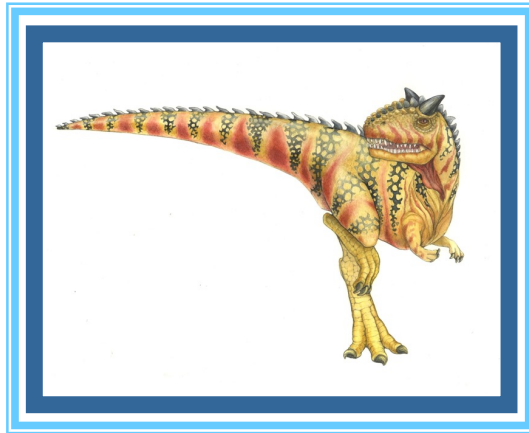


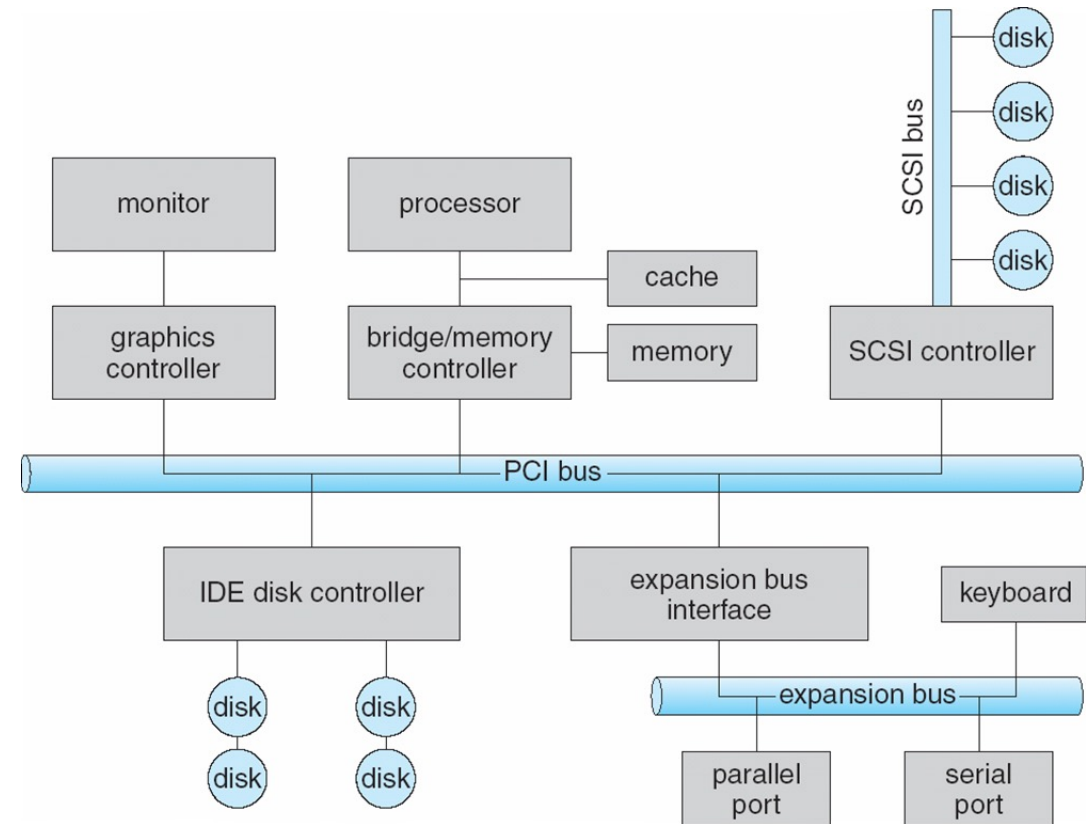
Chapter 11: Mass-Storage Systems





Chapter 11: Mass-Storage Systems

- Overview of Mass Storage Structure
- Disk Structure
- Disk Attachment
- Disk Scheduling
- Disk Management
- Swap-Space Management
- RAID Structure
- Stable-Storage





Objectives

- To describe the physical structure of secondary storage devices and its effects on the uses of the devices
- To explain the performance characteristics of mass-storage devices
- To evaluate disk scheduling algorithms
- To discuss operating-system services provided for mass storage, including RAID





File Systems

- First we'll discuss properties of physical disks
 - Structure
 - Performance
 - Scheduling

- Then we'll discuss how we build file systems on them
 - Files
 - Directories
 - Sharing
 - Protection
 - File System Layouts
 - File Buffer Cache
 - Read Ahead





Disks and the OS

- Disks are messy physical devices:
 - Errors, bad blocks, missed seeks, etc.
- The job of the OS is to hide this mess from higher level software
 - Low-level device control (initiate a disk read, etc.)
 - Higher-level abstractions (files, databases, etc.)
- The OS may provide different levels of disk access to different clients
 - Physical disk (surface, cylinder, sector)
 - Logical disk (disk block #)
 - Logical file (file block, record, or byte #)
- Logically, disk broken down into sectors
 - Addressed by cylinder, head, sector (CHS)





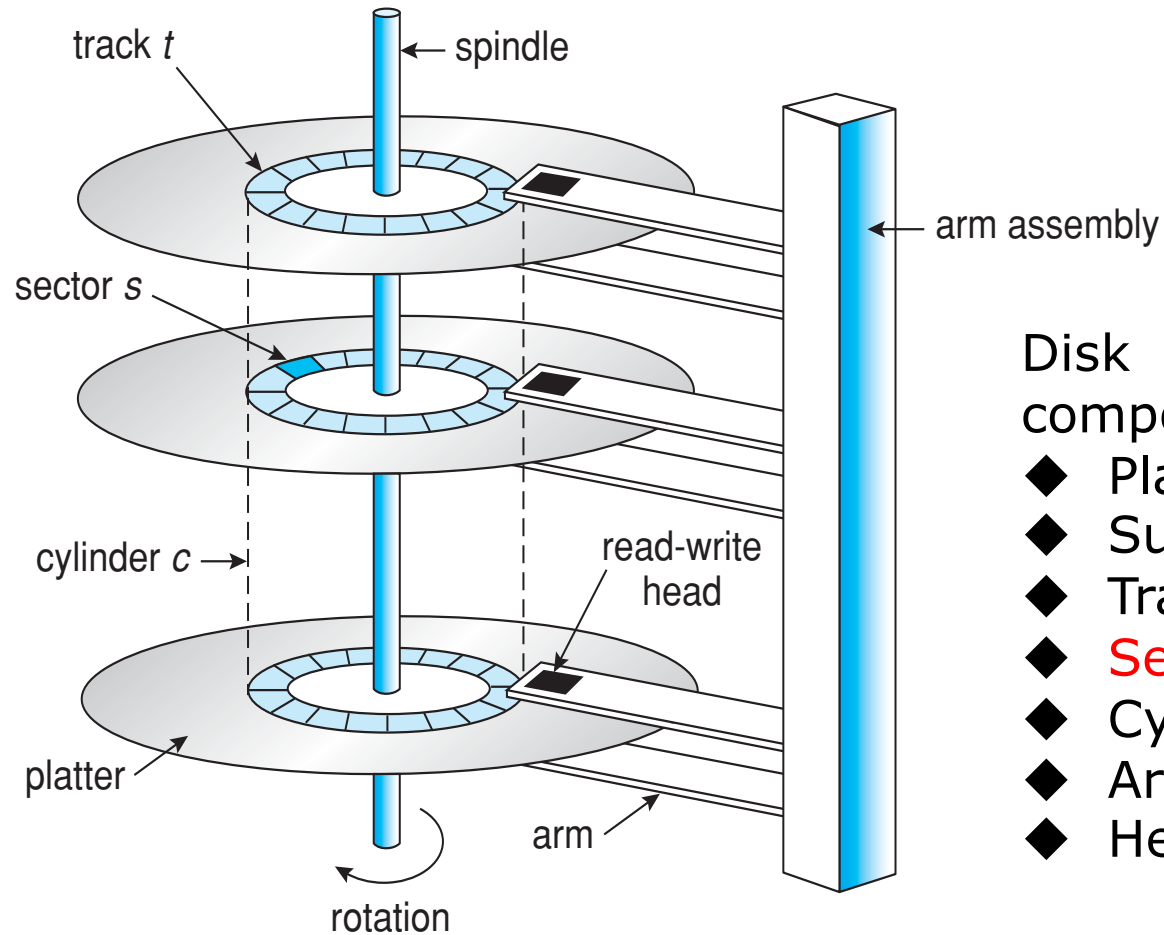
Overview of Mass Storage Structure

- **Magnetic disks** provide bulk of secondary storage of modern computers
 - 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
- Drive attached to computer via **I/O bus**
 - Busses vary, including **EIDE**, **ATA**, **SATA**, **USB**, **Fibre Channel**, **SCSI**, **SAS**, **Firewire**
 - **Host controller** in computer uses bus to talk to **disk controller** built into drive or storage array





Moving-head Disk Mechanism



Disk components

- ◆ Platters
- ◆ Surfaces
- ◆ Tracks
- ◆ Sectors
- ◆ Cylinders
- ◆ Arm
- ◆ Heads





Disk Interaction

- Specifying disk requests requires a lot of info:
 - Cylinder #, surface #, track #, sector #, transfer size...
- Older disks required the OS to specify all of this
 - The OS needed to know all disk parameters
- Modern disks are more complicated
 - Not all sectors are the same size, sectors are remapped, etc.
- Current disks provide a higher-level interface (SCSI)
 - The disk exports its data as a logical array of blocks [0...N]
 - ▶ Disk maps logical blocks to cylinder/surface/track/sector
 - Only need to specify the logical block # to read/write
 - But now the disk parameters are hidden from the OS





Hard Disks

- Platters range from .85" to 14" (historically)
 - Commonly 3.5", 2.5", and 1.8"
- Range from 30GB to 3TB per drive
- Performance
 - Transfer Rate – theoretical – 6 Gb/sec
 - Effective Transfer Rate – real – 1Gb/sec
 - Seek time from 3ms to 12ms – 9ms common for desktop drives
 - Average seek time measured or calculated based on 1/3 of tracks
 - Latency based on spindle speed
 - ▶ $1 / (\text{RPM} / 60) = 60 / \text{RPM}$
 - Average latency = $\frac{1}{2}$ latency

Spindle [rpm]	Average latency [ms]
4200	7.14
5400	5.56
7200	4.17
10000	3
15000	2

(From Wikipedia)





Disk Parameters

- Seagate Barracuda ES.2 (largest SATA disc available)
 - Form factor: 3.5"
 - Capacity: 1 TB
 - Rotation rate: 7,200 RPM
 - Platters: 4
 - Surfaces: 8
 - Sector size: 512 bytes
 - Cache: 32 MB
 - Transfer rate: 53 MB/s (inner) – 104 MB/s (outer)
 - Average seek: 9.5 ms





Hard Disk Performance (I)

- Disk request performance depends upon a number of steps
 - Seek – moving the disk arm to the correct cylinder
 - ▶ Depends on how fast disk arm can move (increasing very slowly)
 - Rotation – waiting for the sector to rotate under the head
 - ▶ Depends on rotation rate of disk (increasing, but slowly)
 - Transfer – transferring data from surface into disk controller electronics, sending it back to the host
 - ▶ Depends on density (increasing quickly)
- When the OS uses the disk, it tries to minimize the cost of all of these steps
 - Particularly seeks and rotation





Hard Disk Performance (II)

- **Access Latency** = **Average access time** = average seek time + average latency
 - For fastest disk $3\text{ms} + 2\text{ms} = 5\text{ms}$
 - For slow disk $9\text{ms} + 5.56\text{ms} = 14.56\text{ms}$
- 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 (average seek time) + 4.17ms (average latency for rotation, 7200 RPM) + 0.1ms (overhead) + transfer time
 - Transfer time = $4\text{KBytes} / 1\text{Gbits/sec} = (32 * 1024 \text{ bits}) / (1024^3 \text{ bits/sec}) = 32 / 1024^2 = 0.031 \text{ ms}$
 - Average I/O time for 4KB block = $9.27\text{ms} + .031\text{ms}$ (transfer time) = 9.301ms

Spindle [rpm]	Average latency [ms]
4200	7.14
5400	5.56
7200	4.17
10000	3
15000	2





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





Magnetic Tape

- Was early secondary-storage medium
 - Evolved from open spools to cartridges
- Relatively permanent and holds large quantities of data
- Access time slow
- Random access ~1000 times slower than disk
- Mainly used for backup, storage of infrequently-used data, transfer medium between systems
- Kept in spool and wound or rewound past read-write head
- Once data under head, transfer rates comparable to disk
 - 140MB/sec and greater
- 200GB to 1.5TB typical storage
- Common technologies are LTO-{3,4,5} and T10000





Solid-State Disks

- Nonvolatile memory used like a hard drive
 - Many technology variations
- Can be more reliable than HDDs
- More expensive per MB
- Maybe have shorter life span
- 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





Disk Structure

- 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





Disk Attachment

- Host-attached storage accessed through I/O ports talking to I/O busses
- SCSI itself is a bus, up to 16 devices on one cable, **SCSI initiator** requests operation and **SCSI targets** perform tasks
 - Each target can have up to 8 **logical units** (disks attached to device controller)
- FC is high-speed serial architecture
 - Can be switched fabric with 24-bit address space – the basis of **storage area networks (SANs)** in which many hosts attach to many storage units
- I/O directed to bus ID, device ID, logical unit (LUN)





Storage Array

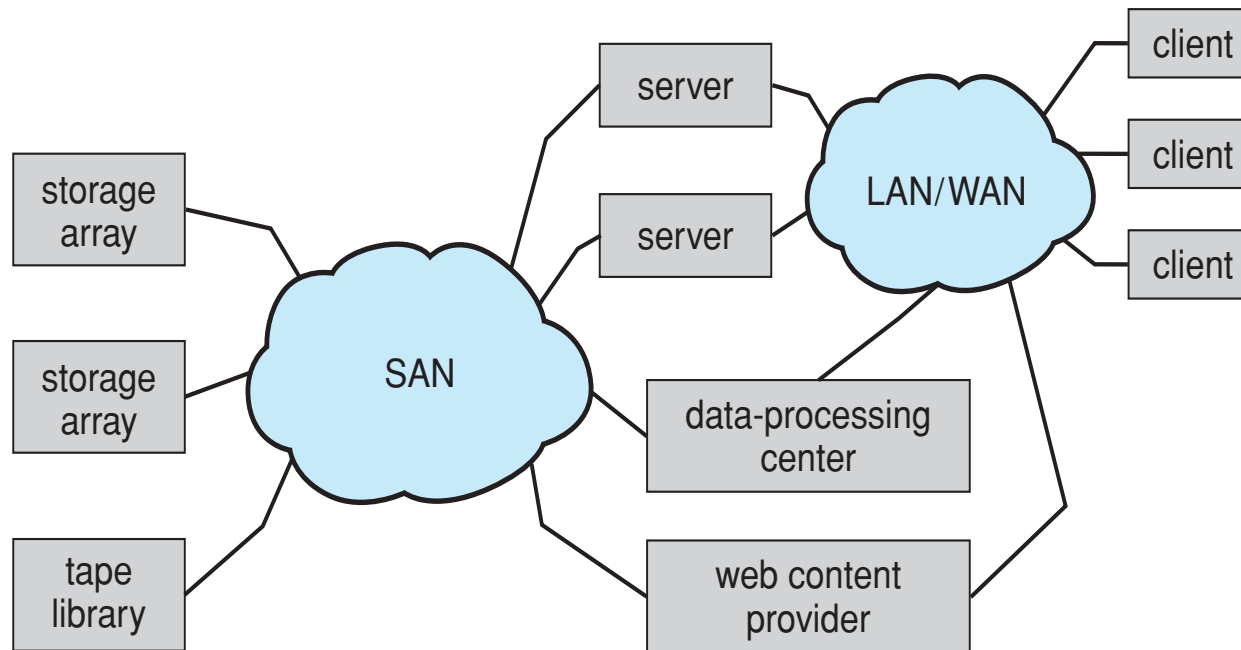
- Can just attach disks, or arrays of disks
- Storage Array has controller(s), provides features to attached host(s)
 - Ports to connect hosts to array
 - Memory, controlling software (sometimes NVRAM, etc)
 - A few to thousands of disks
 - RAID, hot spares, hot swap (discussed later)
 - Shared storage -> more efficiency
 - Features found in some file systems
 - ▶ Snapshots, clones, thin provisioning, replication, deduplication, etc





Storage Area Network

- Common in large storage environments
- Multiple hosts attached to multiple storage arrays - flexible





Storage Area Network (Cont.)

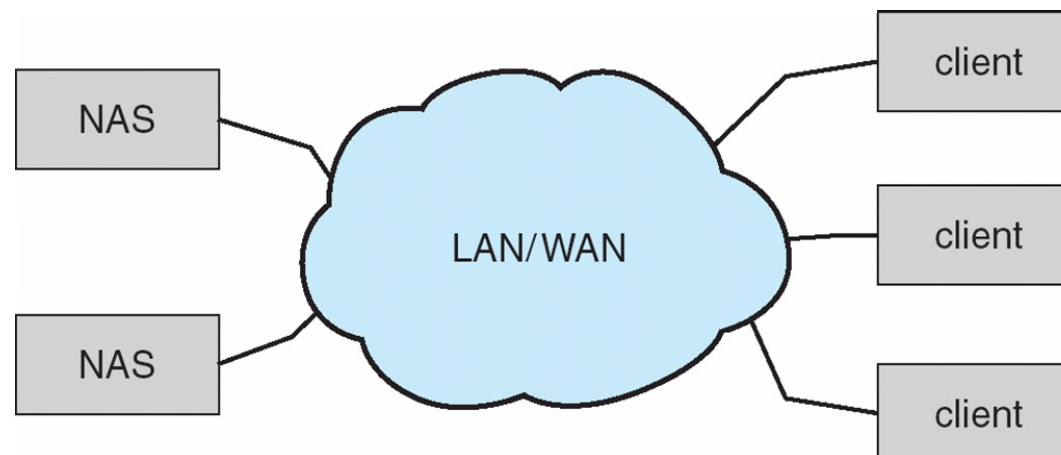
- SAN is one or more storage arrays
 - Connected to one or more Fibre Channel switches
- Hosts also attach to the switches
- Storage made available via **LUN Masking** from specific arrays to specific servers
- Easy to add or remove storage, add new host and allocate it storage
 - Over low-latency Fibre Channel fabric
- Why have separate storage networks and communications networks?
 - Consider iSCSI, FCOE





Network-Attached Storage

- Network-attached storage (**NAS**) is storage made available over a network rather than over a local connection (such as a bus)
 - Remotely attaching to file systems
- NFS and CIFS are common protocols
- Implemented via remote procedure calls (RPCs) between host and storage over typically TCP or UDP on IP network
- **iSCSI** protocol uses IP network to carry the SCSI protocol
 - Remotely attaching to devices (blocks)





Disk 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





Disk Scheduling (Cont.)

- Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying “depth”)
- Because seeks are so expensive (milliseconds!), it helps to schedule disk requests that are queued waiting for the disk
- Several algorithms exist to schedule the servicing of disk I/O requests
 - FCFS (do nothing)
 - ▶ Reasonable when load is low
 - ▶ Long waiting times for long request queues
 - SSTF (shortest seek time first)
 - ▶ Minimize arm movement (seek time), maximize request rate
 - ▶ Favors middle blocks
 - SCAN (elevator)
 - ▶ Service requests in one direction until done, then reverse
 - C-SCAN
 - ▶ Like SCAN, but only go in one direction (typewriter)
 - C-LOOK (a variant of C-SCAN)

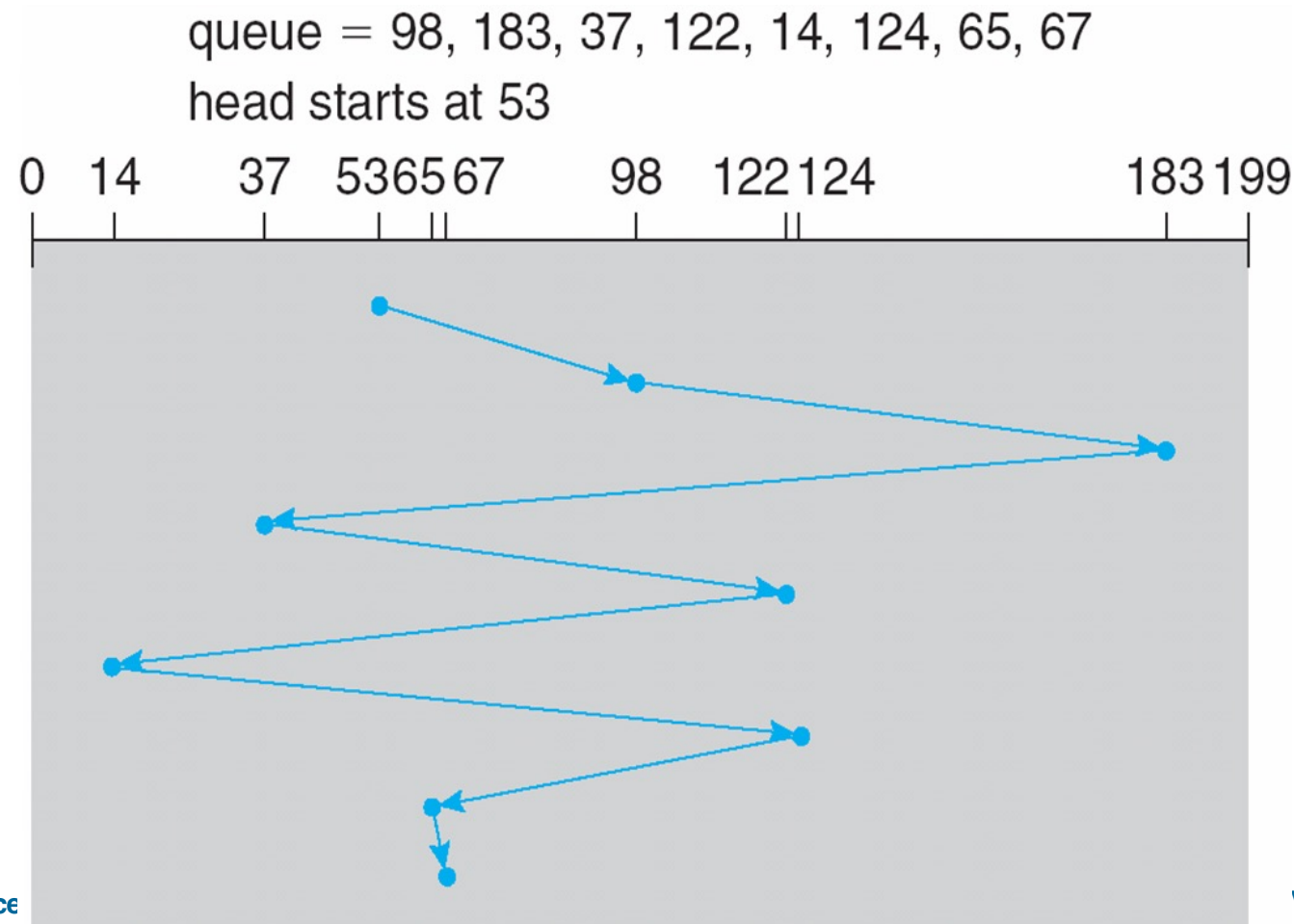




FCFS

We illustrate scheduling algorithms with a request queue (0-199)
98, 183, 37, 122, 14, 124, 65, 67
Head pointer 53

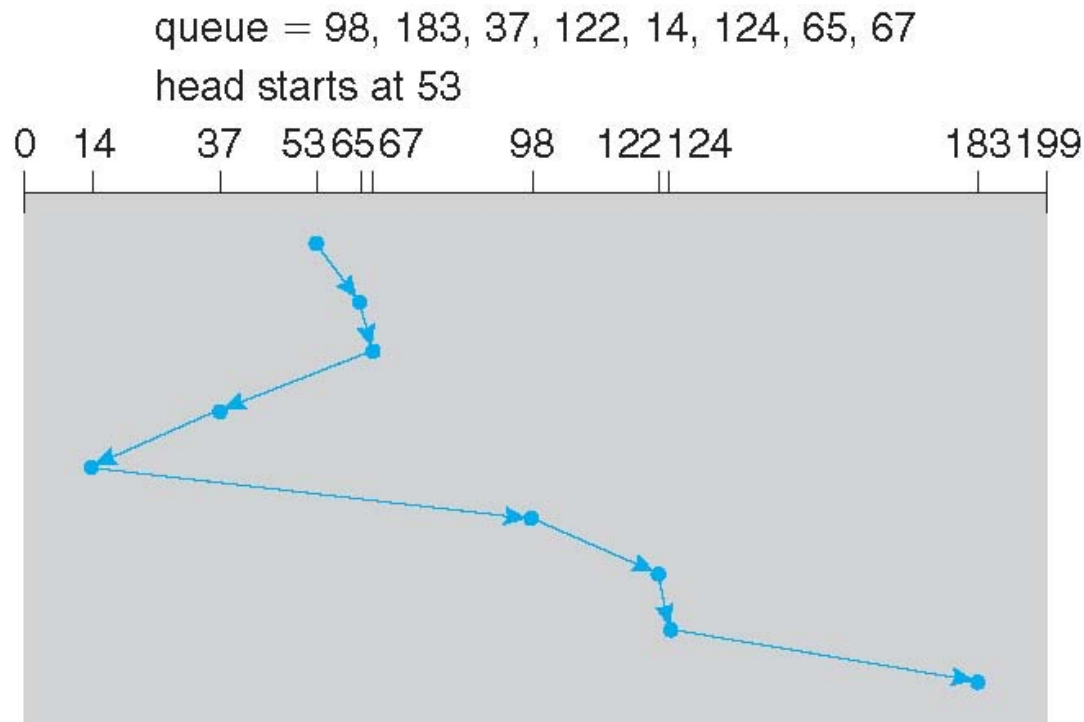
Illustration shows total head movement of 640 cylinders





SSTF

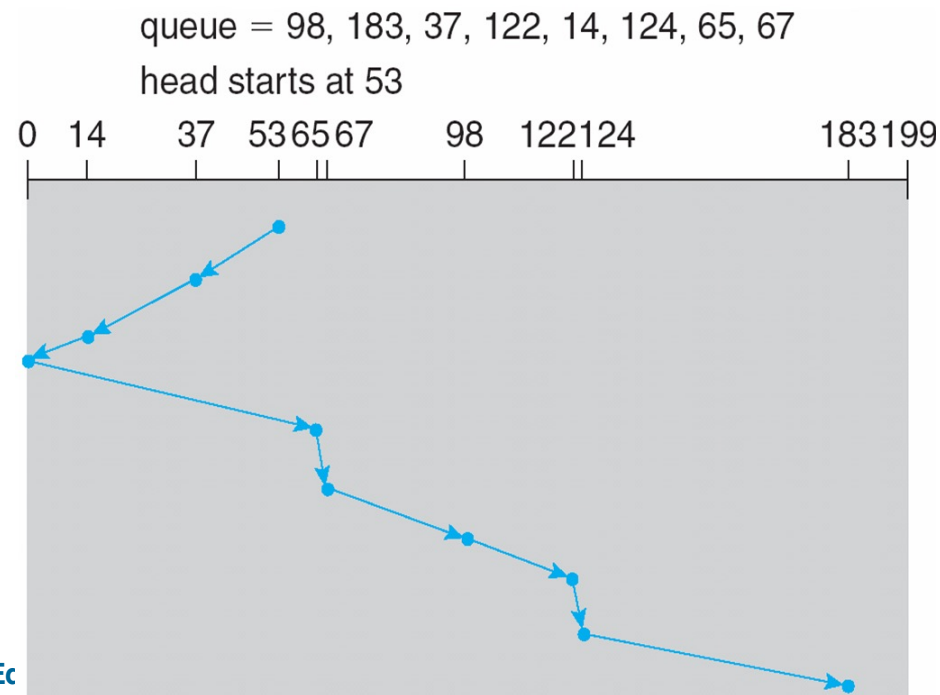
- Shortest Seek Time First selects the request with the minimum seek time from the current head position
- SSTF scheduling is a form of SJF scheduling; may cause starvation of some requests
- Illustration shows total head movement of 236 cylinders





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 236 cylinders
- But note that if requests are uniformly dense, largest density at other end of disk and those wait the longest



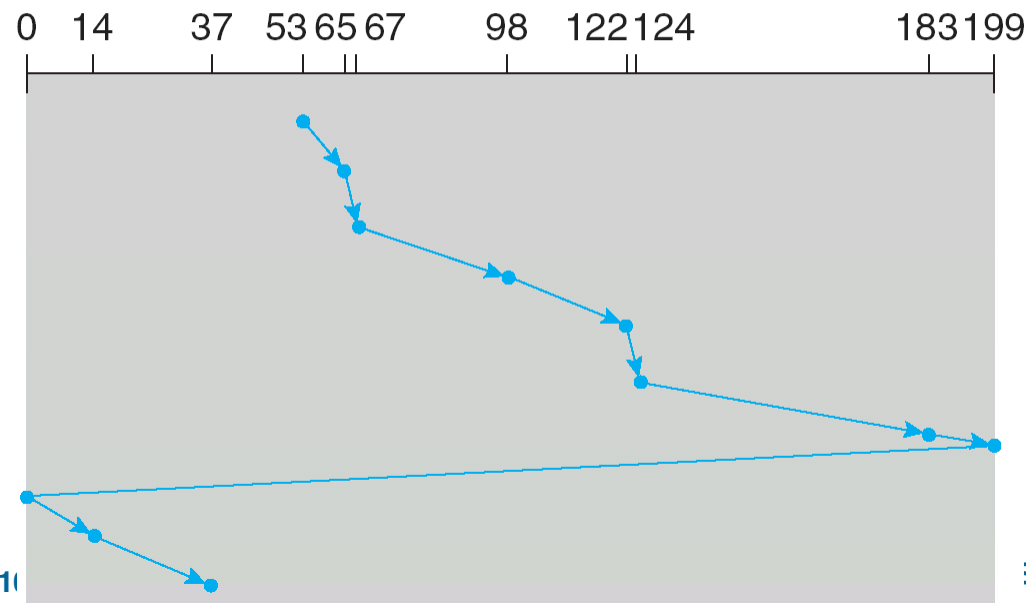


C-SCAN

- 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

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

head starts at 53



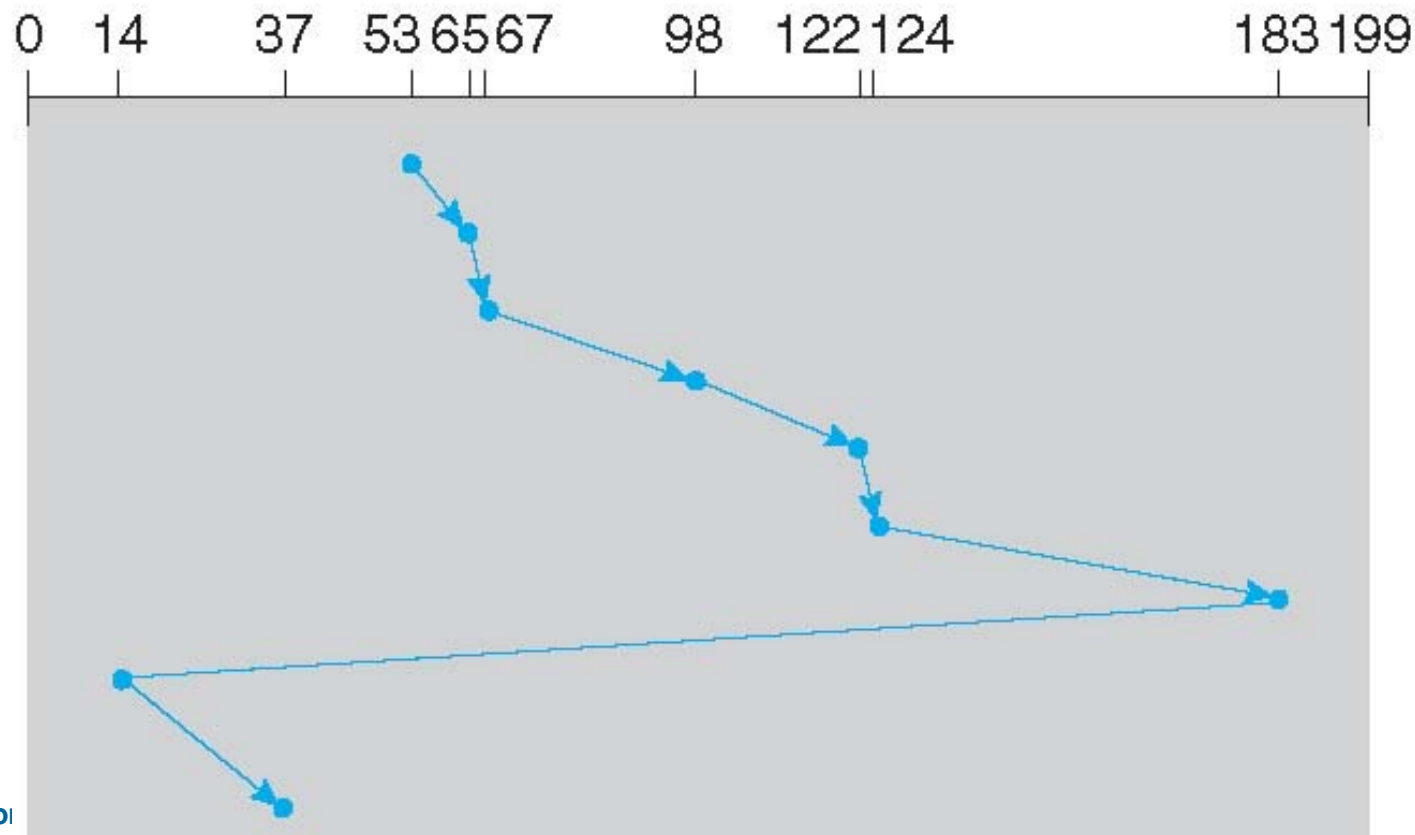


C-LOOK

- LOOK a version of SCAN, C-LOOK a version of C-SCAN
- Arm only goes as far as the last request in each direction, then reverses direction immediately, without first going all the way to the end of the disk.

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
- Performance depends on the number and types of requests
- Requests for disk service can be influenced by the file-allocation method
 - And metadata layout
- The disk-scheduling algorithm should be written as a separate module of the operating system, allowing it to be replaced with a different algorithm if necessary
- Either SSTF or LOOK is a reasonable choice for the default algorithm
- What about rotational latency?
 - Difficult for OS to calculate
- In general, modern disks often do the disk scheduling themselves
 - Disks know their layout better than OS, can optimize better
 - Ignores, undoes any scheduling done by OS





Disk 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





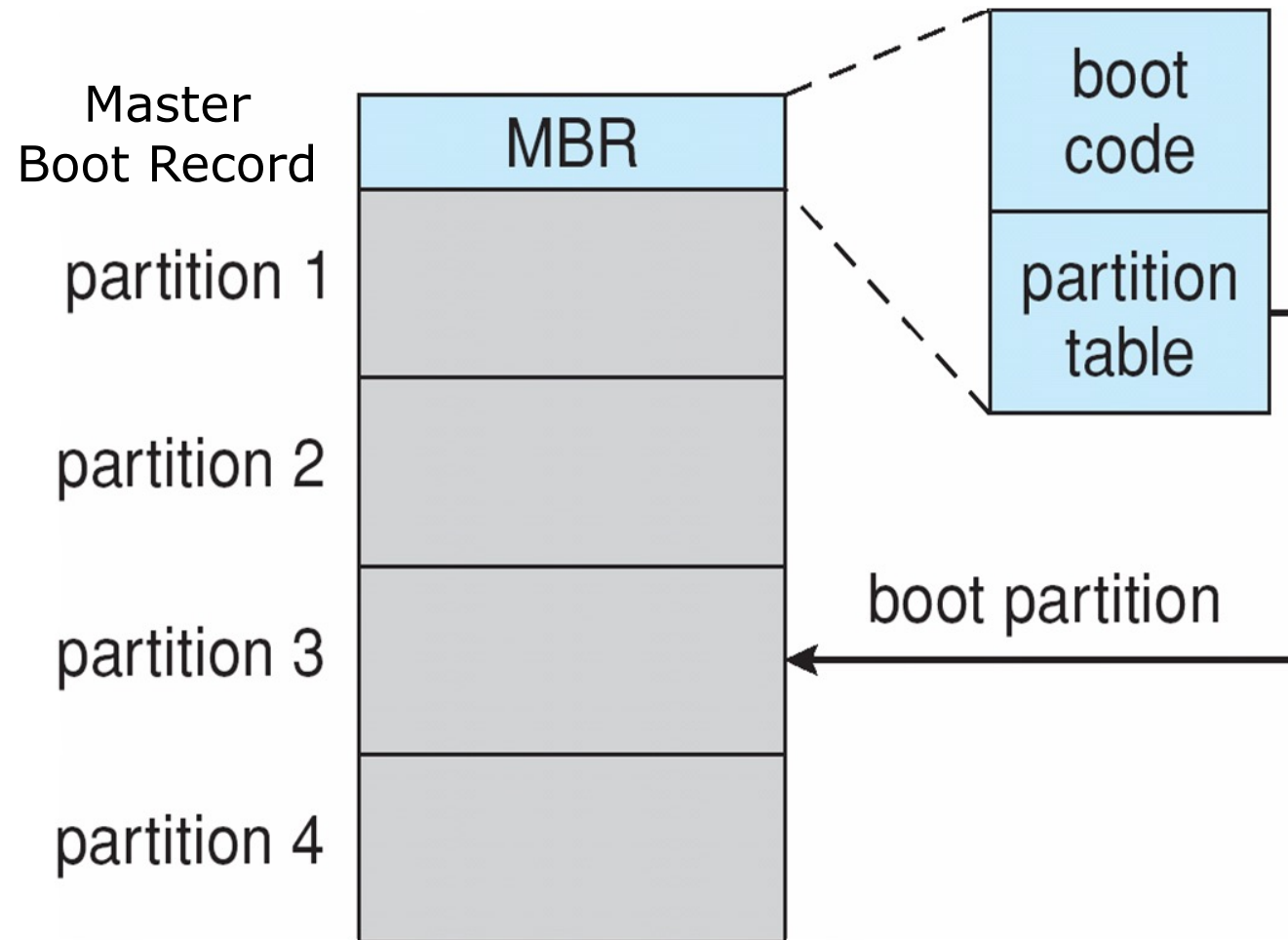
Disk Management (Cont.)

- 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
 - **Bootstrap loader** program stored in boot blocks of boot partition
- Methods such as **sector sparing** used to handle bad blocks





Booting from a Disk in Windows





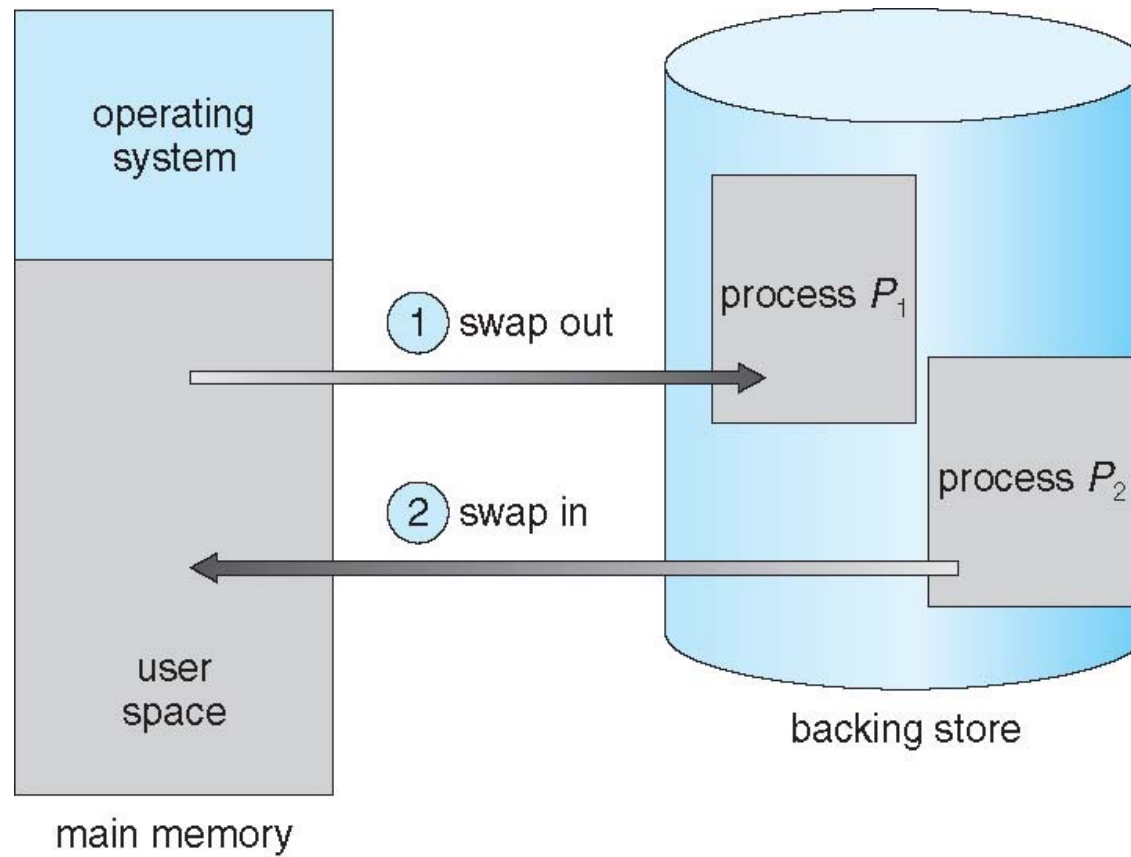
Swap-Space in Disk

- Swap-space — Virtual memory uses disk space as an extension of main memory
 - Less common now due to memory capacity increases
- A process can be **swapped** temporarily out of memory to a backing store, and then brought back into memory for continued execution
 - Total physical memory space of processes can exceed physical memory
- Major part of swap time is transfer time; total transfer time is directly proportional to the amount of memory swapped
- System maintains a **ready queue** of ready-to-run processes which have memory images on disk
- What if a system runs out of swap space? 1) Buy more RAM; 2) Reduce Memory Use; 3) Some systems allow multiple swap spaces



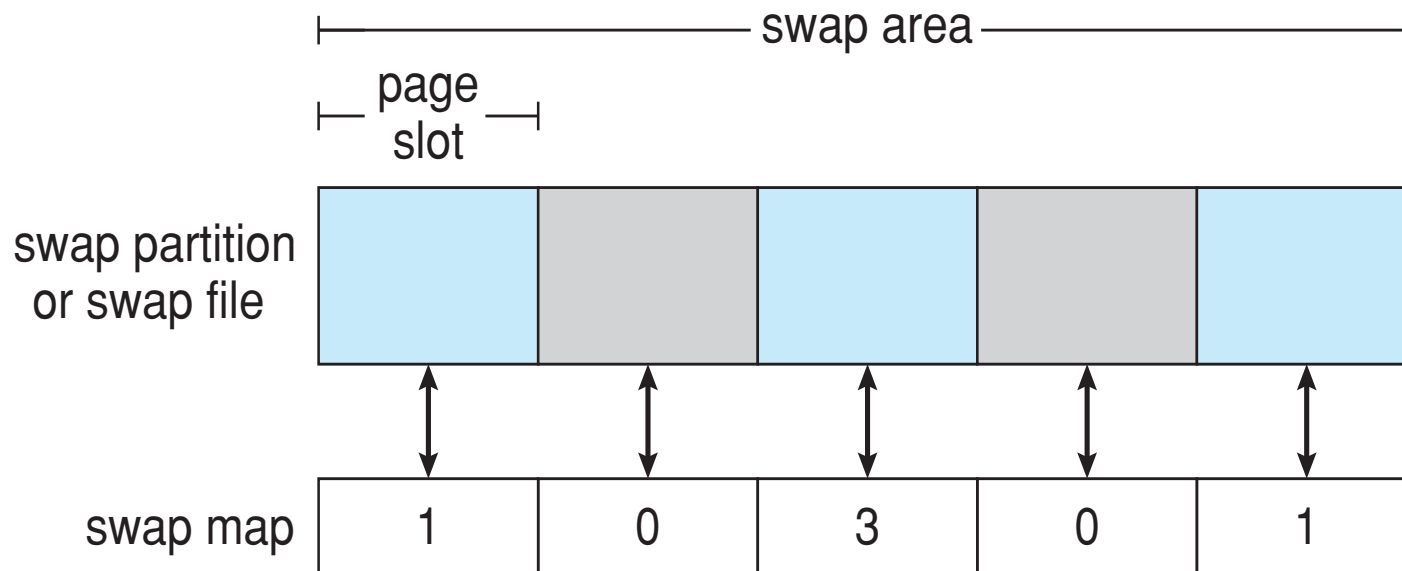


Schematic View of Swapping





Data Structures for Swapping on Linux Systems



Swap map—an array of integer counters. The value of the counter indicates the number of mappings to the swapped page. For example, a value of 3 indicates that the swapped page is mapped to three different processes (which can occur if the swapped page is storing a region of memory shared by three processes).





Context Switch Time including Swapping

- If next processes to be put on CPU is not in memory, need to swap out a process and swap in target process
- Context switch time can then be very high
- 100MB process swapping to hard disk with transfer rate of 50MB/sec
 - Swap out time of 2000 ms
 - Plus swap in of same sized process
 - Total context switch swapping component time of 4000ms (4 seconds)
- Can reduce if reduce size of memory swapped – by knowing how much memory really being used
 - System calls to inform OS of memory use via `request_memory()` and `release_memory()`





Swapping on Mobile Systems

- Not typically supported
 - Flash memory based
 - ▶ Small amount of space
 - ▶ Limited number of write cycles
 - ▶ Poor throughput between flash memory and CPU on mobile platform
- Instead use other methods to free memory if low
 - iOS **asks** apps to voluntarily relinquish allocated memory
 - ▶ Read-only data thrown out and reloaded from flash if needed
 - ▶ Failure to free can result in termination
 - Android terminates apps if low free memory, but first writes **application state** to flash for fast restart
 - Both OSes support paging as discussed below





RAID: Reliability via Redundancy

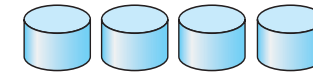
- RAID – redundant array of inexpensive disks
 - multiple disk drives provides reliability via **redundancy**
 - RAID is arranged into six different levels
- Idea: Use many disks in parallel to increase storage bandwidth, improve reliability
 - Files are striped across disks
 - Each stripe portion is read/written in parallel
 - Bandwidth increases with more disks
- 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!





RAID: Performance via Parallelism

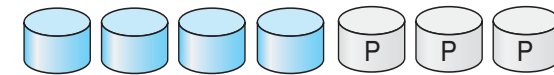
- RAID schemes improve performance and improve the reliability of the storage system by storing redundant data
- For example, use one disk to store parity data and XOR of all data blocks in stripe. Then it can recover any data block from all others + parity block.
- Hence “redundant” in name
- Introduces overhead, but, hey, disks are “inexpensive”
 - **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



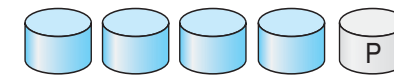
(a) RAID 0: non-redundant striping.



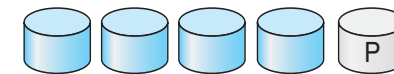
(b) RAID 1: mirrored disks.



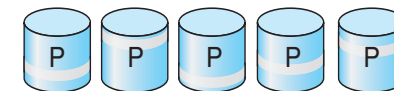
(c) RAID 2: memory-style error-correcting codes.



(d) RAID 3: bit-interleaved parity.



(e) RAID 4: block-interleaved parity.



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



(g) RAID 6: P + Q redundancy.

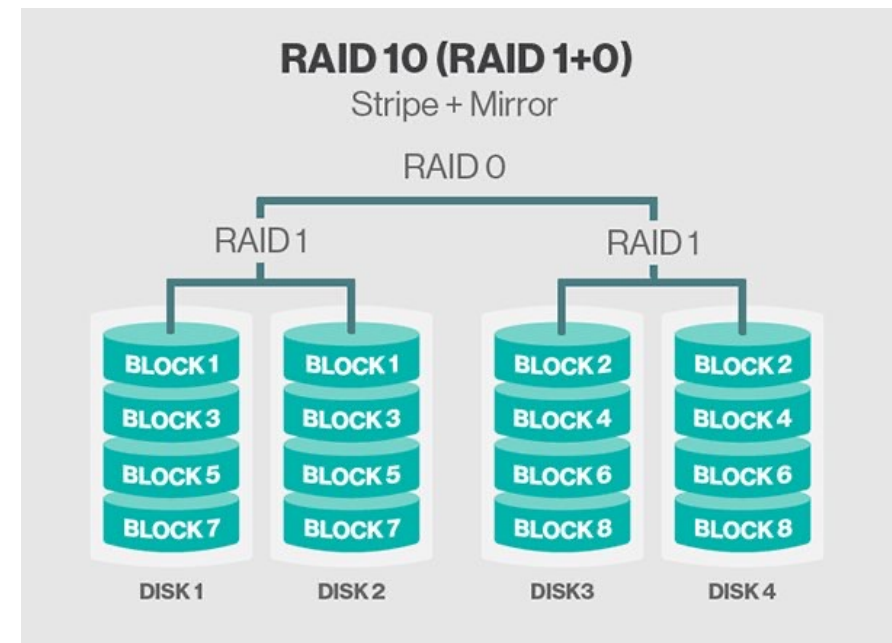
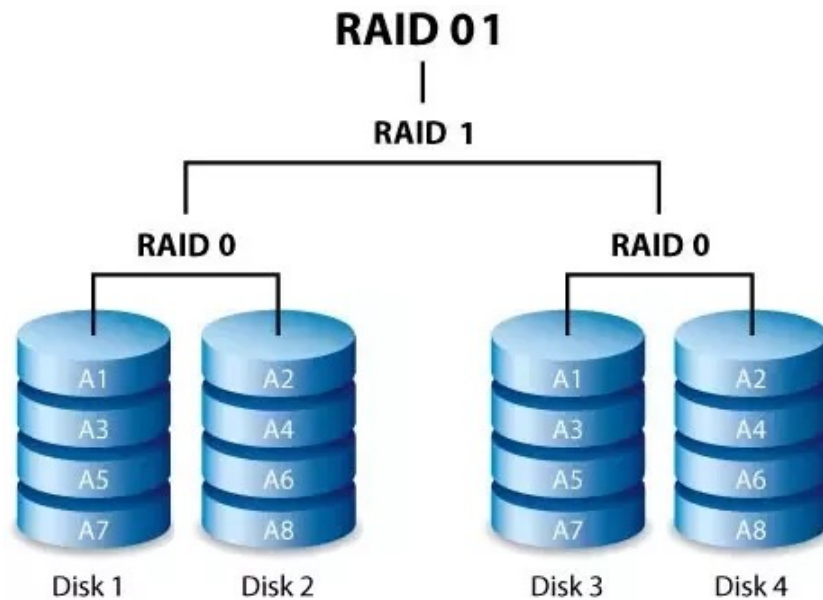
Bottle
necks
in P





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

Disk **striping** uses a group of disks as one storage unit





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
- Nevertheless, RAID alone does not prevent or detect data corruption or other errors, just disk failures.



End of Chapter 11

