

Overview

- ▶ Continuation from Monday (File system implementation)



Efficiency and Performance

- ▶ Efficiency dependent on:
 - disk allocation and directory algorithms
 - types of data kept in file's directory entry
- ▶ Performance
 - disk cache – separate section of main memory for frequently used blocks
 - free-behind and read-ahead – techniques to optimize sequential access
 - improve PC performance by dedicating section of memory as virtual disk, or RAM disk

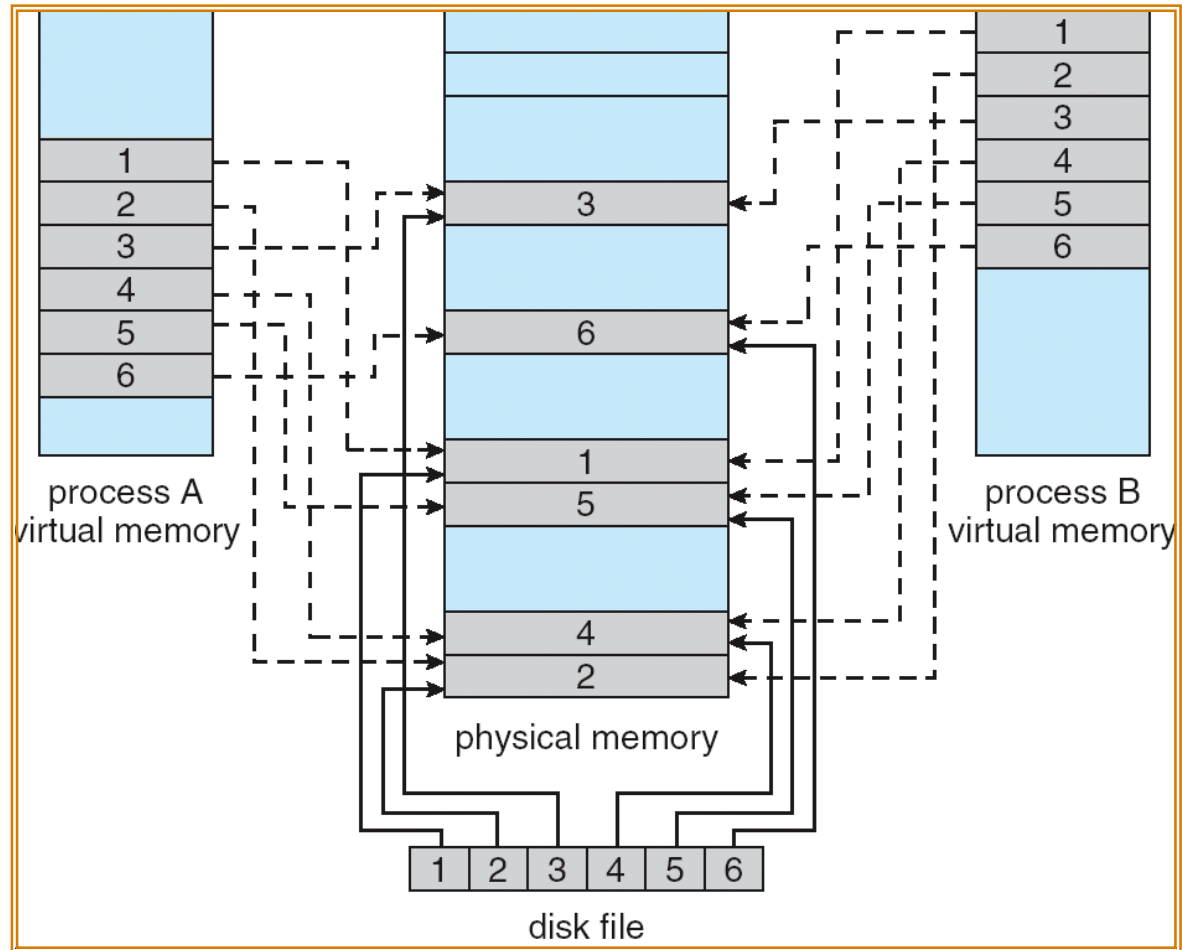


Memory-Mapped Files

- ▶ Memory-mapped file I/O allows file I/O to be treated as routine memory access by mapping a disk block to a page in memory
- ▶ A file is initially read using demand paging. A page-sized portion of the file is read from the file system into a physical page. Subsequent reads/writes to/from the file are treated as ordinary memory accesses.
- ▶ Simplifies file access by treating file I/O through memory rather than `read()` `write()` system calls
- ▶ Also allows several processes to map the same file allowing the pages in memory to be shared



Memory Mapped Files

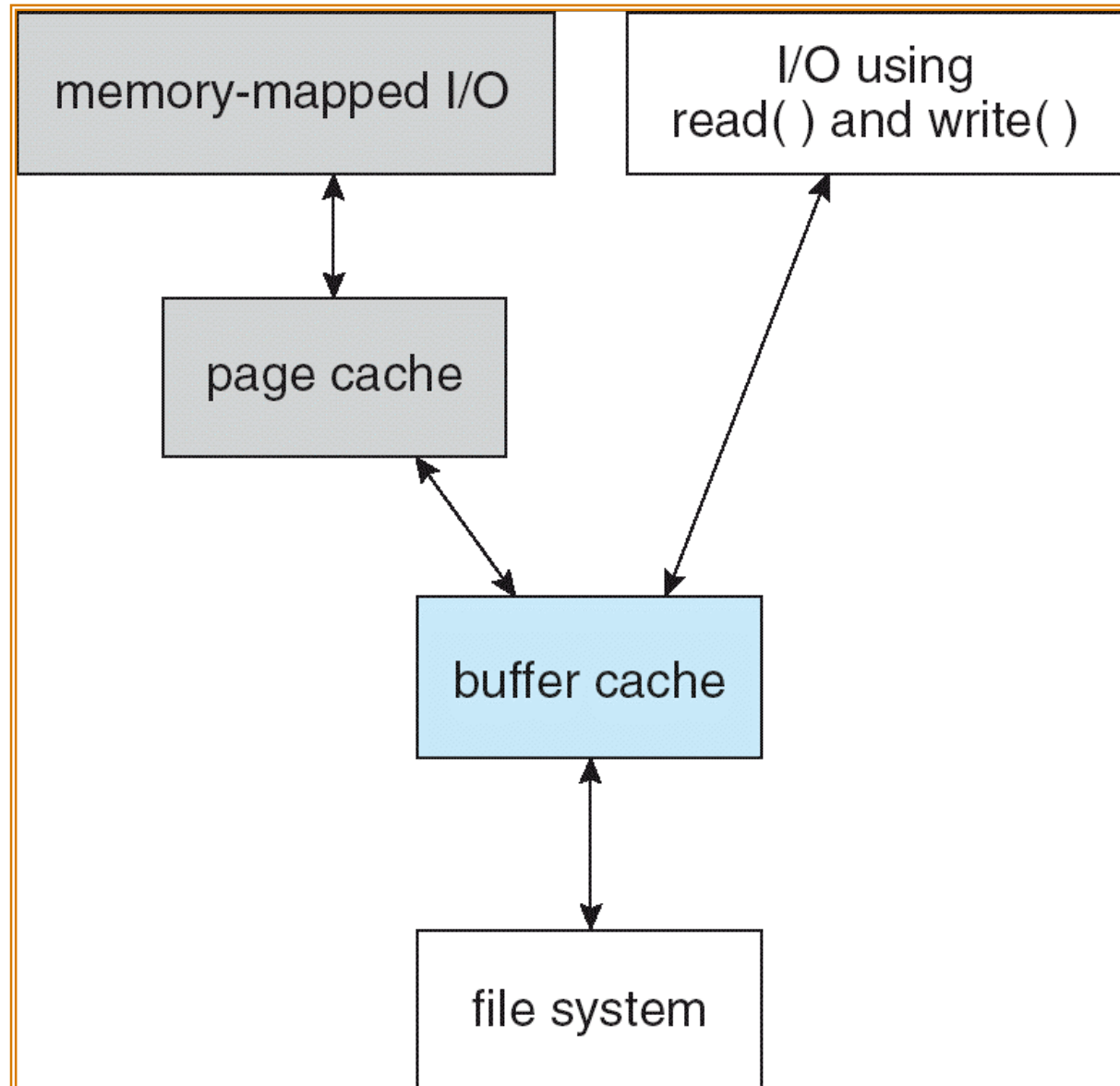


Page Cache

- ▶ A **page cache** caches pages rather than disk blocks using virtual memory techniques
- ▶ Memory-mapped I/O uses a page cache
- ▶ Routine I/O through the file system uses the buffer (disk) cache
- ▶ This leads to the following figure



I/O Without a Unified Buffer Cache

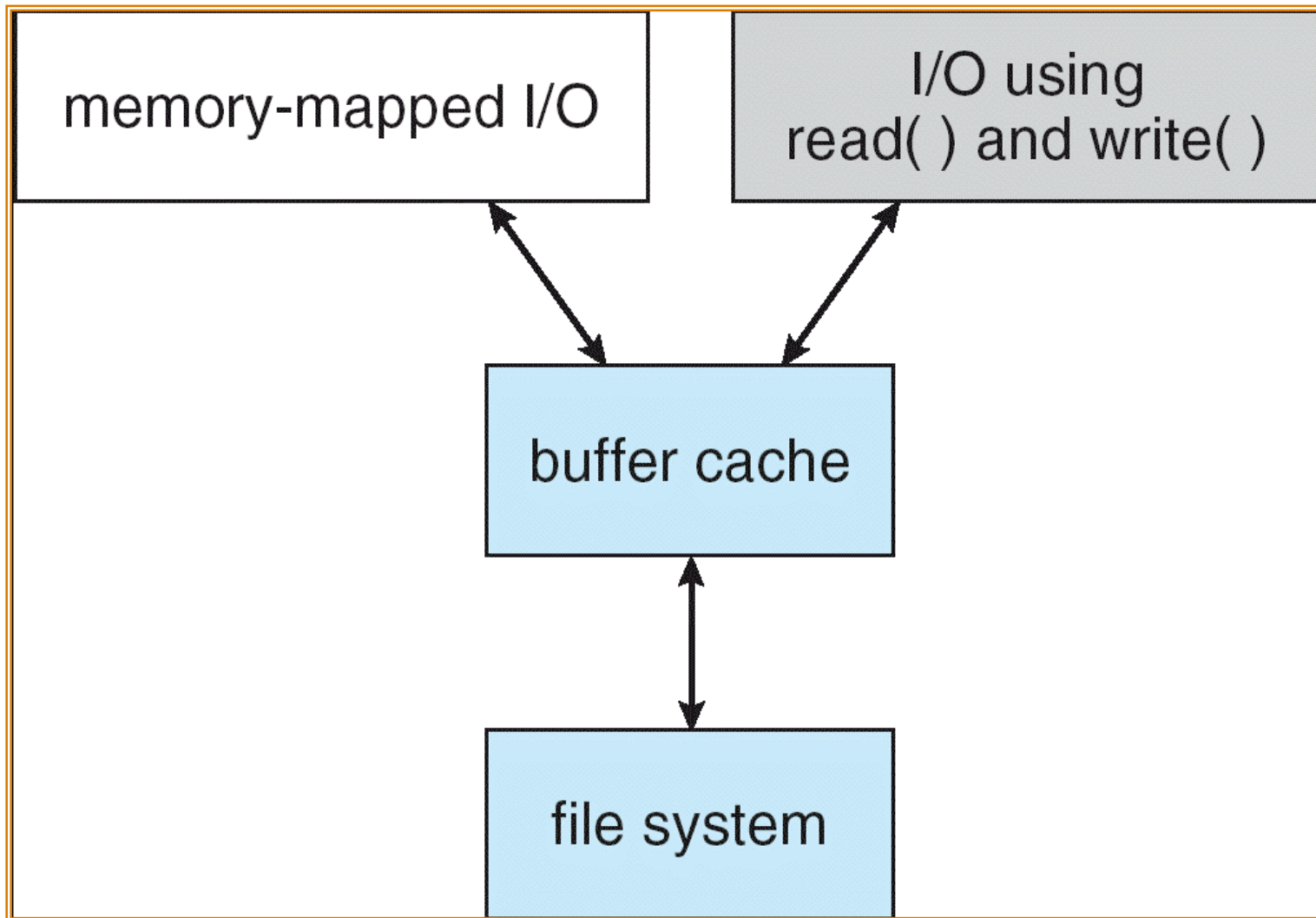


Unified Buffer Cache

- ▶ A unified buffer cache uses the same page cache to cache both memory-mapped pages and ordinary file system I/O



I/O Using a Unified Buffer Cache



Recovery

- ▶ Consistency checking – compares data in directory structure with data blocks on disk, and tries to fix inconsistencies
- ▶ Use system programs to **back up** data from disk to another storage device (floppy disk, magnetic tape, other magnetic disk, optical)
- ▶ Recover lost file or disk by **restoring** data from backup



Log Structured File Systems

- ▶ Log structured (or journaling) file systems record each update to the file system as a transaction
- ▶ All transactions are written to a log
 - A transaction is considered committed once it is written to the log
 - However, the file system may not yet be updated
- ▶ The transactions in the log are asynchronously written to the file system
 - When the file system is modified, the transaction is removed from the log
- ▶ If the file system crashes, all remaining transactions in the log must still be performed



Disk Scheduling (Continuation from Mass storage structure)

- ▶ The operating system is responsible for using hardware efficiently — for the disk drives, this means having a fast access time and disk bandwidth.
- ▶ Access time has two major components
 - *Seek time* is the time for the disk are to move the heads to the cylinder containing the desired sector.
 - *Rotational latency* is the additional time waiting for the disk to rotate the desired sector to the disk head.
- ▶ 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.)

- ▶ Several algorithms exist to schedule the servicing of disk I/O requests.
- ▶ We illustrate them with a request queue (0-199).

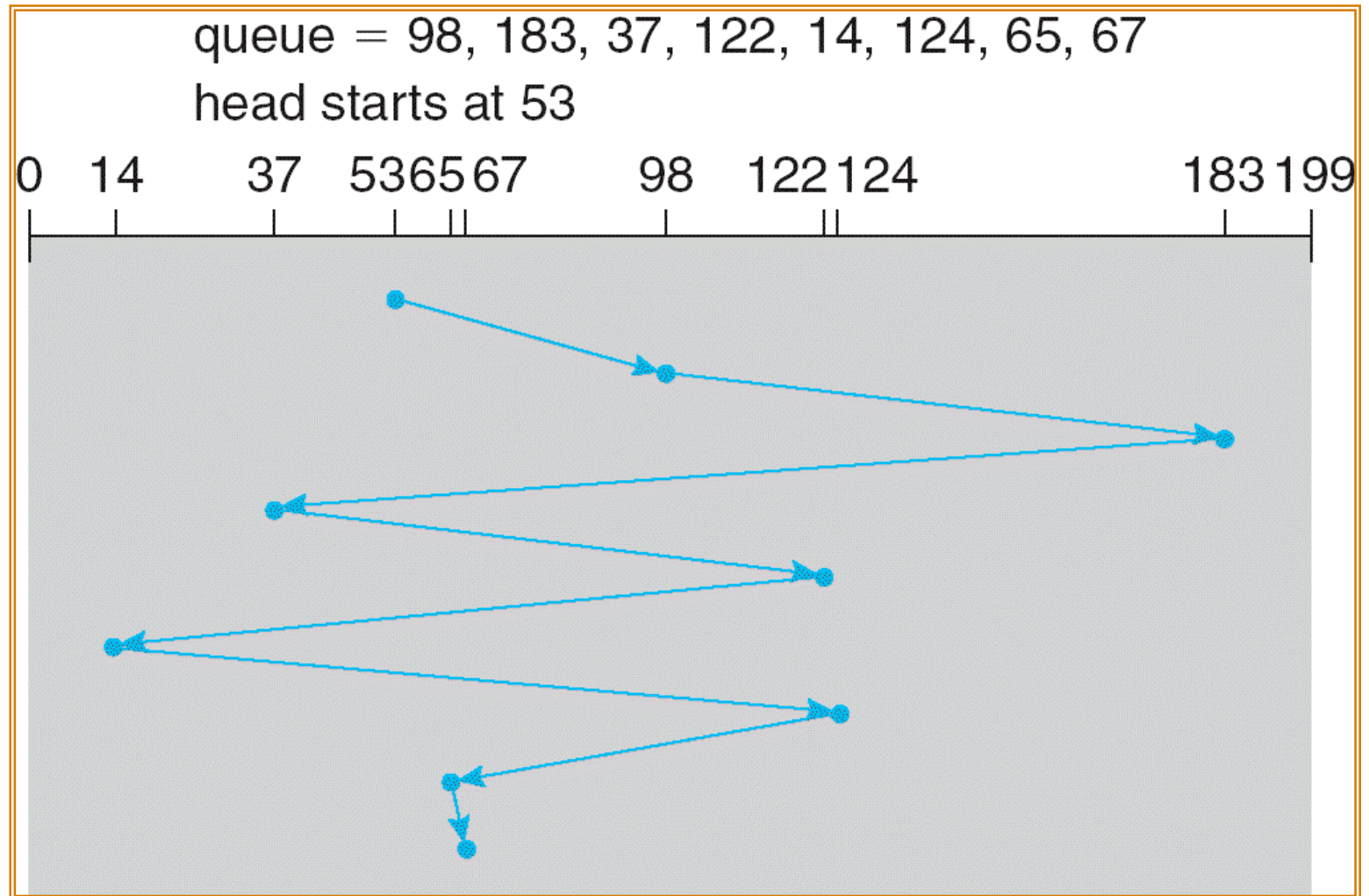
98, 183, 37, 122, 14, 124, 65, 67

Head pointer 53



FCFS

Illustration shows total head movement of 640 cylinders.



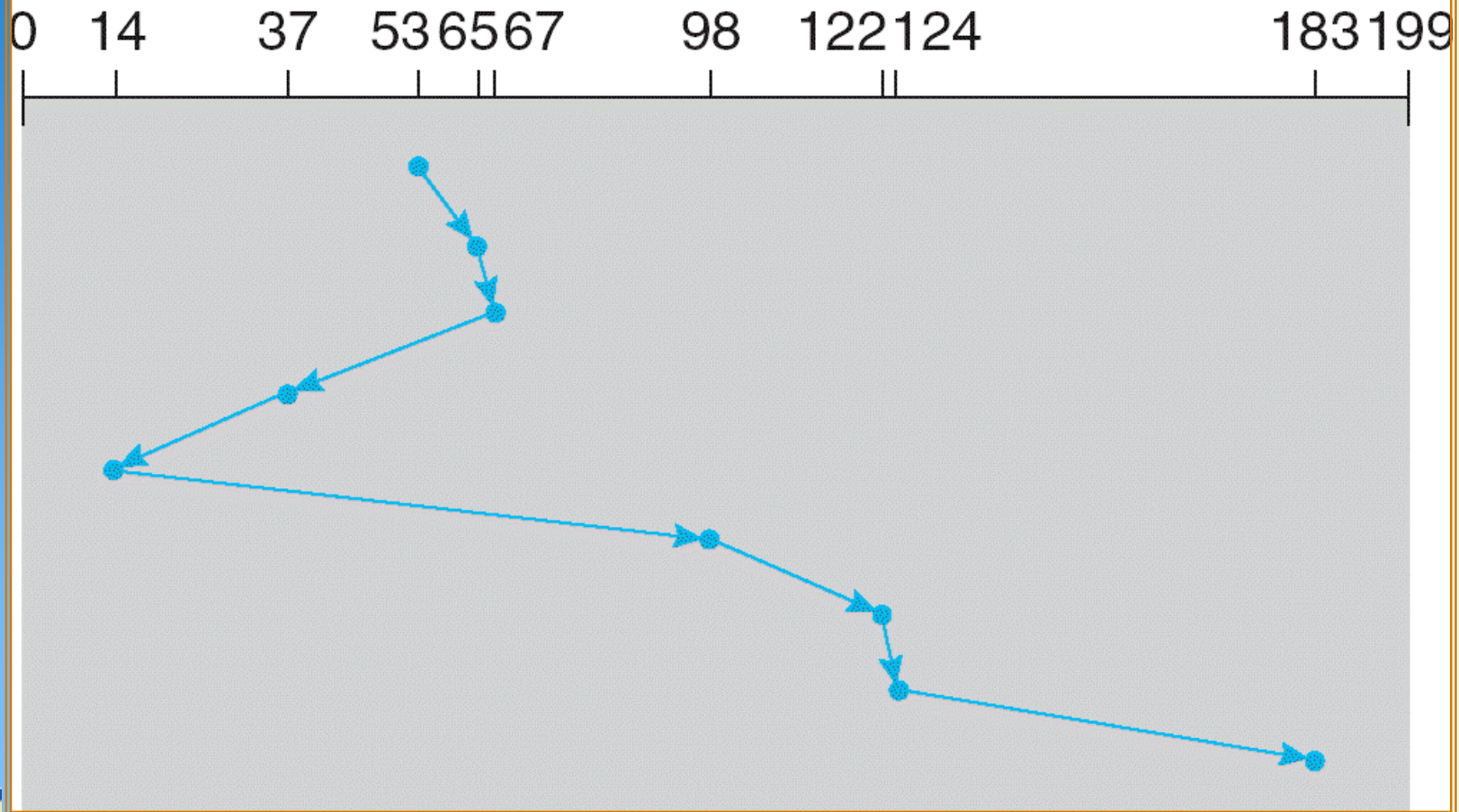
SSTF

- ▶ 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



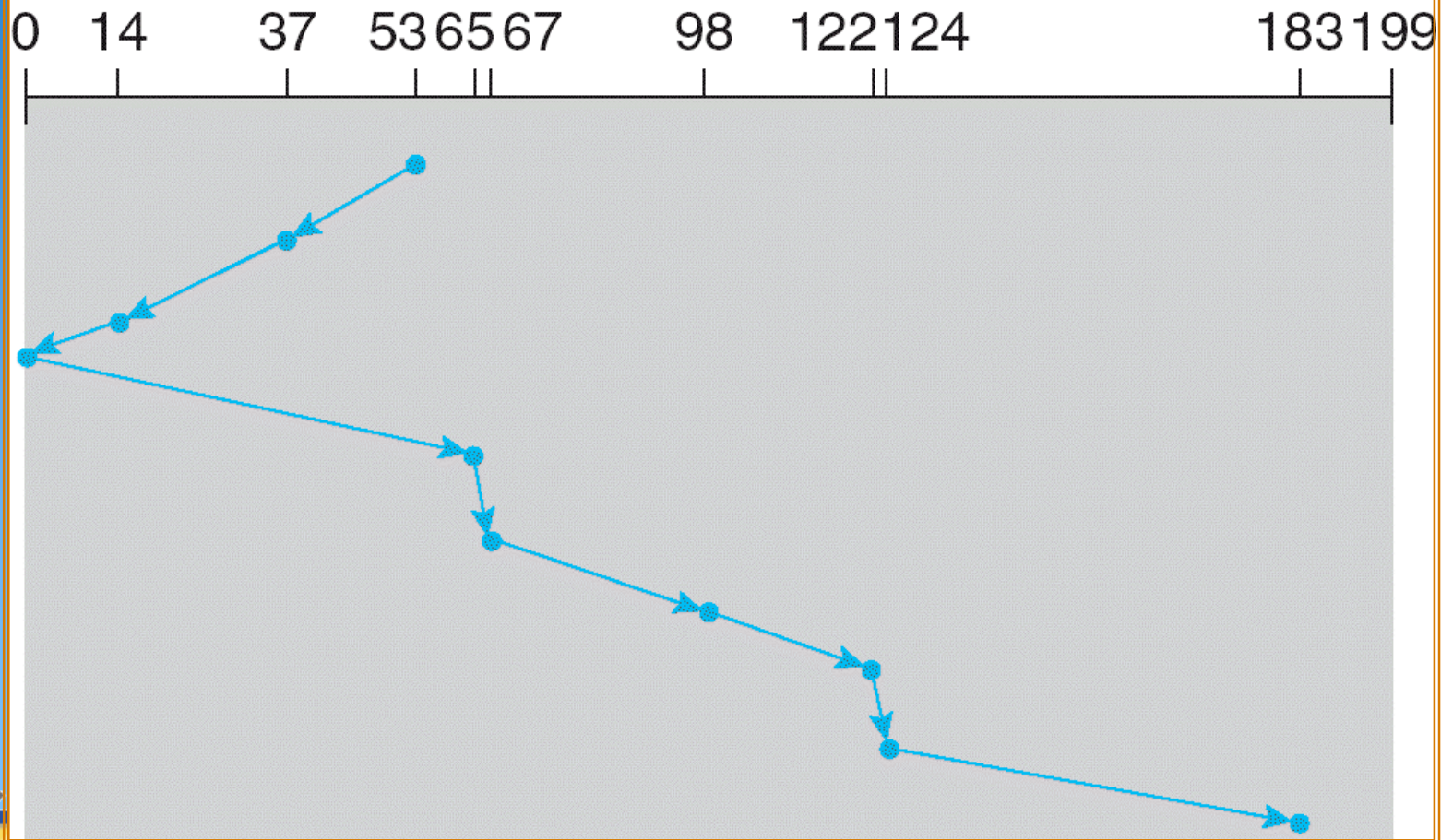
SCAN

- ▶ 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.
- ▶ Sometimes called the *elevator algorithm*.
- ▶ Illustration shows total head movement of 208 cylinders.



SCAN (Cont.)

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



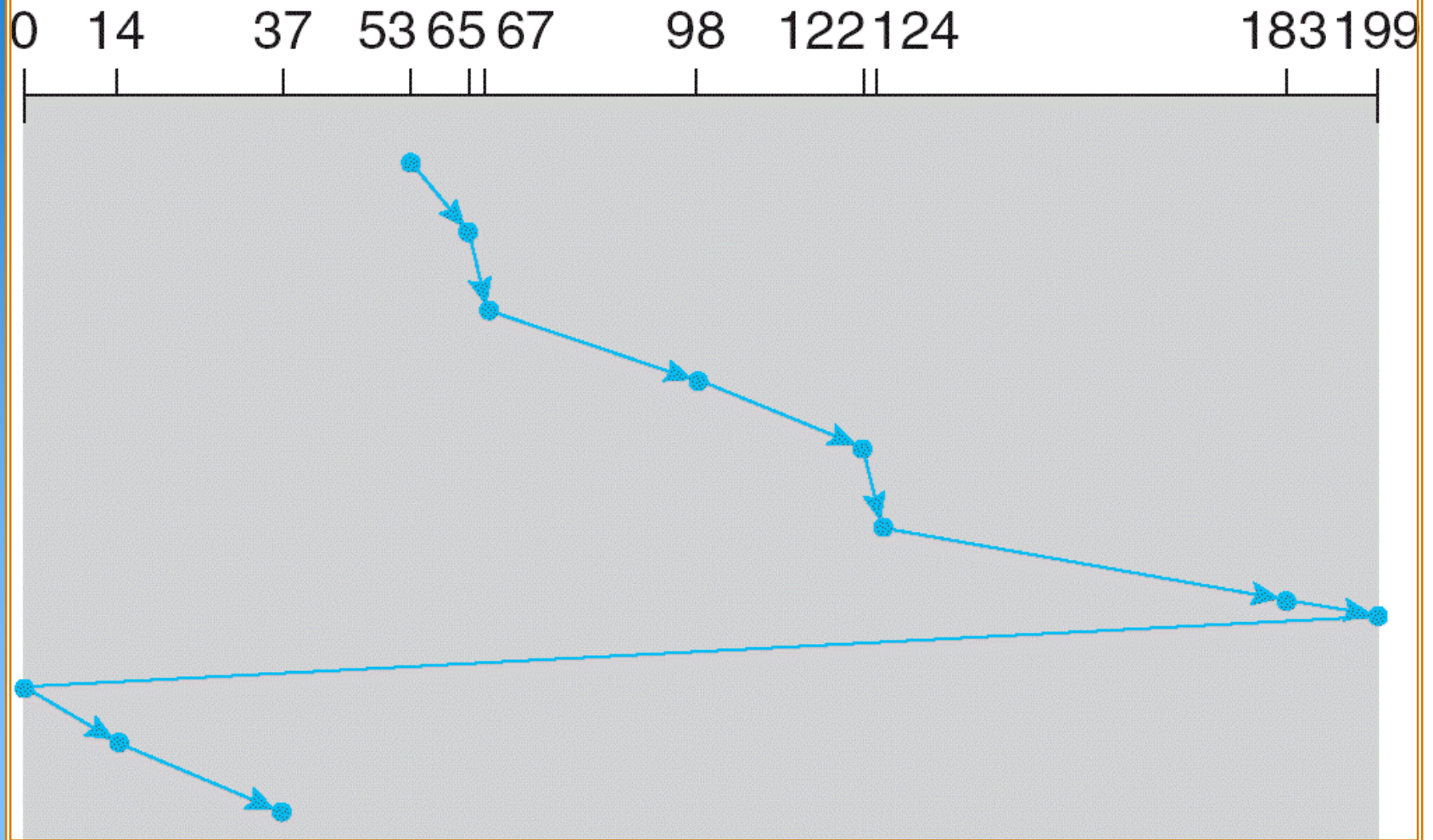
C-SCAN

- ▶ 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.



C-SCAN (Cont.)

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



C-LOOK

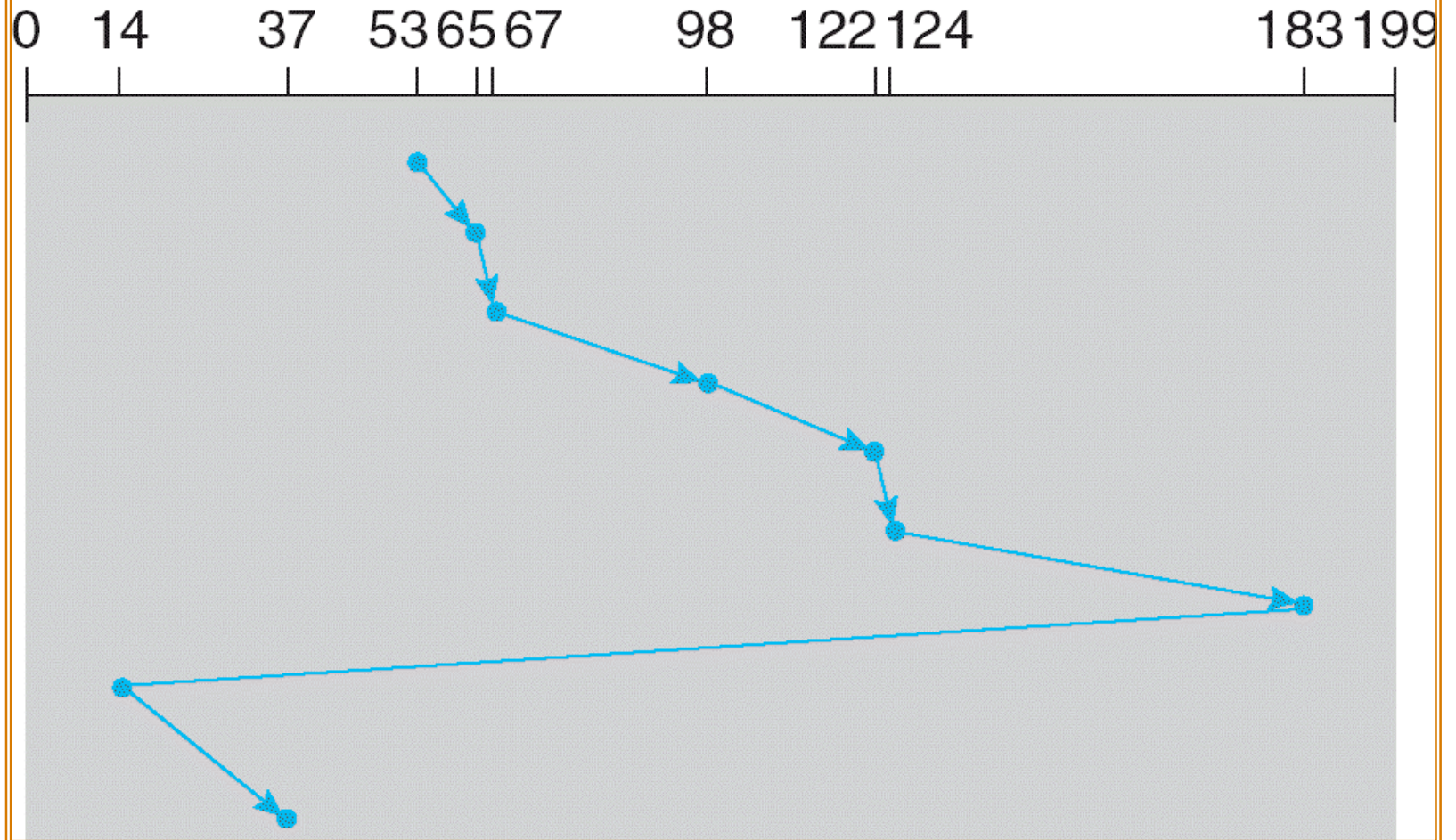
- ▶ Version of C-SCAN
- ▶ Arm only goes as far as the last request in each direction, then reverses direction immediately, without first going all the way to the end of the disk.



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.
- ▶ Performance depends on the number and types of requests.
- ▶ Requests for disk service can be influenced by the file-allocation method.
- ▶ 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.

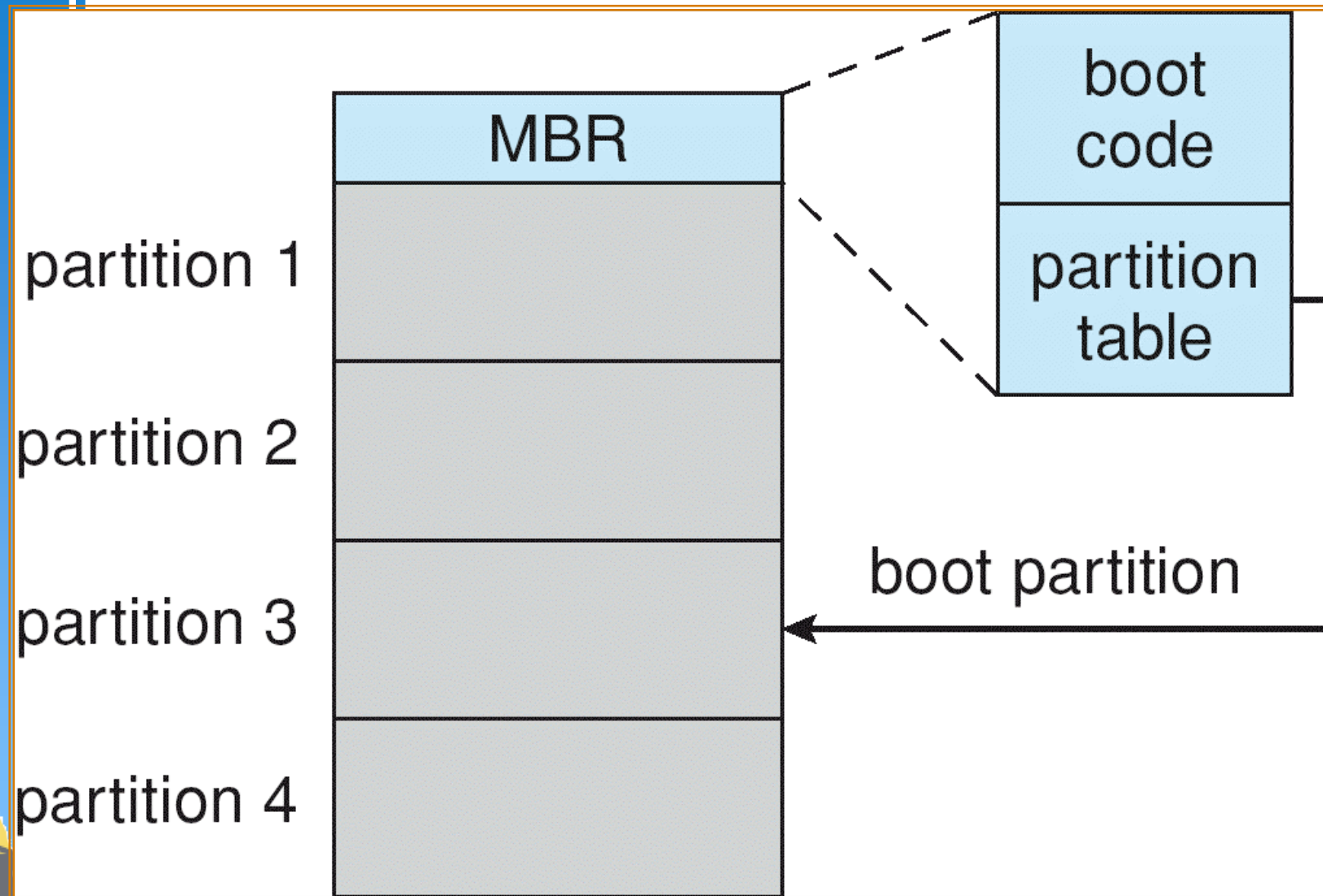


Disk Management

- ▶ *Low-level formatting, or physical formatting* — Dividing a disk into sectors that the disk controller can read and write.
- ▶ 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.
 - *Logical formatting* or “making a file system”.
- ▶ **Boot block initializes system.**
 - The bootstrap is stored in ROM.
 - *Bootstrap loader* program.
- ▶ Methods such as *sector sparing* used to handle bad blocks.



Booting from a Disk in Windows 2000

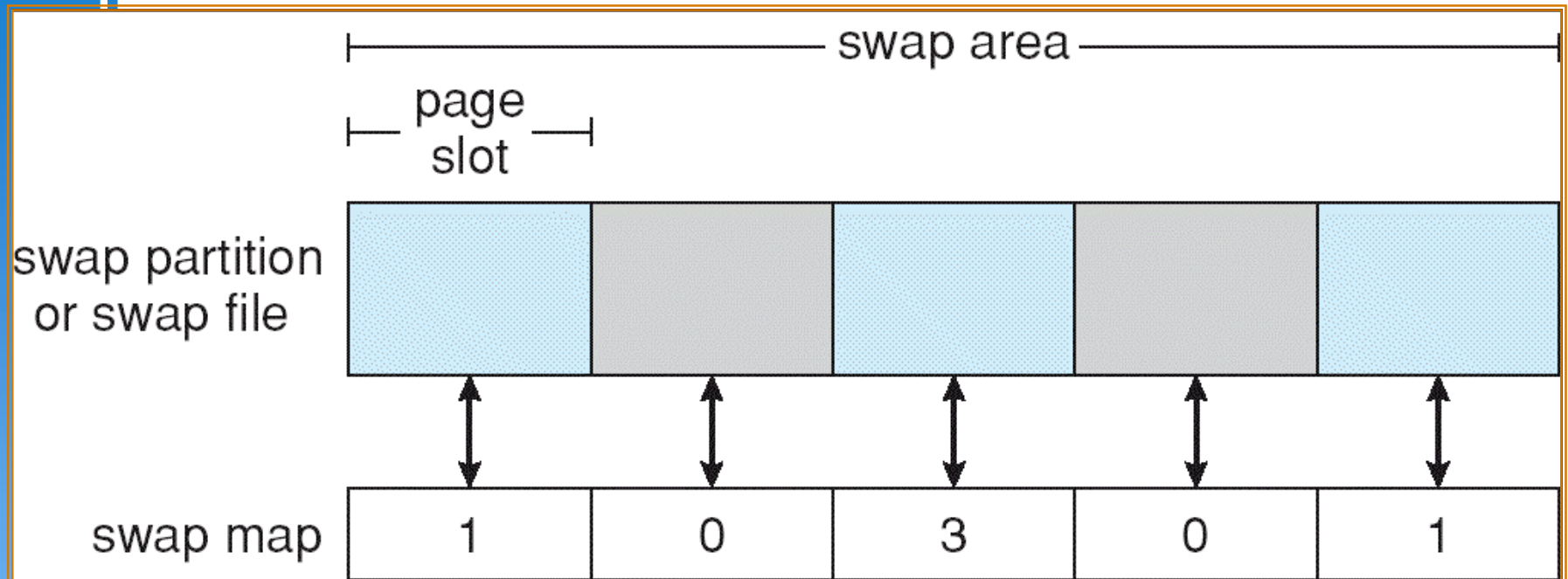


Swap-Space Management

- ▶ Swap-space — Virtual memory uses disk space as an extension of main memory.
- ▶ Swap-space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition.
- ▶ 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 page is forced out of physical memory, not when the virtual memory page is first created.



Data Structures for Swapping on Linux Systems



RAID Structure

- ▶ **RAID** – multiple disk drives provides **reliability** via **redundancy**.
- ▶ RAID is arranged into six different levels.

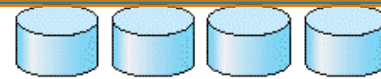


RAID (cont)

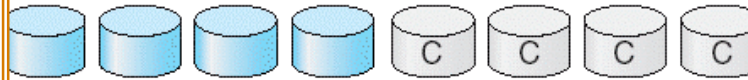
- ▶ Several improvements in disk-use techniques involve the use of multiple disks working cooperatively.
- ▶ Disk striping uses a group of disks as one storage unit.
- ▶ RAID schemes improve performance and improve the reliability of the storage system by storing redundant data.
 - *Mirroring or shadowing* keeps duplicate of each disk.
 - *Block interleaved parity* uses much less redundancy.



RAID Levels



(a) RAID 0: non-redundant striping.



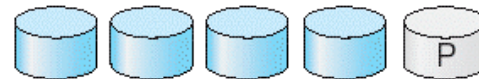
(b) RAID 1: mirrored disks.



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



(d) RAID 3: bit-interleaved parity.



(e) RAID 4: block-interleaved parity.



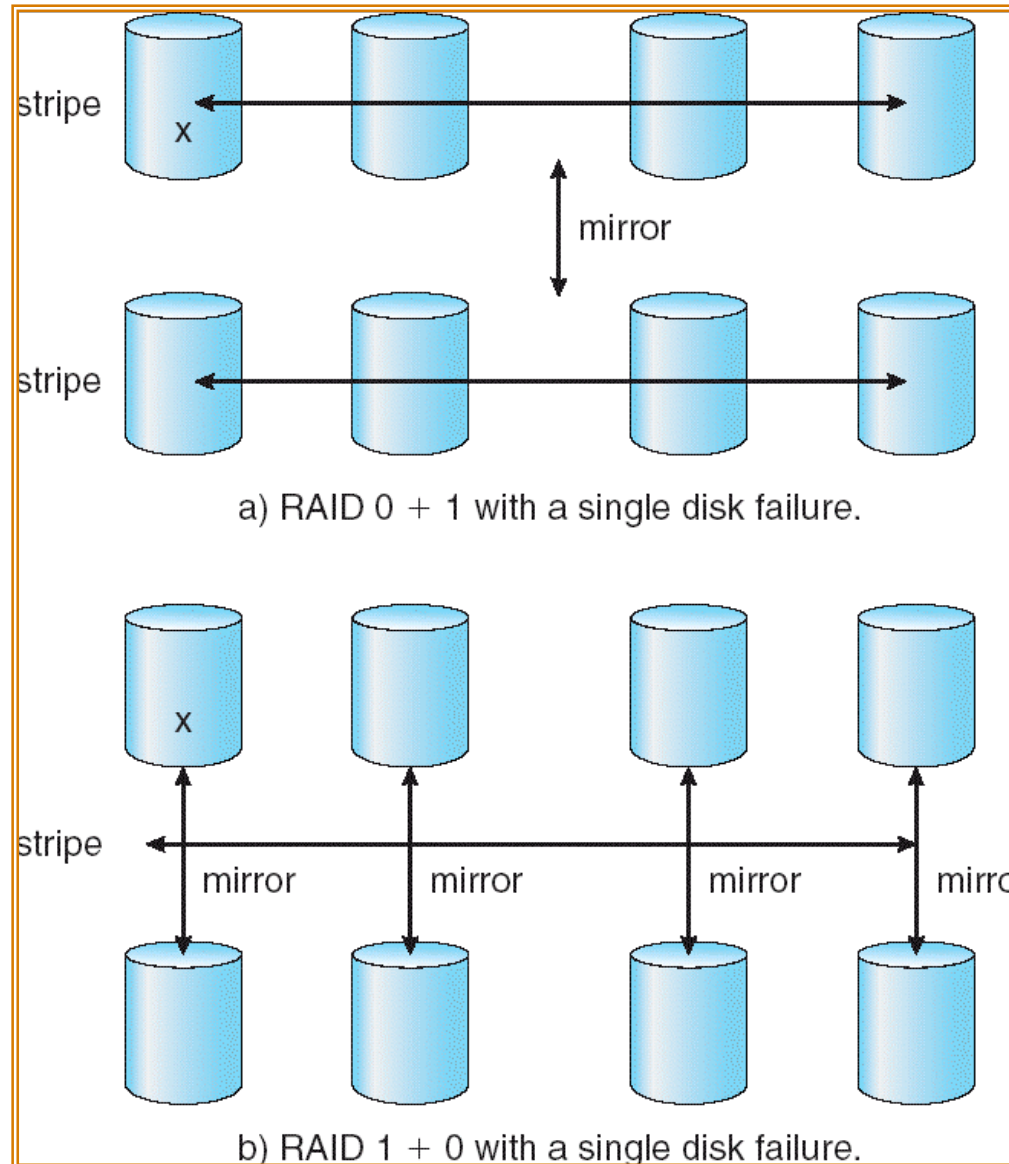
(f) RAID 5: block-interleaved distributed parity.



(g) RAID 6: P + Q redundancy.



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



Stable-Storage Implementation

- ▶ Write-ahead log scheme requires stable storage.
- ▶ To implement stable storage:
 - Replicate information on more than one nonvolatile storage media with independent failure modes.
 - Update information in a controlled manner to ensure that we can recover the stable data after any failure during data transfer or recovery.

