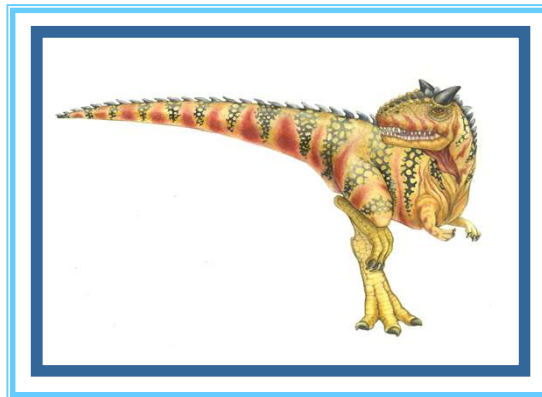
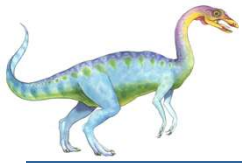


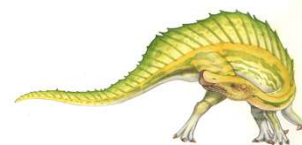
Chapter 12: Mass-Storage Systems

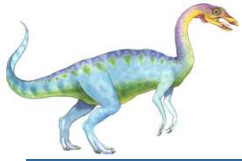




Chapter 12: Mass-Storage Systems

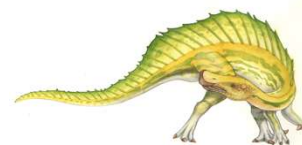
- Overview of Mass Storage Structure
- Disk Structure
- Disk Attachment
- Disk Scheduling
- Disk Management
- Swap-Space Management
- RAID Structure
- Stable-Storage Implementation
- Tertiary Storage Devices

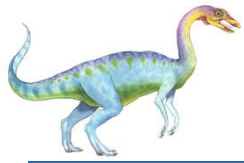




Overview of Mass Storage Structure

- **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



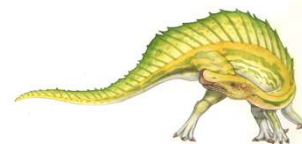


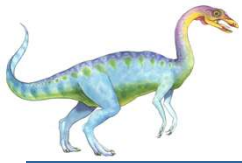
Magnetic Disks

- 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
 - ▶ $1/(RPM * 60)$
 - Average latency = $\frac{1}{2}$ latency

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

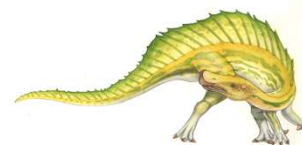
(From Wikipedia)





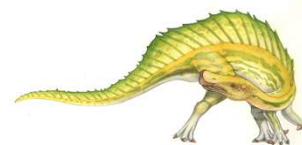
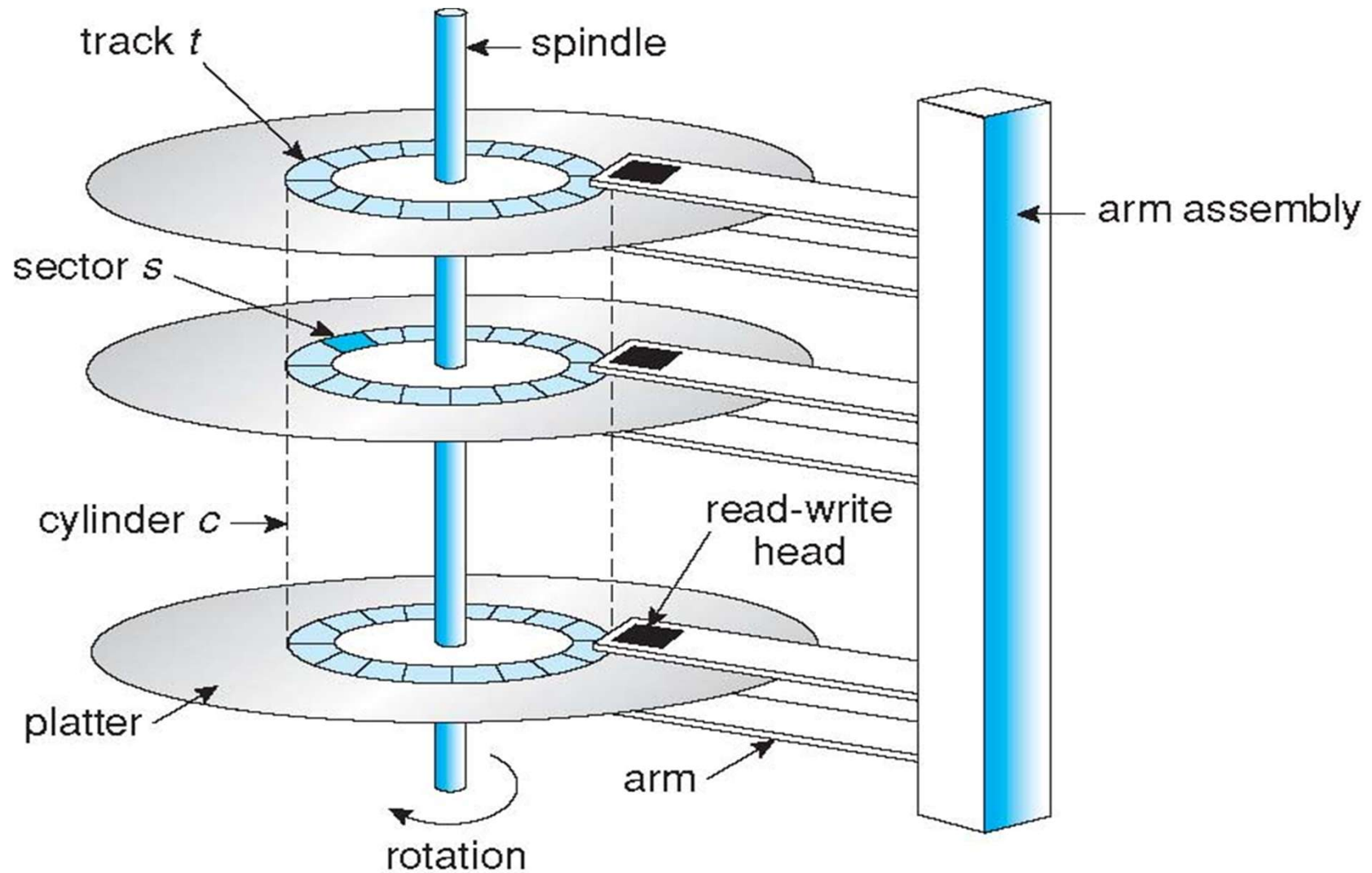
Magnetic Disk Performance

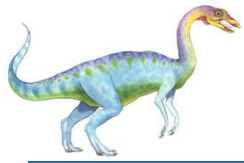
- **Access Latency = Average access time** = average seek time + average latency
 - For fastest disk $3\text{ms} + 2\text{ms} = 5\text{ms}$
 - For slow disk $9\text{ms} + 5.56\text{ms} = 14.56\text{ms}$
- Average I/O time = average access time + (amount to transfer / transfer rate) + controller overhead
- 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 =
 - $5\text{ms} + 4.17\text{ms} + 4\text{KB} / 1\text{Gb/sec} + 0.1\text{ms} =$
 - $9.27\text{ms} + 4 / 131072 \text{ sec} =$
 - $9.27\text{ms} + .12\text{ms} = 9.39\text{ms}$





Moving-head Disk Mechanism

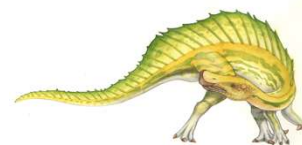


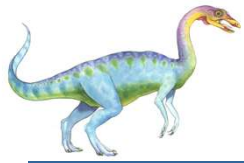


Disk Structure

- Disk drives are addressed as large 1-dimensional arrays of **logical blocks**, where the logical block is the smallest unit of transfer

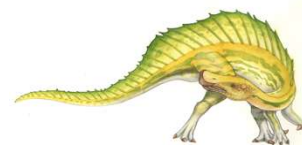
- The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially
 - Sector 0 is the first sector of the first track on the outermost cylinder
 - Mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the 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

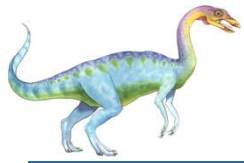




Disk Attachment

- 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)

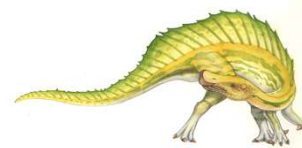


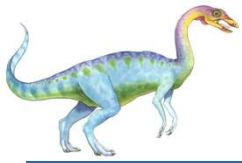


Storage Array

- Can just attach disks, or arrays of disks

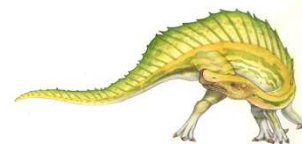
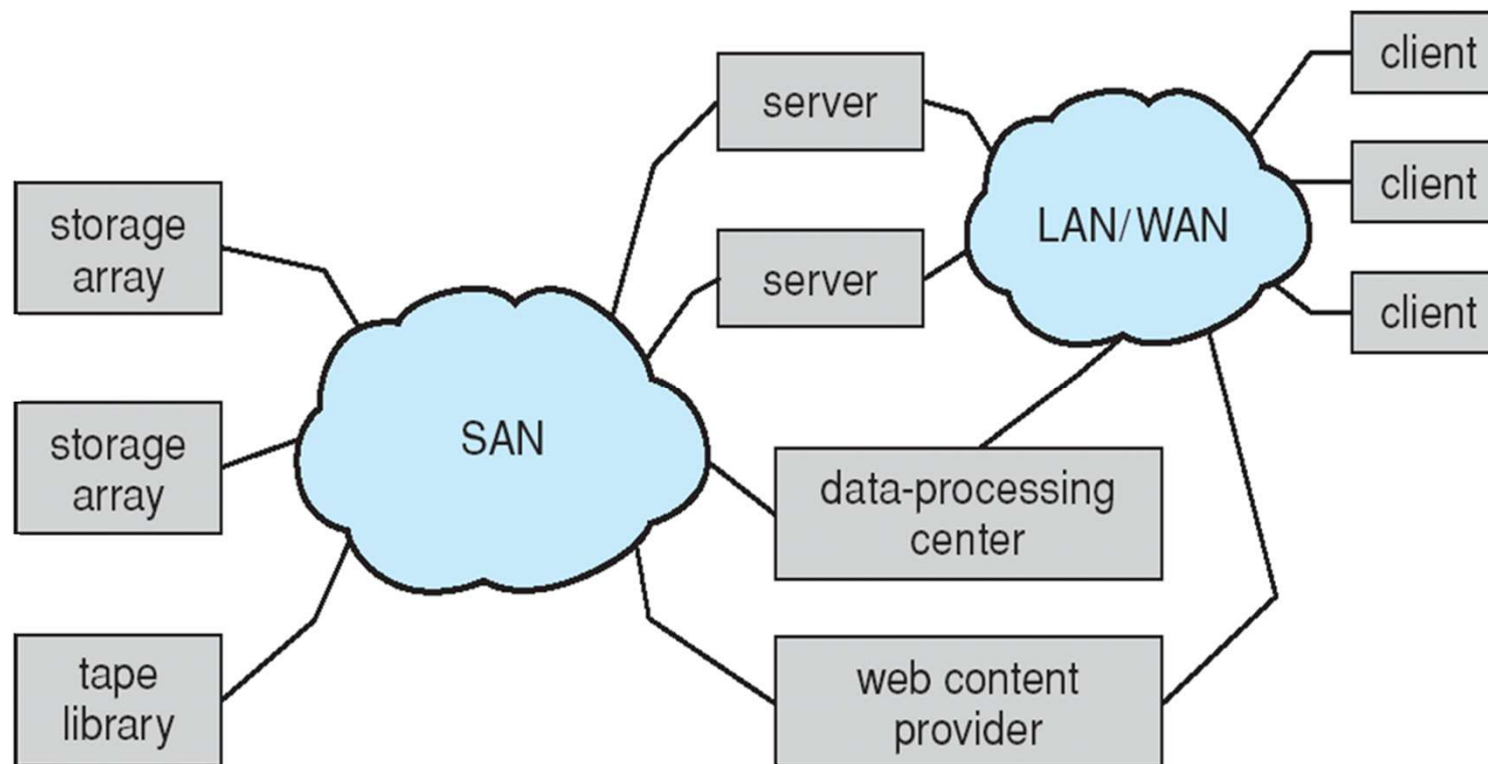
- Storage Array has controller(s), provides features to attached host(s)
 - Ports to connect hosts to array
 - Memory, controlling software (sometimes NVRAM, etc)
 - A few to thousands of disks
 - RAID, hot spares, hot swap (discussed later)
 - Shared storage -> more efficiency
 - Features found in some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

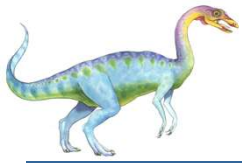




Storage Area Network

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





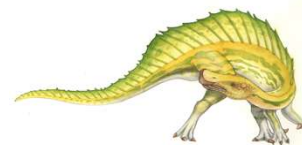
Storage Area Network (Cont.)

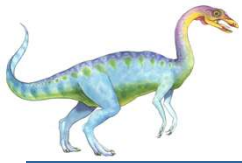
■ Vantagens

- Permite aumentar a capacidade de armazenamento, pois múltiplos servidores podem compartilhar a reserva de crescimento
- Acesso rápido a informação
- Permite aos servidores fazerem o boot diretamente pelo SAN, permitindo a rápida substituição de servidores defeituosos – aumentando assim a disponibilidade do sistema, tão almejada nos dias de hoje
- Permite o processo de recuperação de dados, pois o SAN pode replicar dados de vários servidores para uma área de armazenamento secundária.
- Aumenta a disponibilidade do sistema ao balancear a carga da rede, permitindo transferências de rápidas de grandes volumes de dados, reduz também a latência de I/O

■ Desvantagens

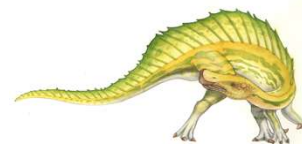
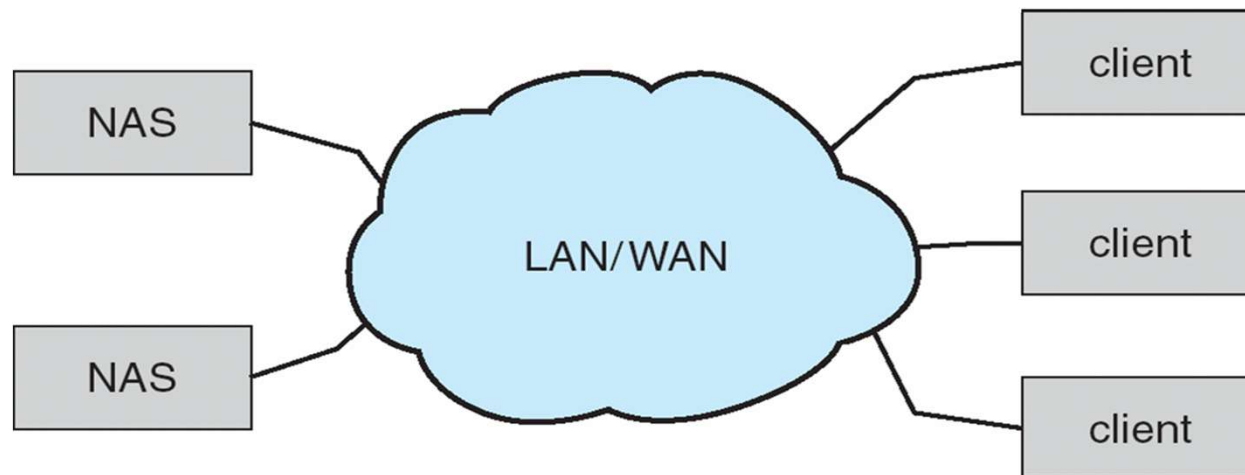
- Elevado custo de implementação e de gestão
- Falta de um “standard”, para a sua implementação é necessário vários produtos de software e de hardware de diferentes vendedores, o que complica e aumenta o custo da implementação.

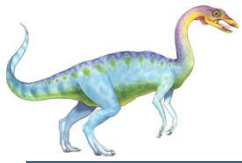




Network-Attached Storage

- O NAS, é um dispositivo dedicado ao armazenamento de ficheiros dentro de uma rede disponibilizando o acesso aos dados aos clientes dessa rede. Basicamente o NAS é um computador conectado á rede com o objetivo e funcionalidade única de fornecer serviços de armazenamento de dados a outros dispositivos de rede.
- 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)





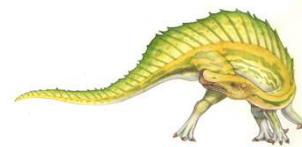
Network-Attached Storage cont

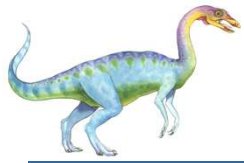
■ Vantagens

- Melhora o desempenho dos outros servidores na rede, pois o compartilhamento de ficheiros é feito exclusivamente pelo NAS, poupando trabalho aos outros servidores
- Permite aumentar a disponibilidade de dados usando [clustering](#) e [RAID](#)
- Promove a heterogeneidade, permitindo o acesso a clientes usando diferentes SO's e arquitecturas
- Baixo custo, tanto de implementação como de gestão

■ Desvantagens

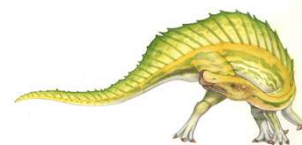
- A principal desvantagem é o facto do NAS estar limitado pelo hardware, isto é, existe um limite de carga ate onde o servidor pode ser sujeito, quando atingir esse nível limitado pelo hardware, por exemplo pelo CPU, o NAS fica comprometido. Este problema em um sistema de servidores seria facilmente solucionado adicionando um servidor no cluster, enquanto que no NAS não é possível
- Uso intensivo do NAS diminui a performance da rede, prejudicando o resto dos clientes da rede
- É um sistema de armazenamento partilhado, o administrador tem que definir quotas para os utilizadores

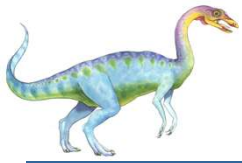




Disk Scheduling

- 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





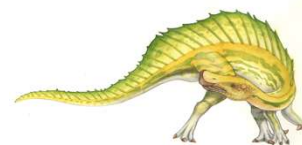
Disk Scheduling (Cont.)

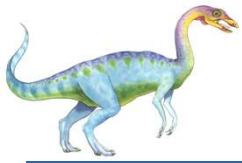
- There are many sources of disk I/O request
 - OS
 - System processes
 - Users processes
- I/O request includes input or output mode, disk address, memory address, number of sectors to transfer
- OS maintains queue of requests, per disk or device
- Idle disk can immediately work on I/O request, busy disk means work must queue
 - Optimization algorithms only make sense when a queue exists
- 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



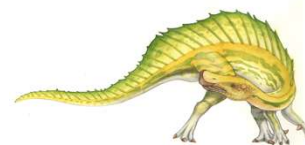
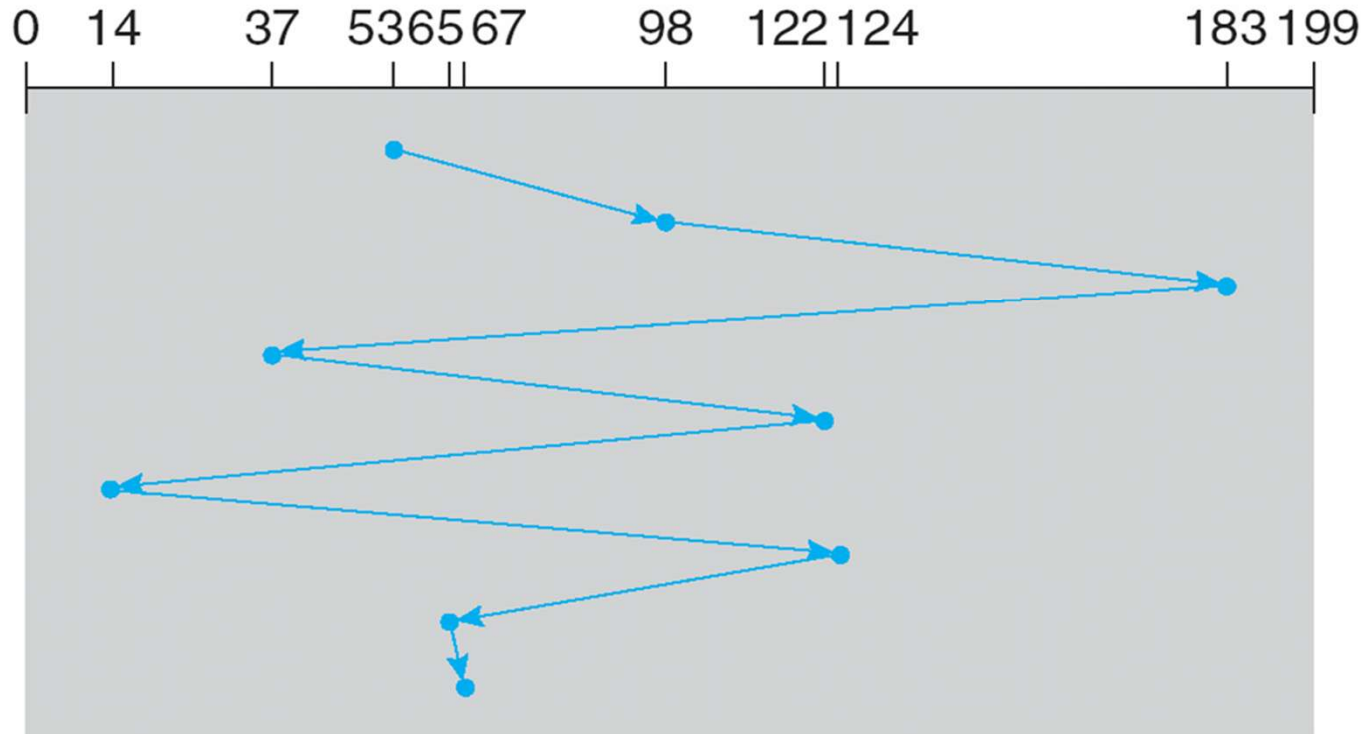


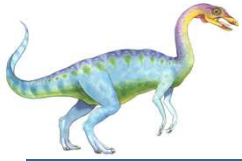
FCFS

Illustration shows total head movement of 640 cylinders

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

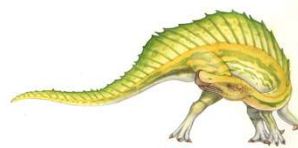
head starts at 53

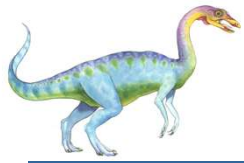




SSTF

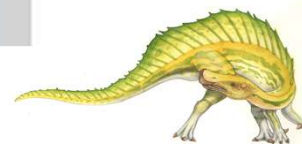
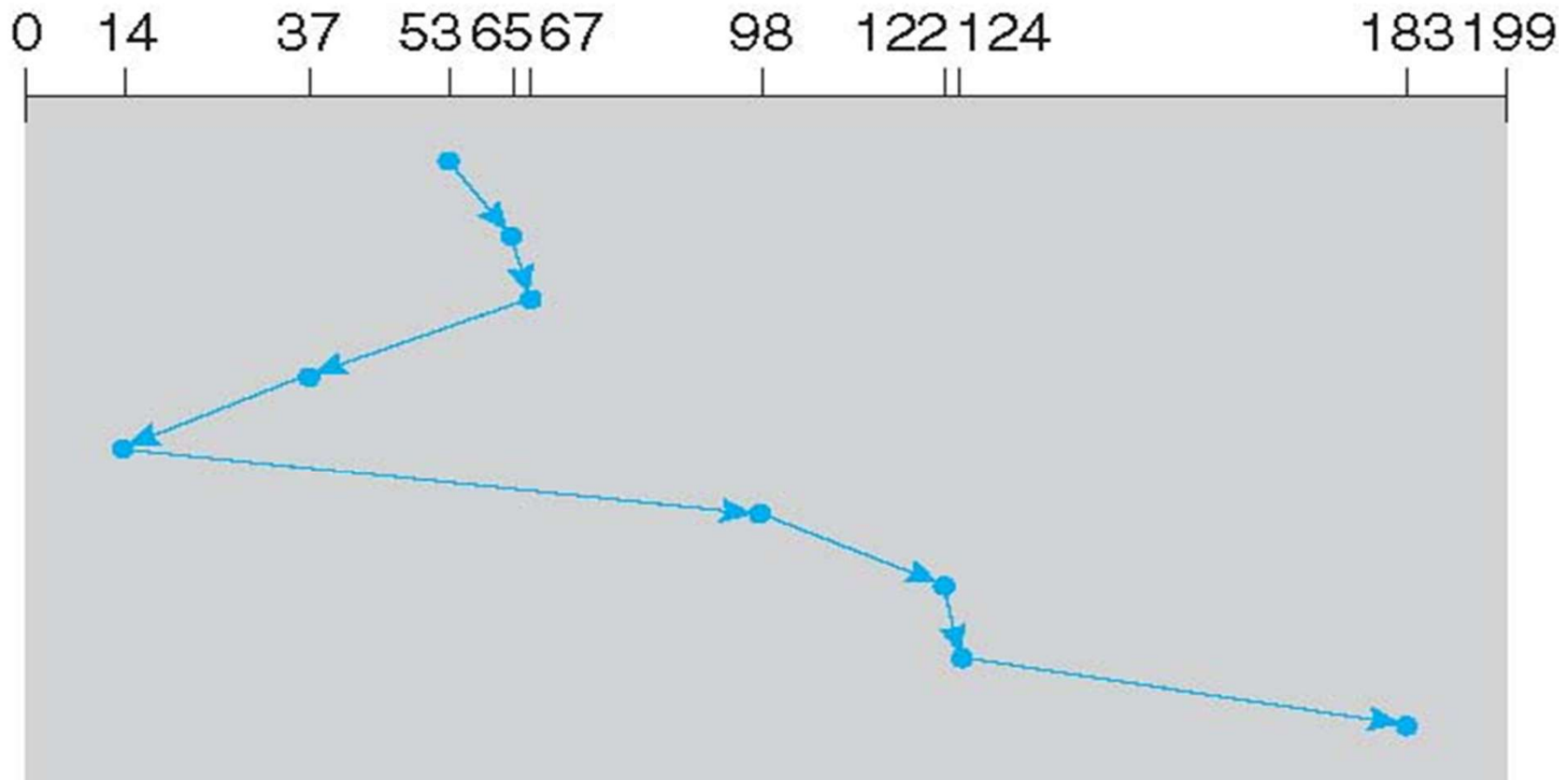
- Shortest Seek Time First selects the request with the minimum seek time from the current head position
- SSTF scheduling is a form of SJF scheduling; may cause starvation of some requests
- Illustration shows total head movement of 236 cylinders

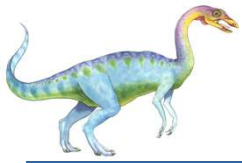




SSTF (Cont.)

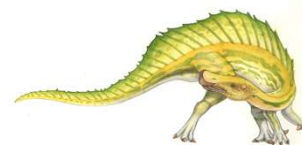
queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53





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

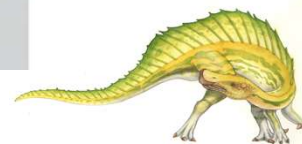
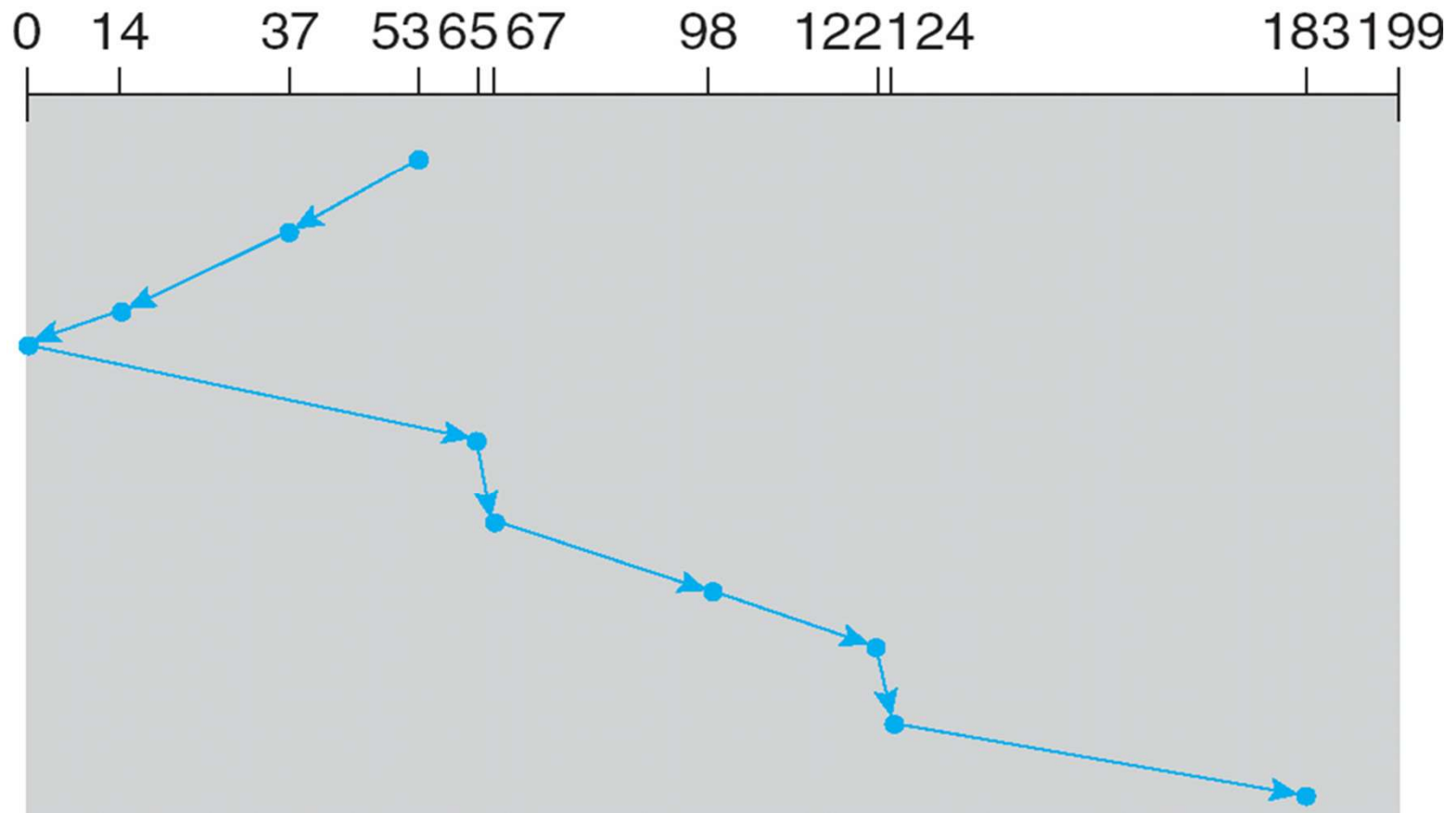


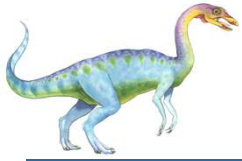


SCAN (Cont.)

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

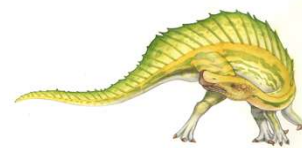
head starts at 53





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?

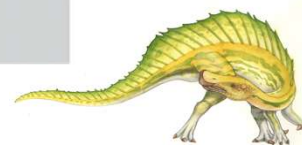
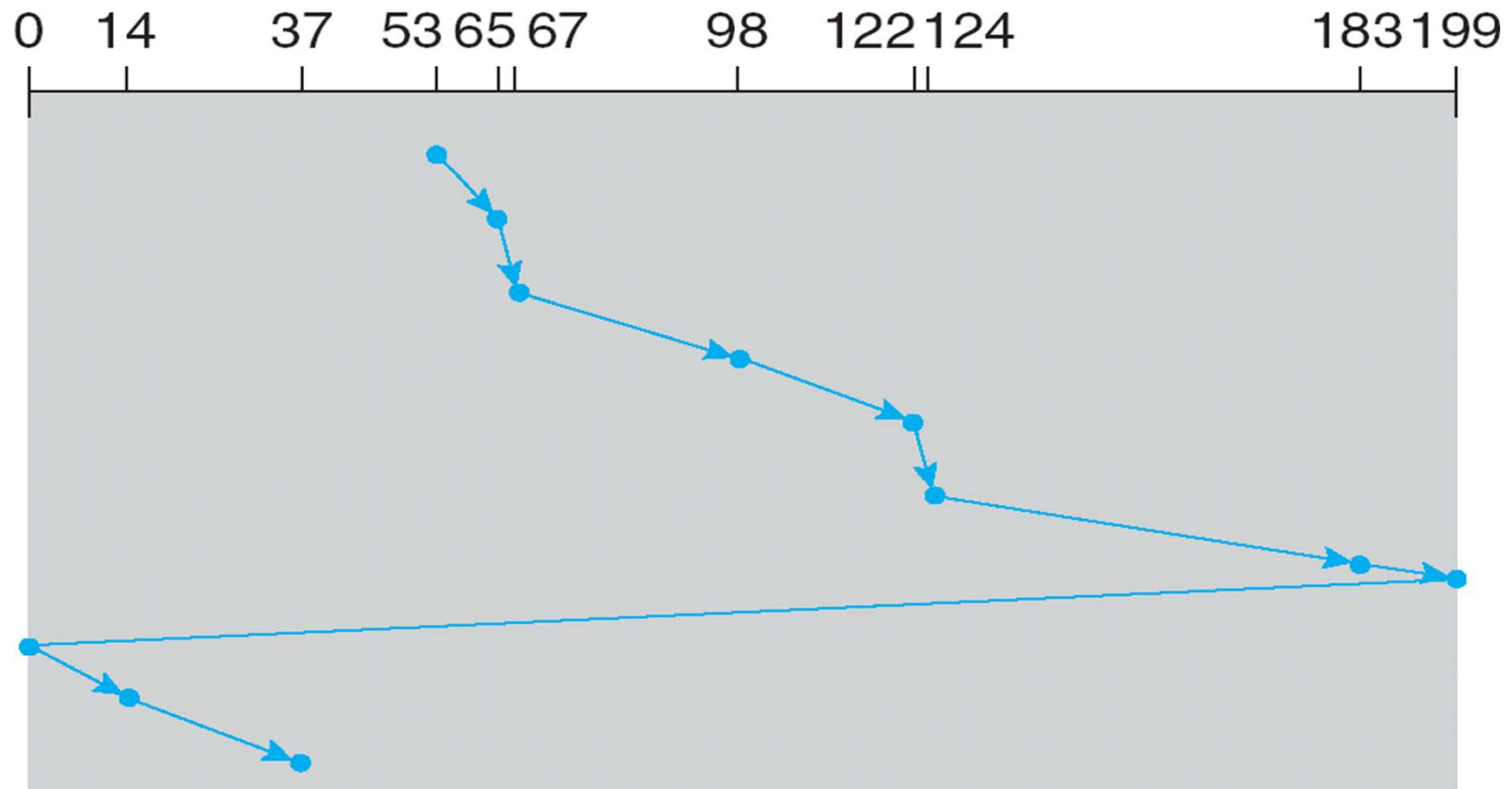


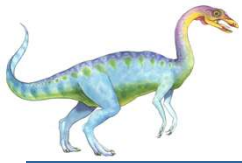


C-SCAN (Cont.)

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

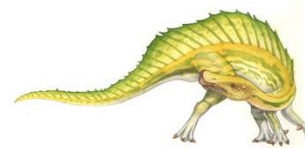
head starts at 53





C-LOOK

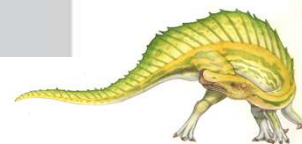
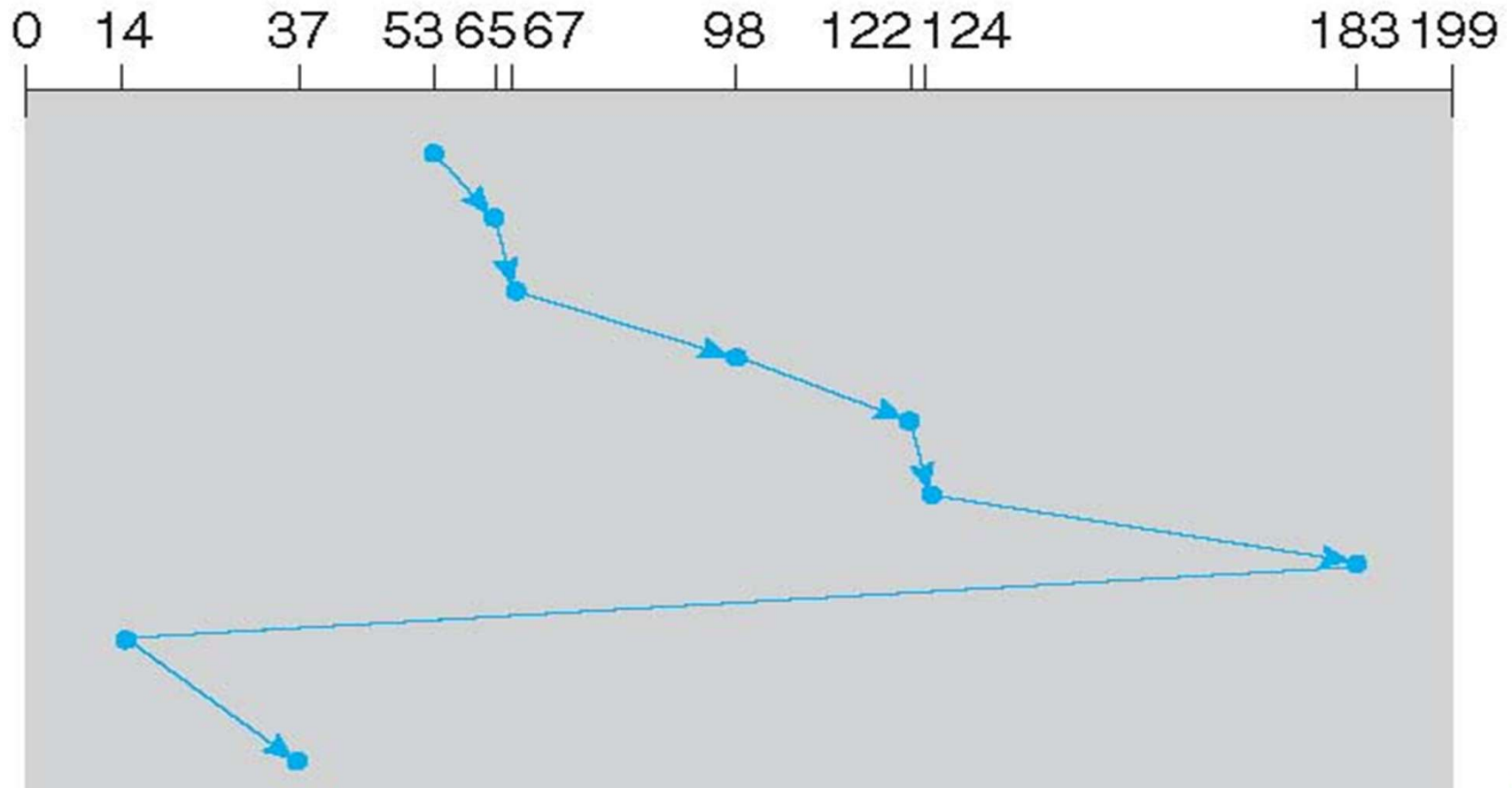
- LOOK a version of SCAN, C-LOOK a version of C-SCAN
- Arm only goes as far as the last request in each direction, then reverses direction immediately, without first going all the way to the end of the disk
- Total number of cylinders?

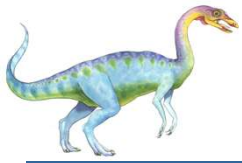




C-LOOK (Cont.)

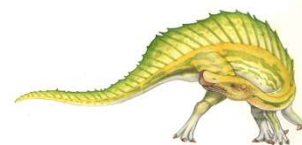
queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53

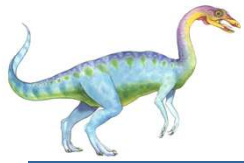




Selecting a Disk-Scheduling Algorithm

- 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
- 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
- How does disk-based queuing effect OS queue ordering efforts?





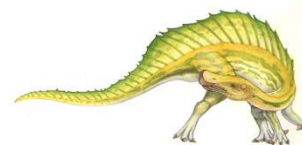
Disk Management

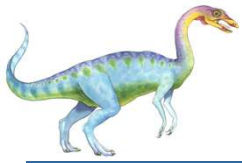
- **Low-level formatting**, or **physical formatting** — Dividing a disk into sectors that the disk controller can read and write
 - Each sector can hold header information, plus data, plus error correction code (**ECC**)
 - Usually 512 bytes of data but can be selectable

- 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

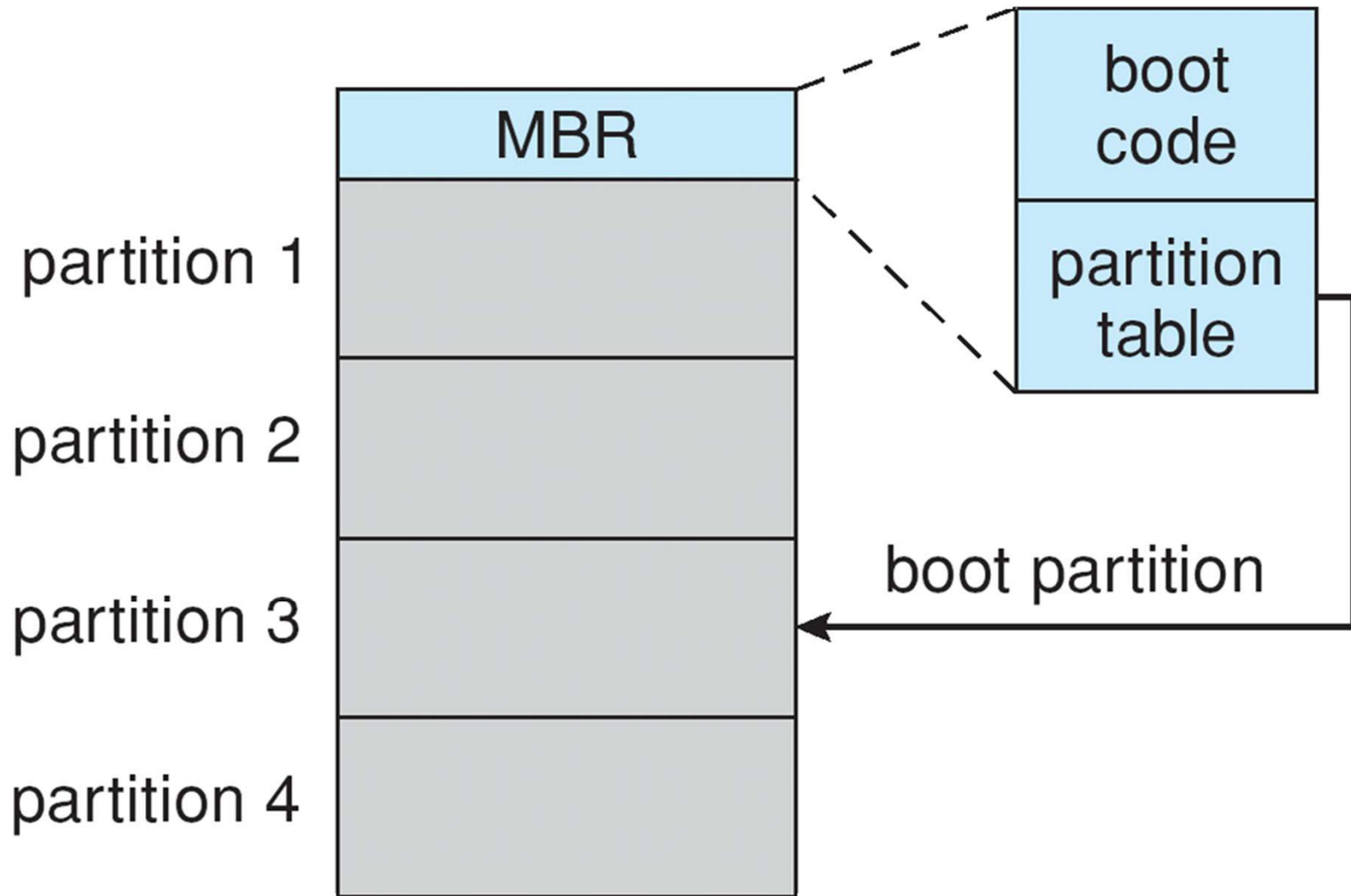
- Boot block initializes system
 - The bootstrap is stored in ROM
 - **Bootstrap loader** program stored in boot blocks of boot partition

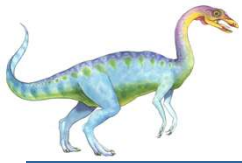
- Methods such as **sector sparing** used to handle bad blocks





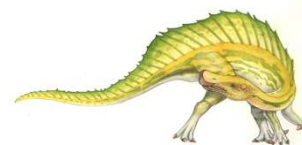
Booting from a Disk in Windows 2000

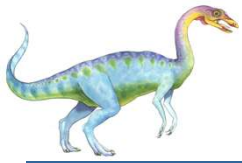




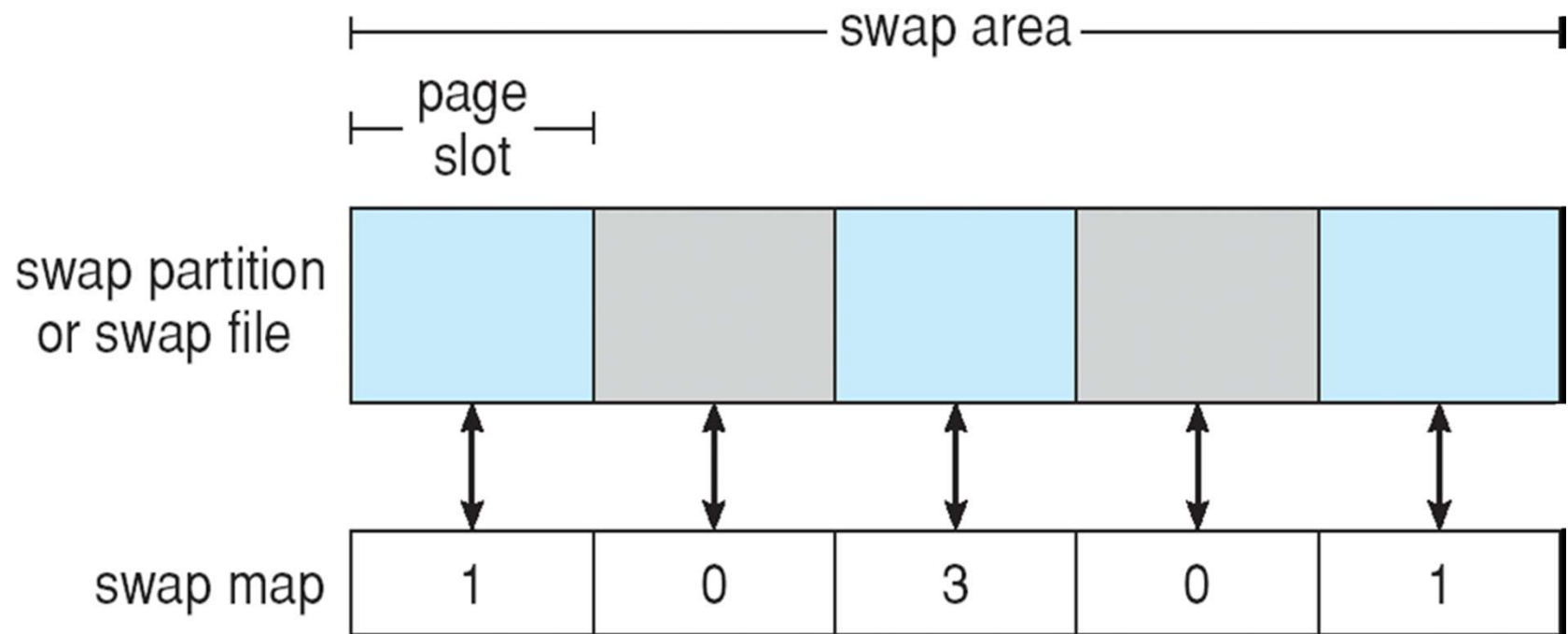
Swap-Space Management

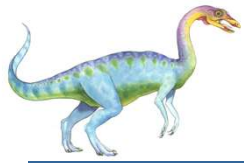
- Swap-space — Virtual memory uses disk space as an extension of main memory
 - Less common now due to memory capacity increases
- Swap-space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition (raw)
- Swap-space management
 - 4.3BSD allocates swap space when process starts; holds text segment (the program) and data segment
 - Kernel uses **swap maps** to track swap-space use
 - Solaris 2 allocates swap space only when a dirty page is forced out of physical memory, not when the virtual memory page is first created
 - ▶ File data written to swap space until write to file system requested
 - ▶ Other dirty pages go to swap space due to no other home
 - ▶ Text segment pages thrown out and reread from the file system as needed
- What if a system runs out of swap space?
- Some systems allow multiple swap spaces





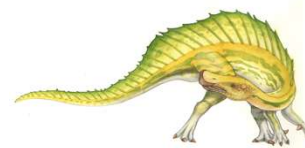
Data Structures for Swapping on Linux Systems

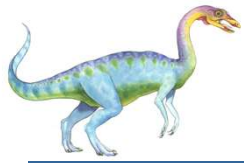




RAID Structure

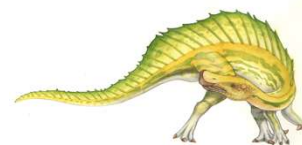
- RAID – multiple disk drives provides reliability via **redundancy**
- Increases the **mean time to failure**
- Frequently combined with **NVRAM** to improve write performance
- RAID is arranged into six different levels





RAID (Cont.)

- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively
- Disk **striping** uses a group of disks as one storage unit
- 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

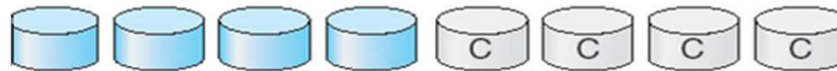




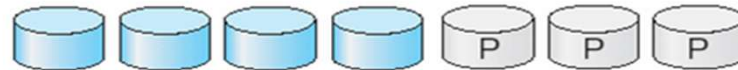
RAID Levels



(a) RAID 0: non-redundant striping.



(b) RAID 1: mirrored disks.



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



(d) RAID 3: bit-interleaved parity.



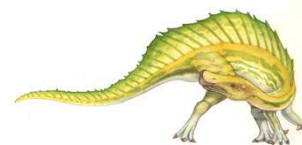
(e) RAID 4: block-interleaved parity.

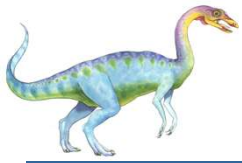


(f) RAID 5: block-interleaved distributed parity.

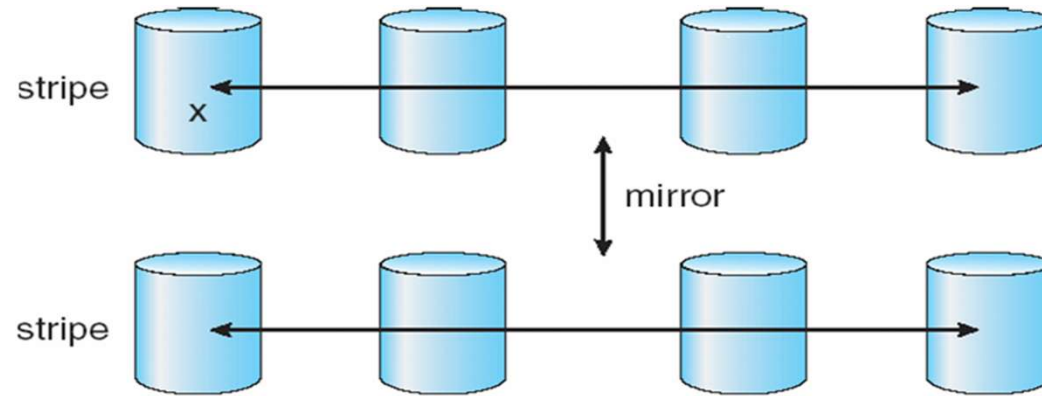


(g) RAID 6: P + Q redundancy.

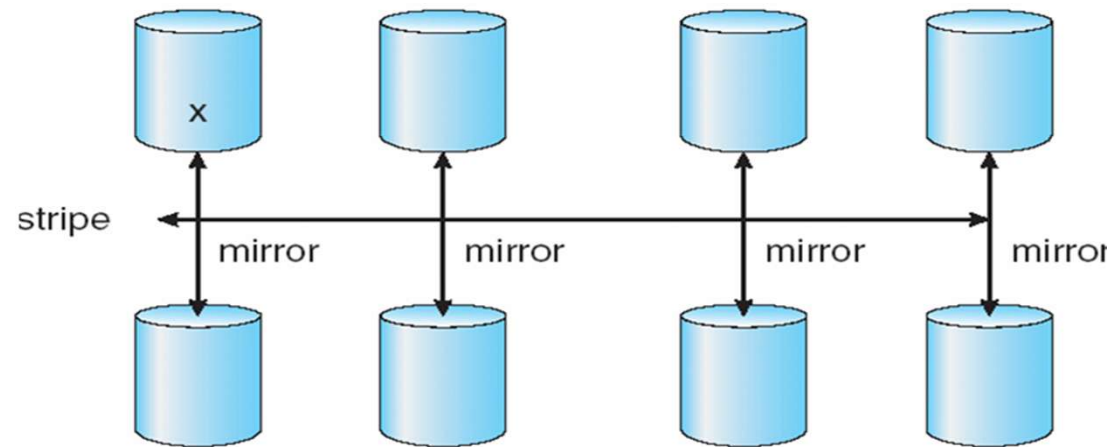




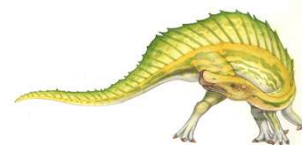
RAID (0 + 1) and (1 + 0)

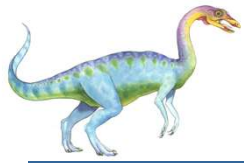


a) RAID 0 + 1 with a single disk failure.



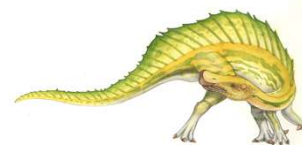
b) RAID 1 + 0 with a single disk failure.

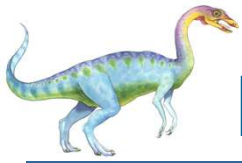




WORM Disks

- The data on read-write disks can be modified over and over
- **WORM** (“Write Once, Read Many Times”) disks can be written only once
- Thin aluminum film sandwiched between two glass or plastic platters
- To write a bit, the drive uses a laser light to burn a small hole through the aluminum; information can be destroyed by not altered
- Very durable and reliable
- **Read-only disks**, such as CD-ROM and DVD, come from the factory with the data pre-recorded



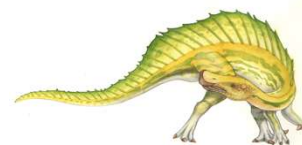


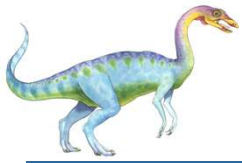
Hierarchical Storage Management (HSM)

- A hierarchical storage system extends the storage hierarchy beyond primary memory and secondary storage to incorporate tertiary storage — usually implemented as a jukebox of tapes or removable disks.

- Usually incorporate tertiary storage by extending the file system
 - Small and frequently used files remain on disk
 - Large, old, inactive files are archived to the jukebox

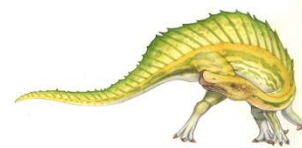
- HSM is usually found in supercomputing centers and other large installations that have enormous volumes of data.

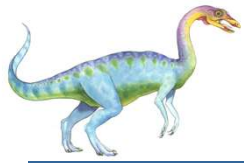




Speed

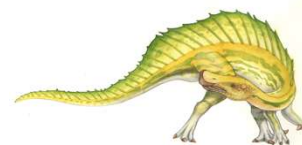
- Two aspects of speed in tertiary storage are bandwidth and latency.
- Bandwidth is measured in bytes per second.
 - **Sustained bandwidth** – average data rate during a large transfer; # of bytes/transfer time
Data rate when the data stream is actually flowing
 - **Effective bandwidth** – average over the entire I/O time, including `seek()` or `locate()`, and cartridge switching
Drive's overall data rate

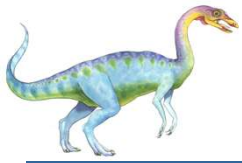




Speed (Cont.)

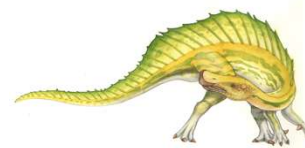
- **Access latency** – amount of time needed to locate data
 - Access time for a disk – move the arm to the selected cylinder and wait for the rotational latency; < 35 milliseconds
 - Access on tape requires winding the tape reels until the selected block reaches the tape head; tens or hundreds of seconds
 - Generally say that random access within a tape cartridge is about a thousand times slower than random access on disk
- The low cost of tertiary storage is a result of having many cheap cartridges share a few expensive drives
- A removable library is best devoted to the storage of infrequently used data, because the library can only satisfy a relatively small number of I/O requests per hour

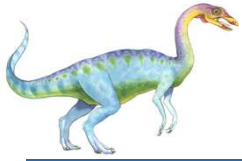




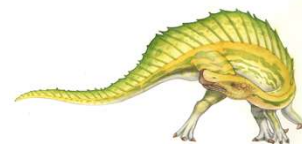
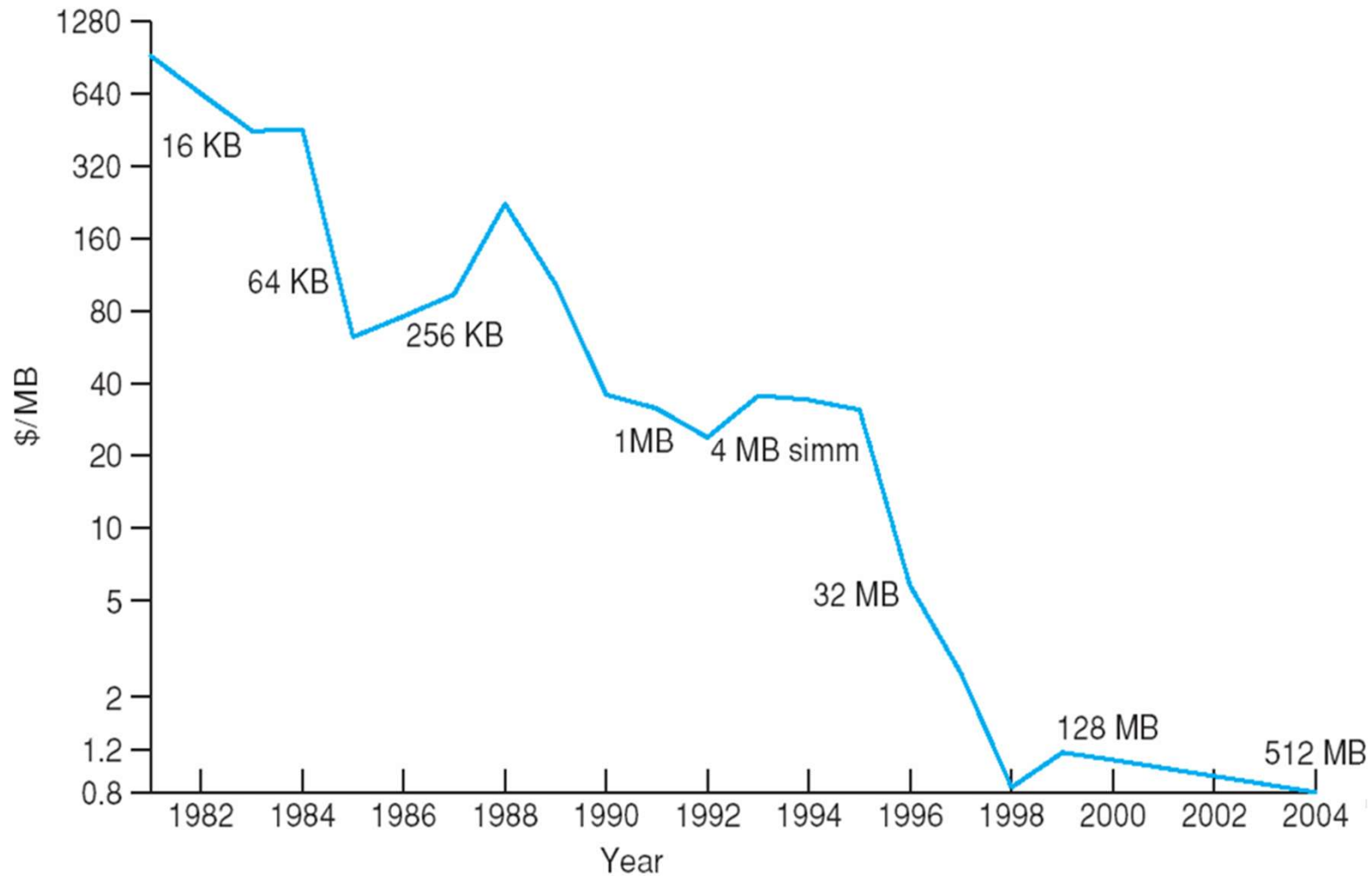
Cost

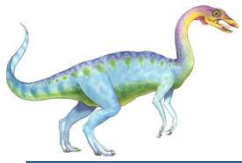
- Main memory is much more expensive than disk storage
- The cost per megabyte of hard disk storage is competitive with magnetic tape if only one tape is used per drive
- The cheapest tape drives and the cheapest disk drives have had about the same storage capacity over the years
- Tertiary storage gives a cost savings only when the number of cartridges is considerably larger than the number of drives



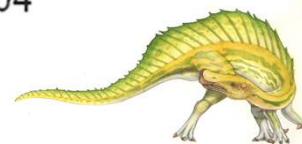
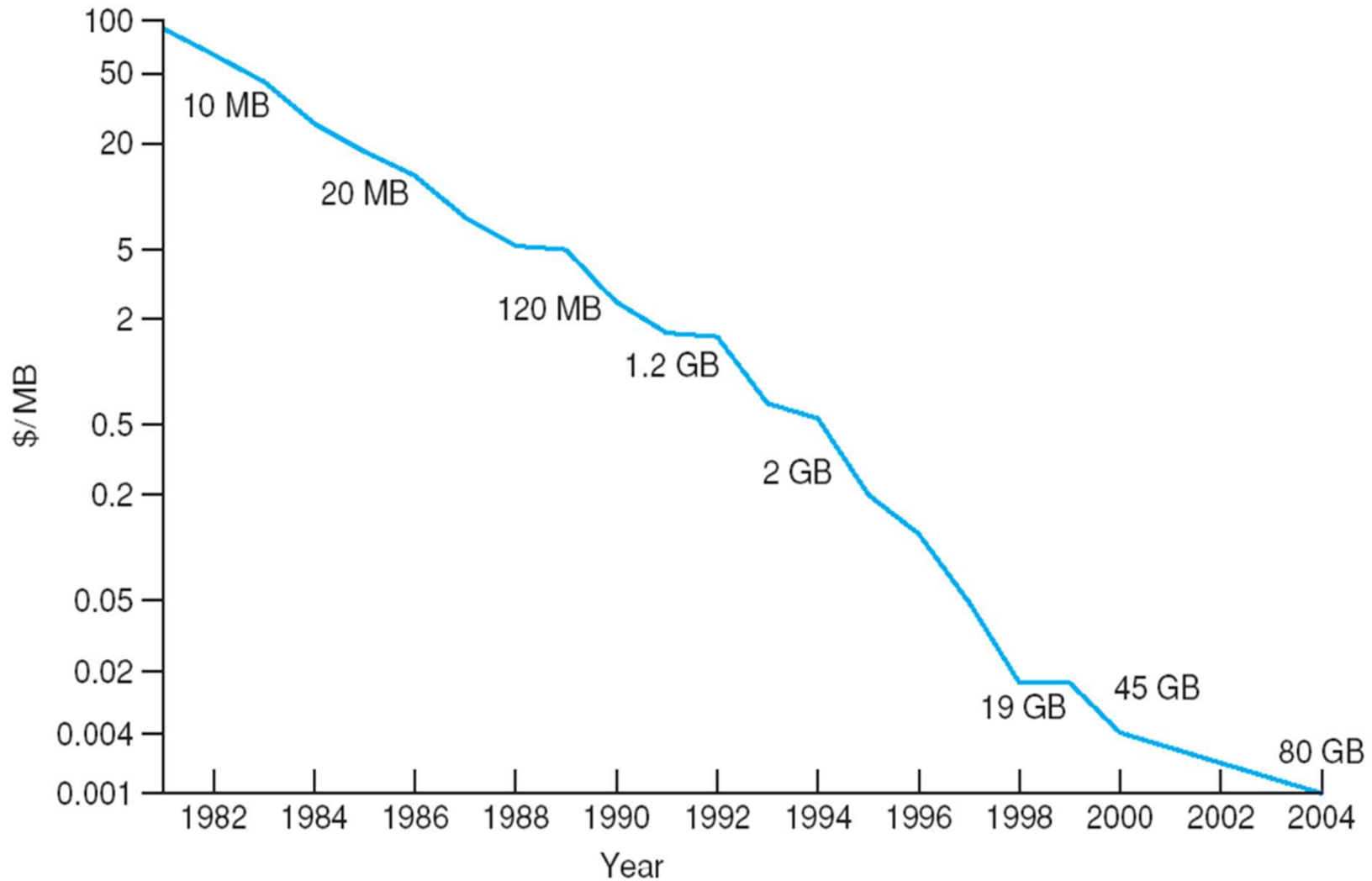


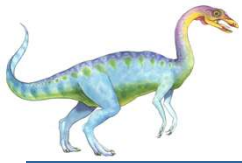
Price per Megabyte of DRAM From 1981 to 2004



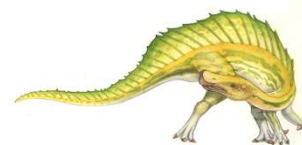
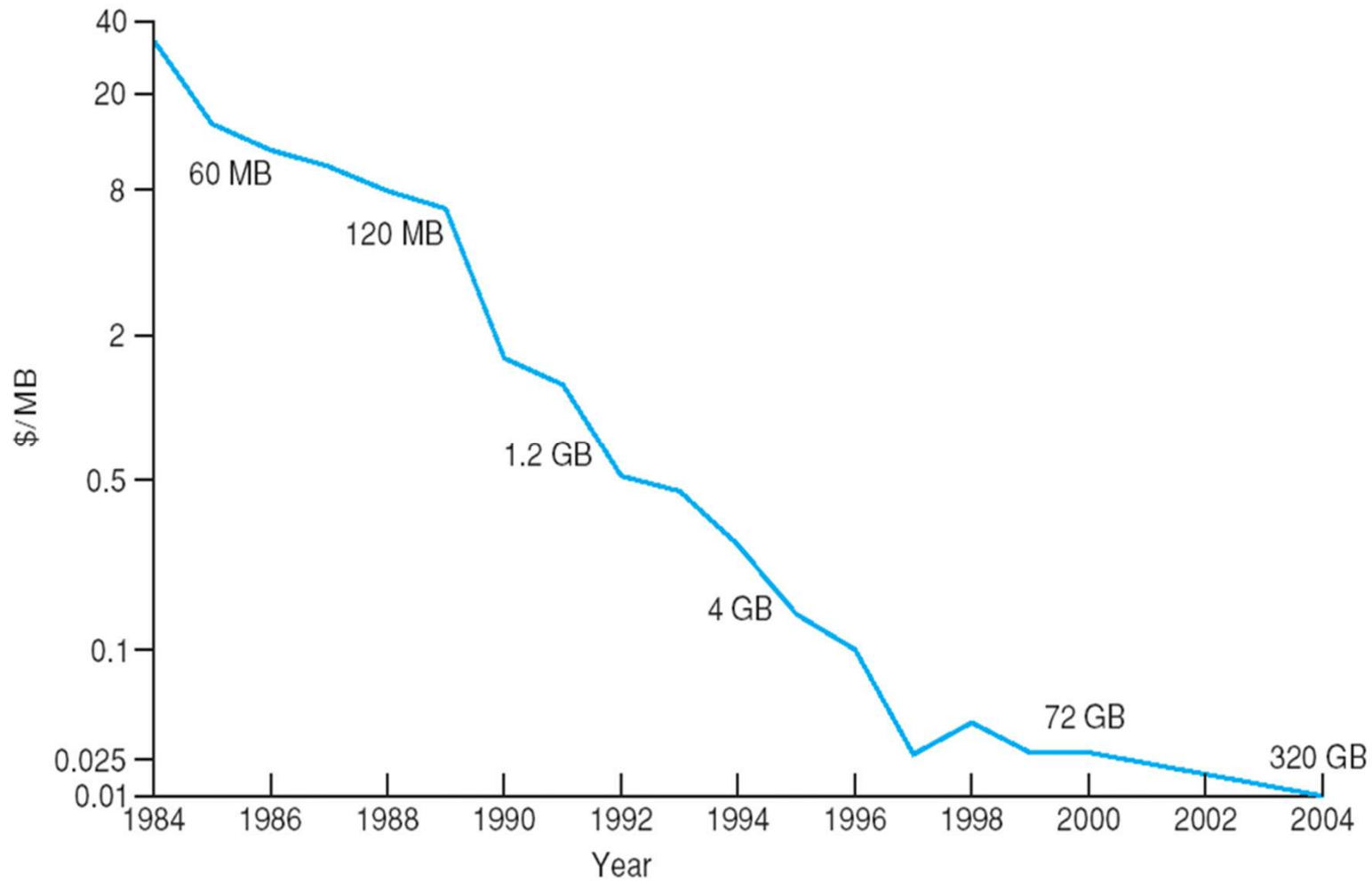


Price per Megabyte of Magnetic Hard Disk From 1981 to 2004





Price per Megabyte of a Tape Drive From 1984-2000



End of Chapter 12

