can be improved by processing requests in a good order.





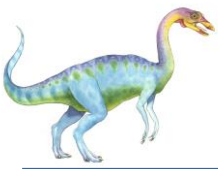
Disk Scheduling (Cont.)

- 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



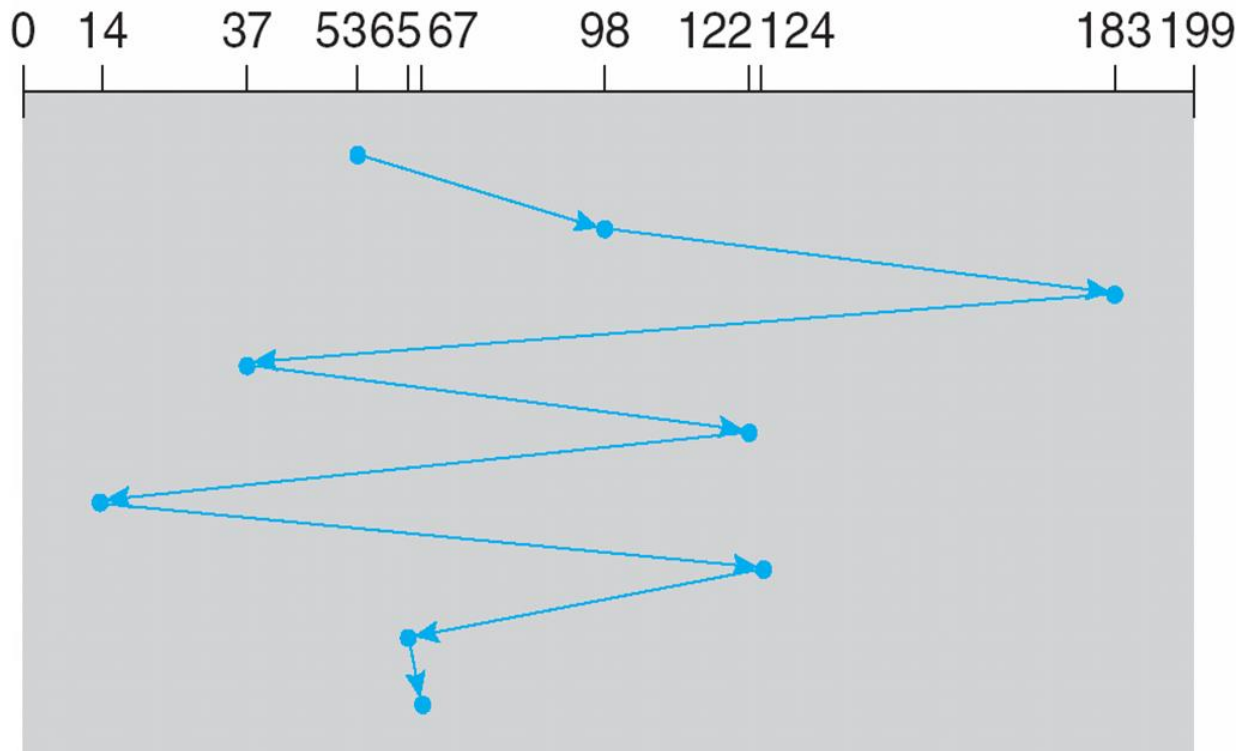


FCFS

Illustration shows total head movement of 640 cylinders

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

head starts at 53





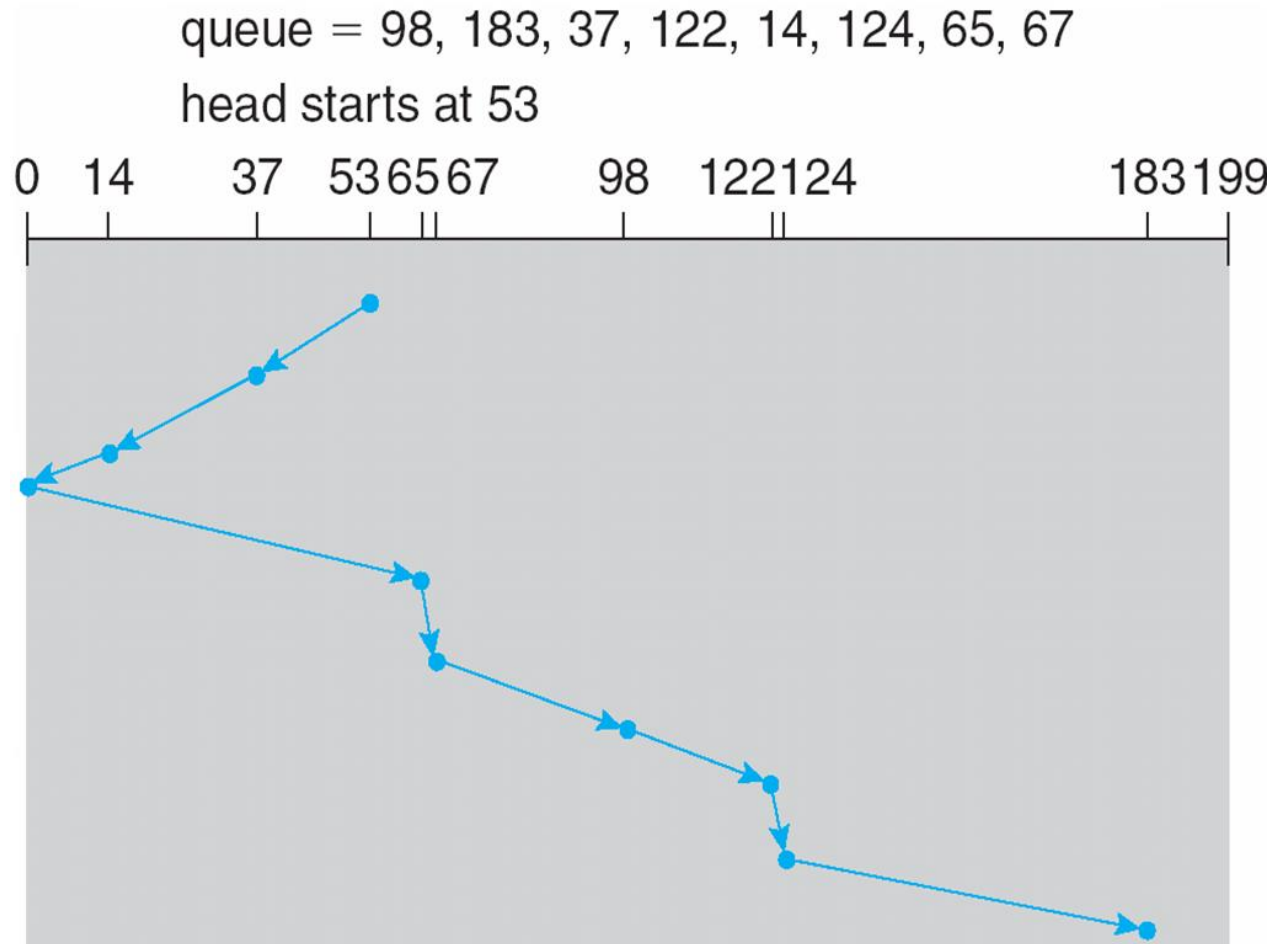
SCAN

- 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?

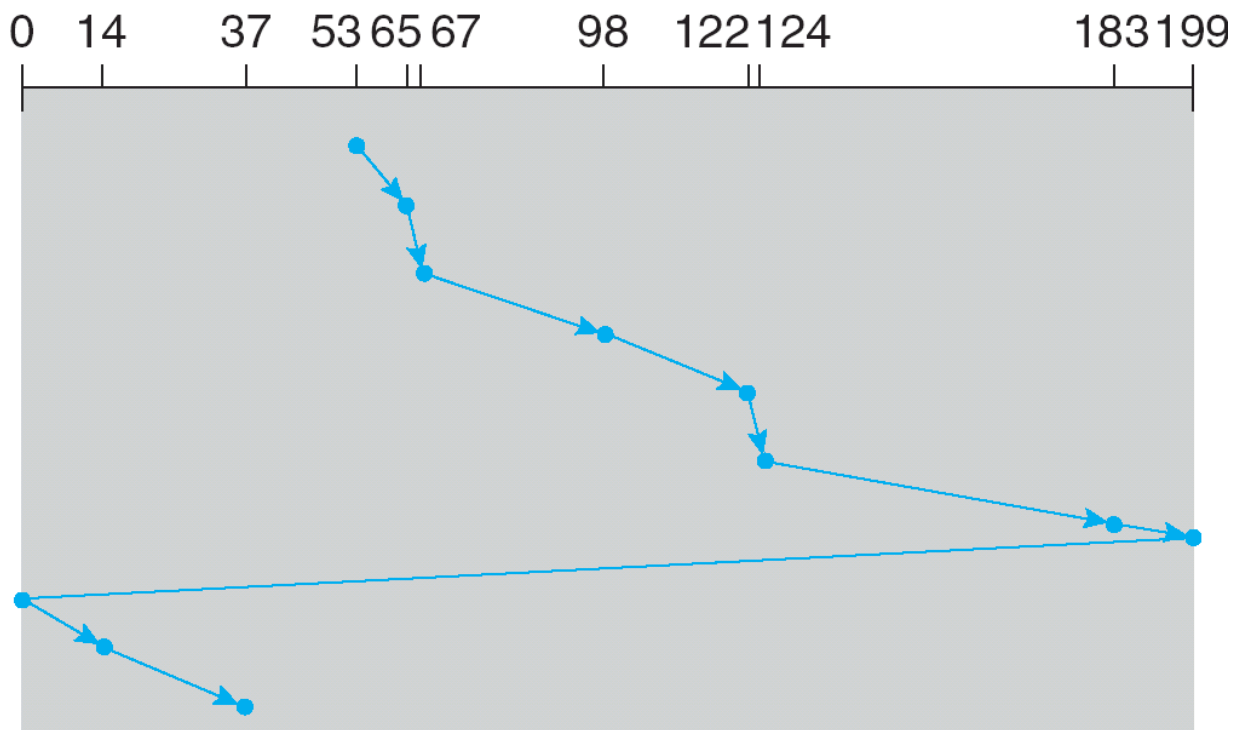




C-SCAN (Cont.)

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

head starts at 53





Selecting a Disk-Scheduling Algorithm

- SSTF is common and has a natural appeal
- SCAN and C-SCAN perform better for systems that place a **heavy load on the disk**
 - Less starvation, but still possible
- To avoid starvation Linux implements **deadline** scheduler
 - Maintains separate **read and write queues**, **gives read priority**
 - ▶ Because processes more likely to block on read than write
 - Implements four queues: 2 x read and 2 x write
 - ▶ 1 read and 1 write queue sorted in LBA order, essentially implementing C-SCAN
 - ▶ 1 read and 1 write queue sorted in FCFS order
 - ▶ All I/O requests sent in batch sorted in that queue's order
 - ▶ After each batch, checks if any requests in **FCFS older than configured age (default 500ms)**
 - If so, LBA queue containing that request is selected for next batch of I/O
- The **deadline I/O scheduler** is the default in the Linux RedHat 7 distribution





Error Detection and Correction

- Fundamental aspect of many parts of computing (memory, networking, storage)
- **Error detection** determines if there a problem has occurred (for example a bit flipping)
 - If detected, can halt the operation
 - Detection frequently done via parity bit
- Parity one form of **checksum** – uses modular arithmetic to compute, store, compare values of fixed-length words
 - Another error-detection method common in networking is **cyclic redundancy check (CRC)** which uses hash function to detect multiple-bit errors
- **Error-correction code (ECC)** not only detects, but can correct some errors
 - Soft errors correctable, hard errors detected but not corrected

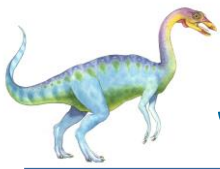




Storage Device Management

- **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





Storage Device Management (cont.)

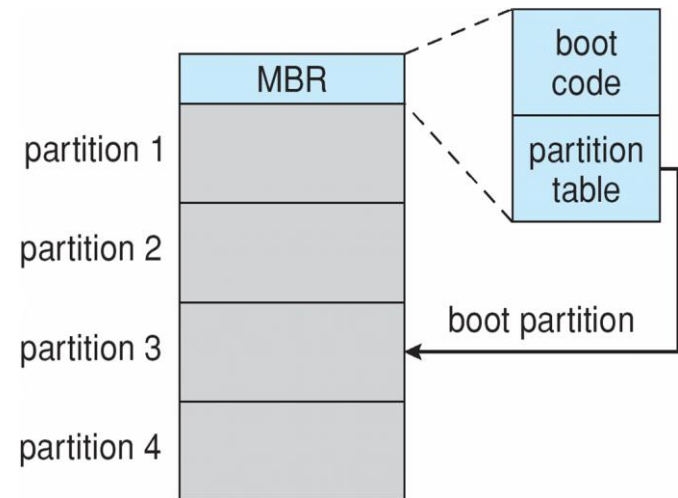
- **Root partition** contains the OS, other partitions can hold other Oses, other file systems, or be raw
 - **Mounted** at boot time
 - Other partitions can mount automatically or manually
- At mount time, file system consistency checked
 - Is all metadata correct?
 - ▶ If not, fix it, **try again**
 - ▶ If yes, add to mount table, allow access
- **Boot block can point to boot volume or boot loader set of blocks that contain enough code** to know how to load the kernel from the file system
 - Or a boot management program for multi-os booting





Device Storage Management (Cont.)

- Raw disk access for apps that want to do their own block management, keep OS out of the way (databases for example)
- Boot block initializes system
 - The bootstrap is stored in ROM, firmware
 - **Bootstrap loader** program stored in boot blocks of boot partition
- Methods such as **sector sparing** used to handle bad blocks



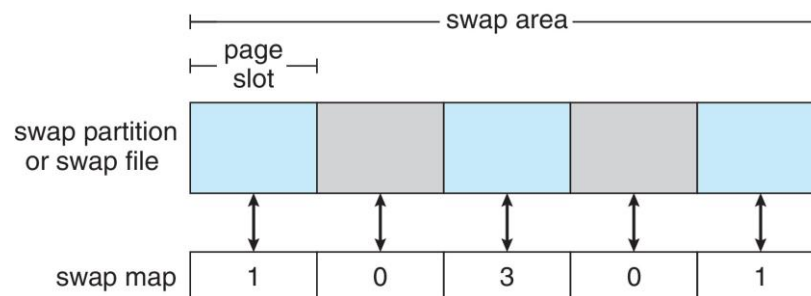
Booting from secondary storage in Windows





Swap-Space Management

- Used for moving entire processes (swapping), or pages (paging), from **DRAM to secondary storage** when DRAM not large enough for all processes
- Operating system provides **swap space management**
 - Secondary storage slower than DRAM, so important to optimize performance
 - Usually **multiple swap spaces possible** – decreasing I/O load on any given device
 - Best to have dedicated devices
 - Can be in **raw partition or a file within a file system** (for convenience of adding)
 - Data structures for swapping on Linux systems:





Storage Attachment

- Computers access storage in three ways
 - host-attached
 - network-attached
 - cloud
- Host attached access through local I/O ports, using one of several technologies
 - To attach many devices, use storage busses such as USB, firewire, thunderbolt
 - High-end systems use **fibre channel (FC)**
 - ▶ High-speed serial architecture using fibre or copper cables
 - ▶ Multiple hosts and storage devices can connect to the FC fabric





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 1300,000 **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



(a) RAID 0: non-redundant striping.



(b) RAID 1: mirrored disks.



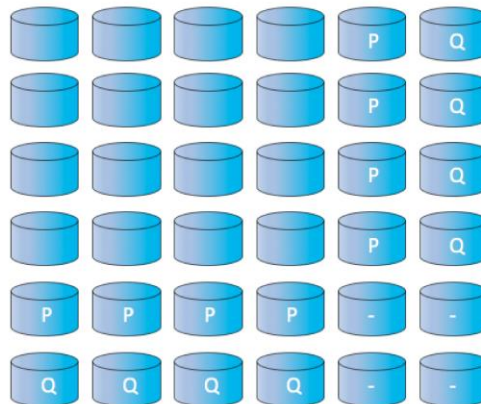
(c) RAID 4: block-interleaved parity.



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



(e) RAID 6: P + Q redundancy.

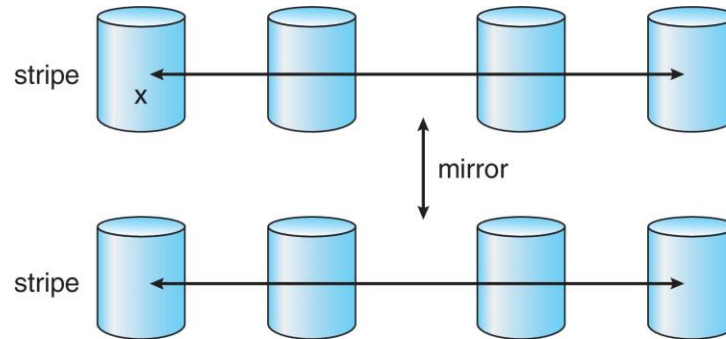


(f) Multidimensional RAID 6.

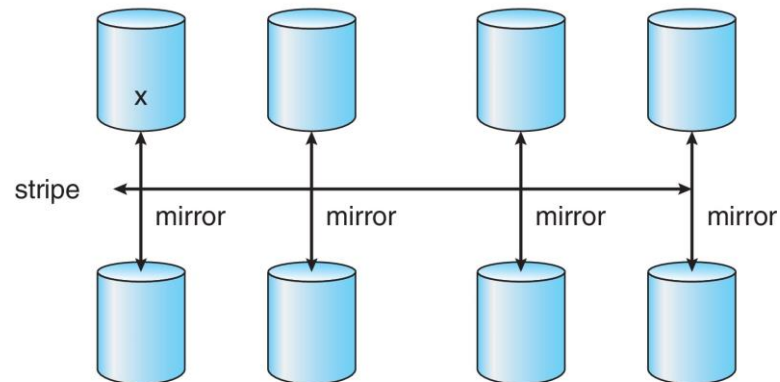




RAID (0 + 1) and (1 + 0)



a) RAID 0 + 1 with a single disk failure.



b) RAID 1 + 0 with a single disk failure.



End of Chapter 11

