

*Operating  
Systems:  
Internals  
and  
Design  
Principles*

# Chapter 11 I/O Management and Disk Scheduling

Ninth Edition  
By William Stallings

# Categories of I/O Devices

External devices that engage in I/O with computer systems can be roughly grouped into three categories:

## **Human readable**

- Suitable for communicating with the computer user
- Printers, terminals, video display, keyboard, mouse

## **Machine readable**

- Suitable for communicating with electronic equipment
- Disk drives, USB keys, sensors, controllers

## **Communication**

- Suitable for communicating with remote devices
- Modems, digital line drivers



# Differences in I/O Devices

- Devices differ in a number of areas:

## *Data Rate*

- There may be differences of magnitude between the data transfer rates

## *Application*

- The use to which a device is put has an influence on the software

## *Complexity of Control*

- The effect on the operating system is filtered by the complexity of the I/O module that controls the device

## *Unit of Transfer*

- Data may be transferred as a stream of bytes or characters or in larger blocks

## *Data Representation*

- Different data encoding schemes are used by different devices

## *Error Conditions*

- The nature of errors, the way in which they are reported, their consequences, and the available range of responses differs from one device to another

**Gigabit Ethernet**

**Graphics display**

**Hard disk**

**Ethernet**

**Optical disk**

**Scanner**

**Laser printer**

**Floppy disk**

**Modem**

**Mouse**

**Keyboard**

$10^1$

$10^2$

$10^3$

$10^4$

$10^5$

$10^6$

$10^7$

$10^8$

$10^9$

Data Rate (bps)

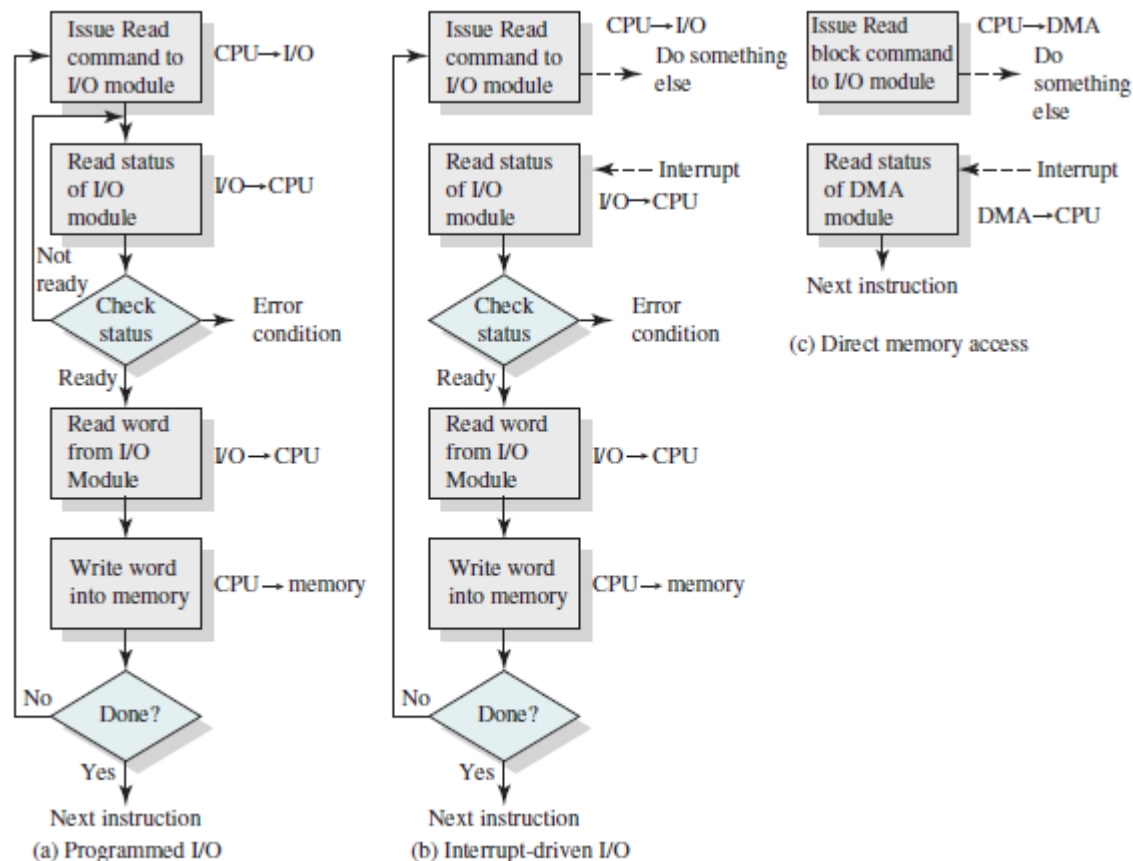
**Figure 11.1 Typical I/O Device Data Rates**

# Organization of the I/O Function

- Three techniques for performing I/O are:
- **Programmed I/O**
  - The processor issues an I/O command on behalf of a process to an I/O module; that process then busy waits for the operation to be completed before proceeding
- **Interrupt-driven I/O**
  - The processor issues an I/O command on behalf of a process
    - If non-blocking – processor continues to execute instructions from the process that issued the I/O command
    - If blocking – the next instruction the processor executes is from the OS, which will put the current process in a blocked state and schedule another process
- **Direct Memory Access (DMA)**
  - A DMA module controls the exchange of data between main memory and an I/O module



# Organization of the I/O Function



**Figure C.1** Three Techniques for Input of a Block of Data

# Evolution of the I/O Function

1

- Processor directly controls a peripheral device

2

- A controller or I/O module is added

3

- Same configuration as step 2, but now interrupts are employed

4

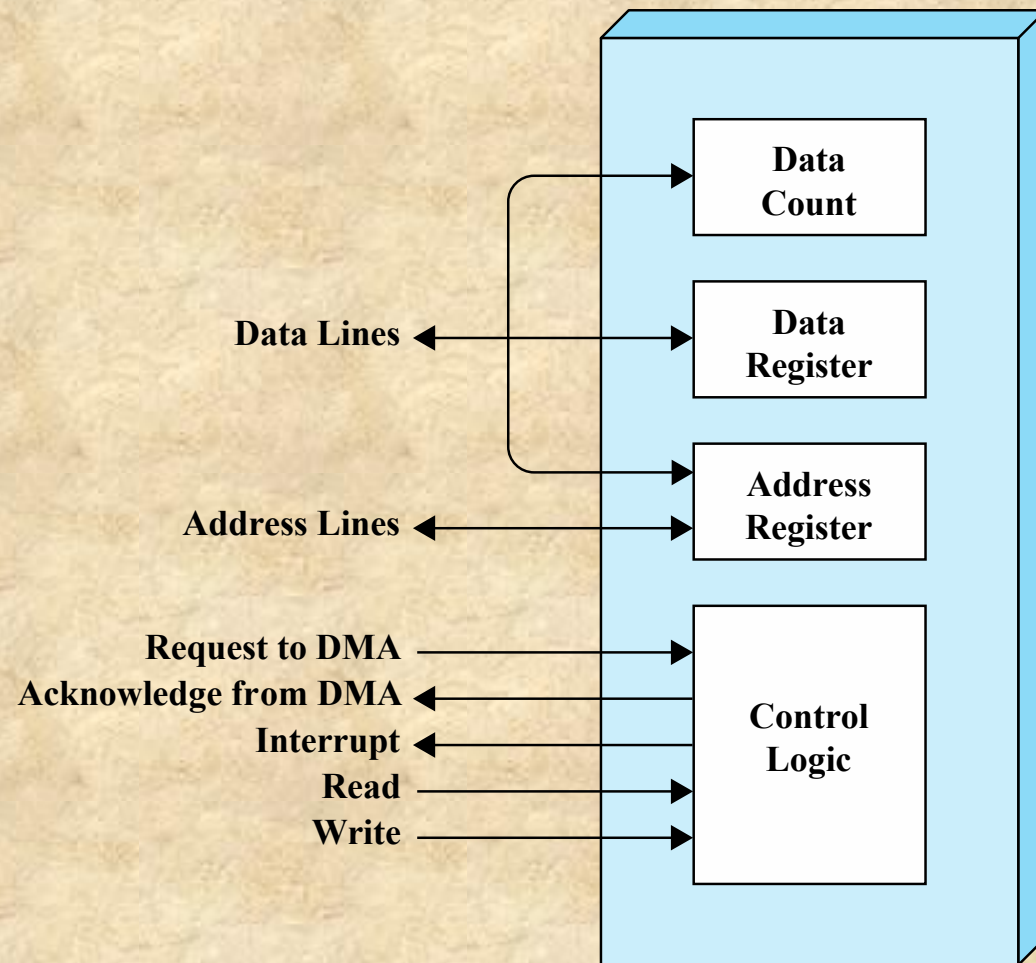
- The I/O module is given direct control of memory via DMA

5

- The I/O module is enhanced to become a separate processor, with a specialized instruction set tailored for I/O

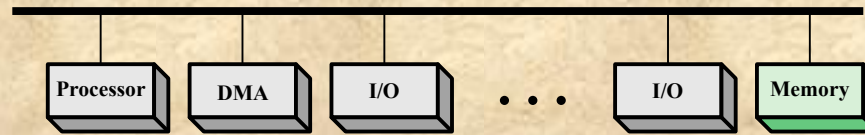
6

- The I/O module has a local memory of its own and is, in fact, a computer in its own right

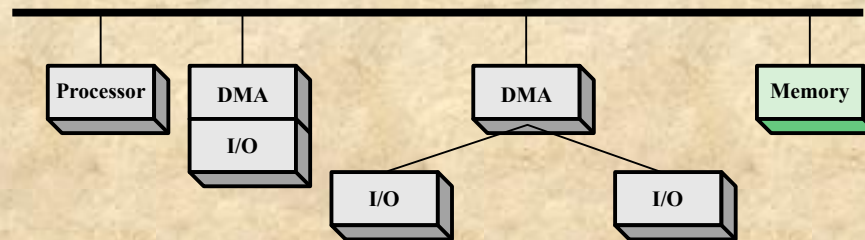


**Figure 11.2 Typical DMA Block Diagram**

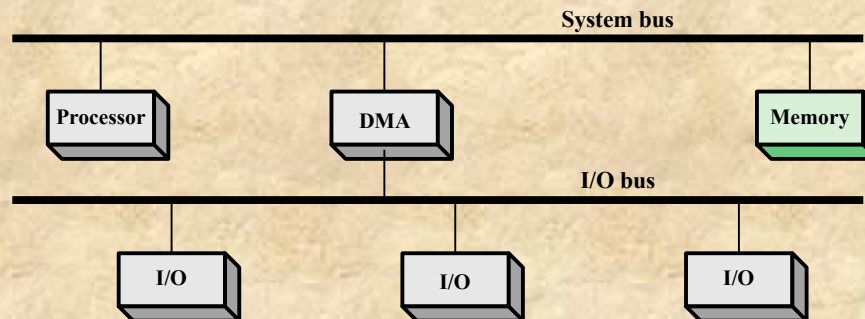




(a) Single-bus, detached DMA



(b) Single-bus, Integrated DMA-I/O



(c) I/O bus

**Figure 11.3 Alternative DMA Configurations**

# Table 11.1

## I/O Techniques

	No Interrupts	Use of Interrupts
<b>I/O-to-memory transfer through processor</b>	Programmed I/O	Interrupt-driven I/O
<b>Direct I/O-to-memory transfer</b>		Direct memory access (DMA)

# Design Objectives

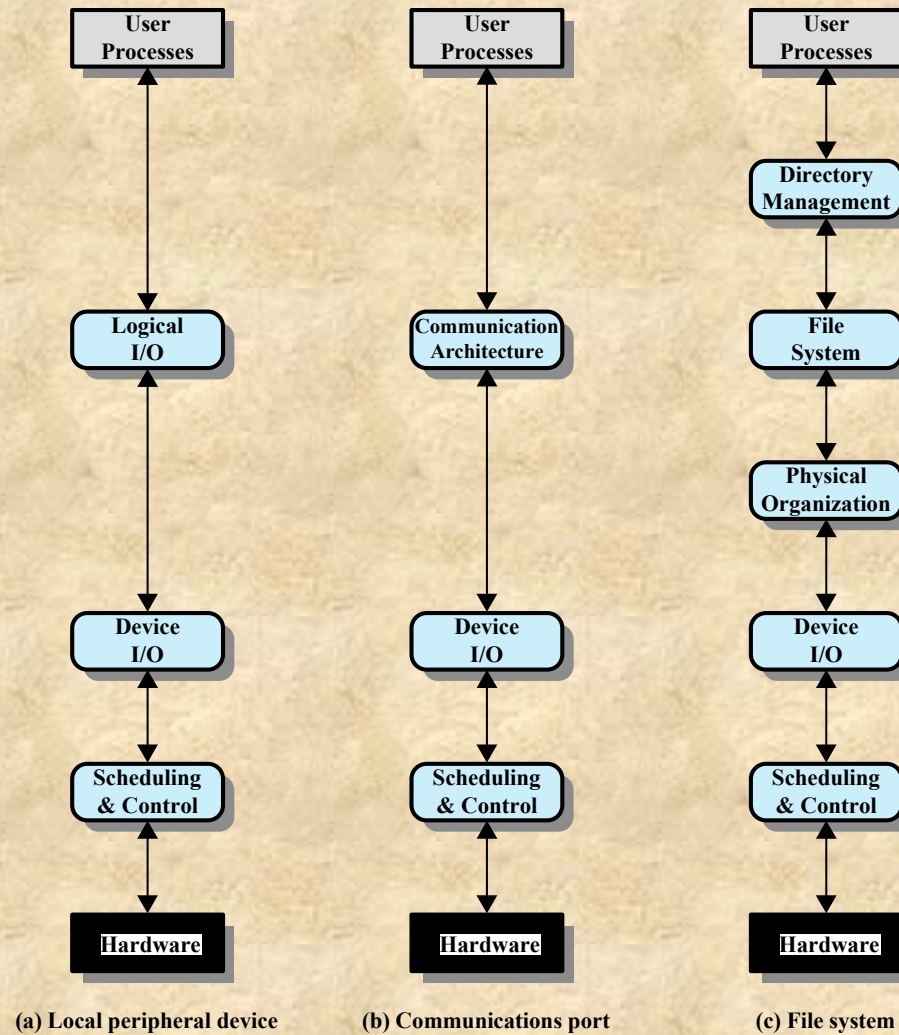
## Efficiency

- Major effort in I/O design
- Important because I/O operations often form a bottleneck
- Most I/O devices are extremely slow compared with main memory and the processor
- The area that has received the most attention is disk I/O

## Generality

- Desirable to handle all devices in a uniform manner
- Applies to the way processes view I/O devices and the way the operating system manages I/O devices and operations
- Diversity of devices makes it difficult to achieve true generality
- Use a hierarchical, modular approach to the design of the I/O function





**Figure 11.4 A Model of I/O Organization**

# Buffering

- To avoid overheads and inefficiencies, it is sometimes convenient to perform input transfers in advance of requests being made, and to perform output transfers some time after the request is made

## Block-oriented device

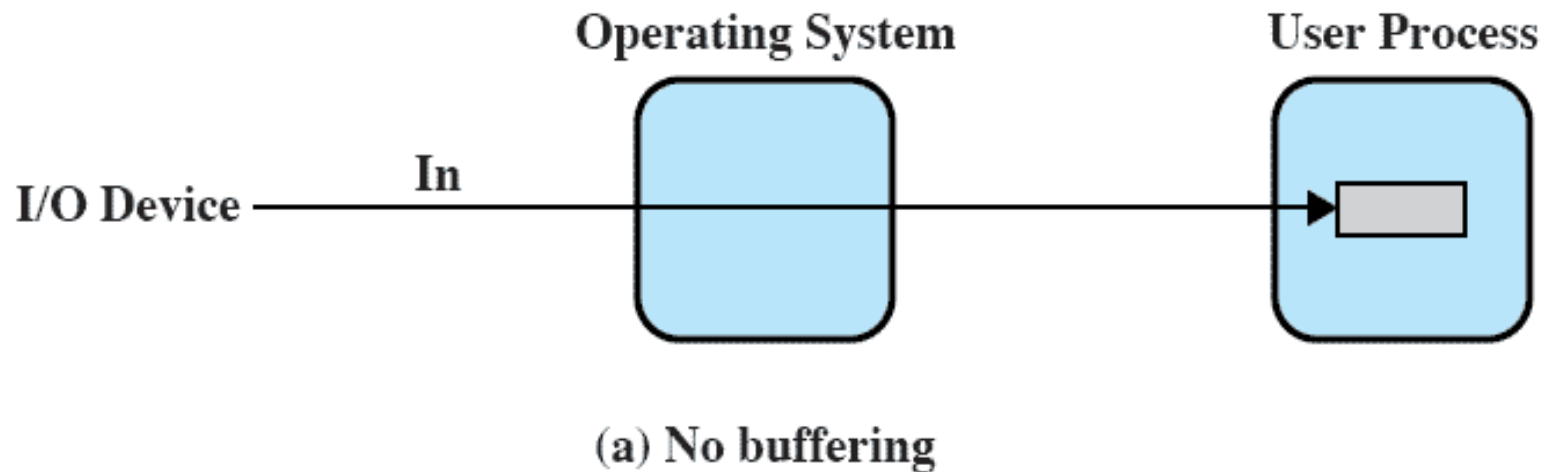
- Stores information in blocks that are usually of fixed size
- Transfers are made one block at a time
- Possible to reference data by its block number
- Disks and USB keys are examples

## Stream-oriented device

- Transfers data in and out as a stream of bytes
- No block structure
- Terminals, printers, communications ports, mouse and other pointing devices, and most other devices that are not secondary storage are examples

# No Buffer

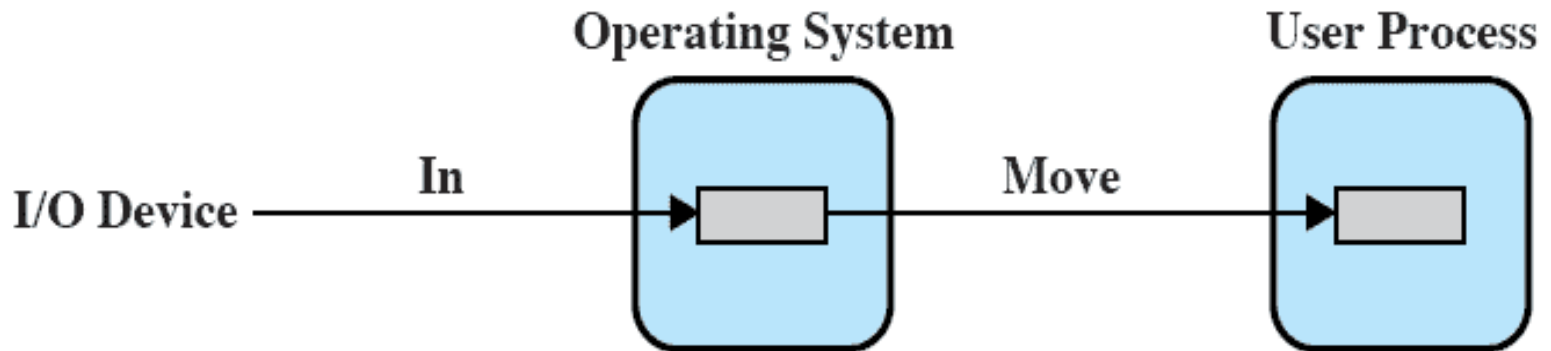
- Without a buffer, the OS directly accesses the device when it needs





# Single Buffer

- The simplest type of support that the operating system can provide
- When a user process issues an I/O request, the OS assigns a buffer in the system portion of main memory to the operation



(b) Single buffering

# Single Buffer

- Input transfers are made to the system buffer
- Reading ahead/anticipated input
  - Is done in the expectation that the block will eventually be needed
  - When the transfer is complete, the process moves the block into user space and immediately requests another block
- Approach generally provides a speedup compared to the lack of system buffering
  - The user process can be processing one block of data while the next block is being read in
  - The OS is able to swap the process out because the input operation is taking place in system memory rather than user process memory
- Disadvantages:
  - Complicates the logic in the operating system
  - Swapping logic is also affected

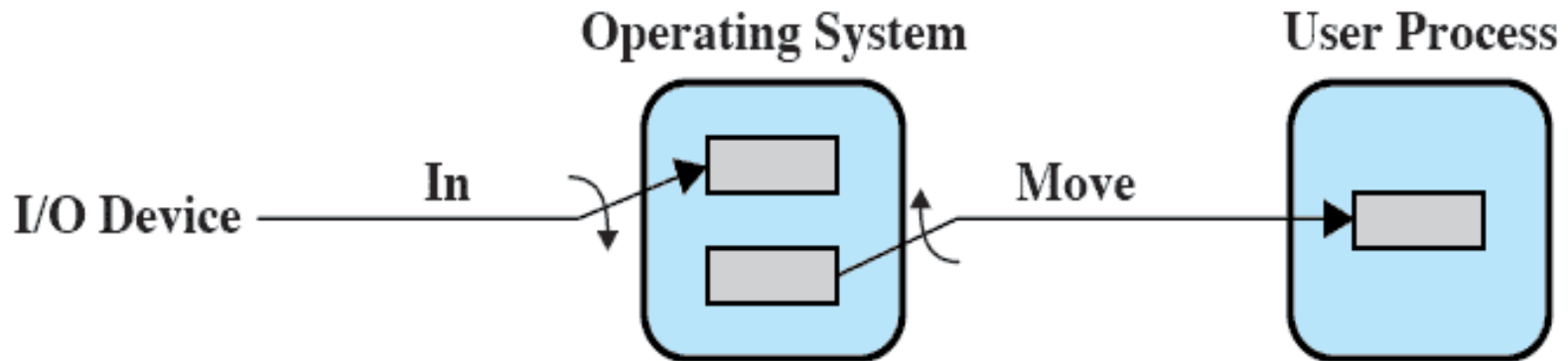
# Single Buffering for Stream-Oriented I/O

- Can be used in a line-at-a-time fashion or a byte-at-a-time fashion
- Line-at-a-time operation is appropriate for scroll-mode terminals (dumb terminals)
  - With this form of terminal, user input is one line at a time, with a carriage return signaling the end of a line
  - Output to the terminal is similarly one line at a time
- Byte-at-a-time operation is used on forms-mode terminals, when each keystroke is significant and for many other peripherals, such as sensors and controllers



# Double Buffer

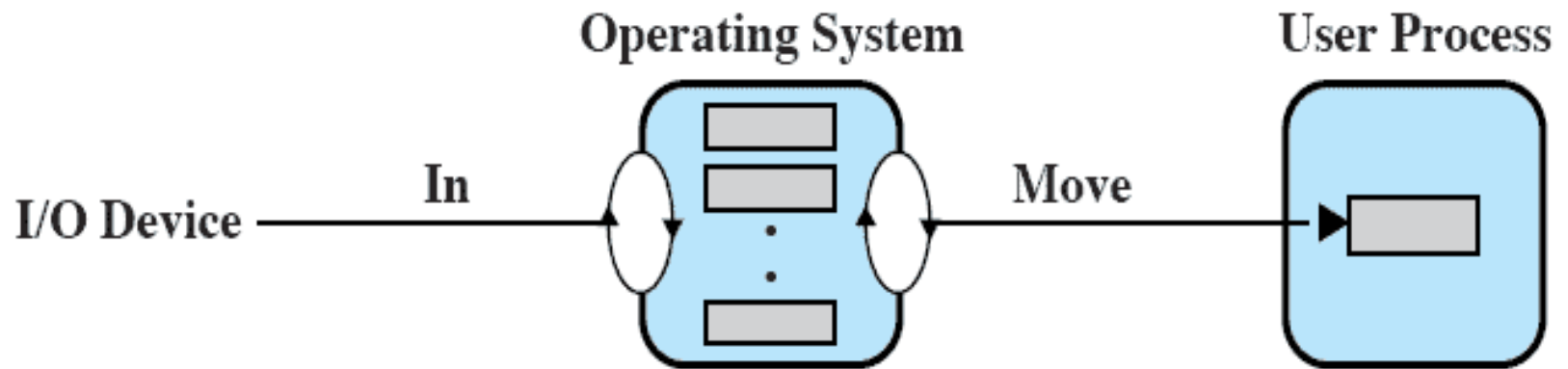
- Assigning two system buffers to the operation
- A process now transfers data to or from one buffer while the operating system empties or fills the other buffer
- Also known as buffer swapping



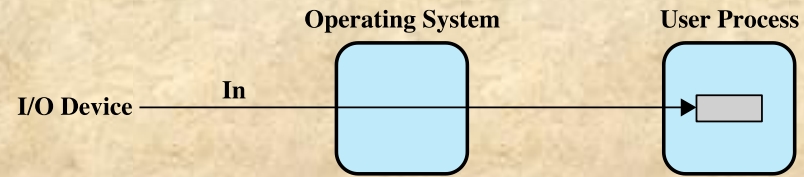
(c) Double buffering

# Circular Buffer

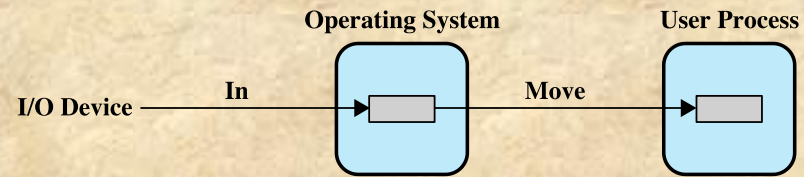
- When more than two buffers are used, the collection of buffers is itself referred to as a circular buffer
- Each individual buffer is one unit in the circular buffer



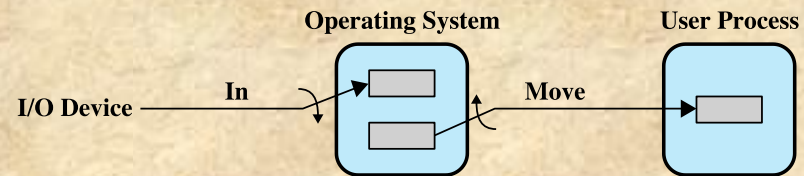
(d) Circular buffering



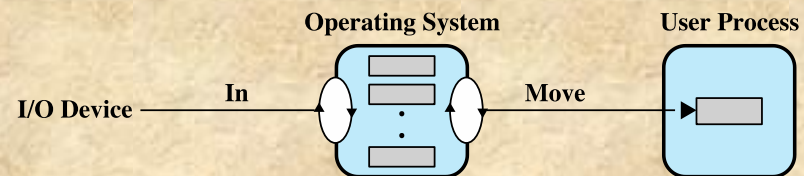
(a) No buffering



(b) Single buffering



(c) Double buffering



(d) Circular buffering

**Figure 11.5 I/O Buffering Schemes (input)**

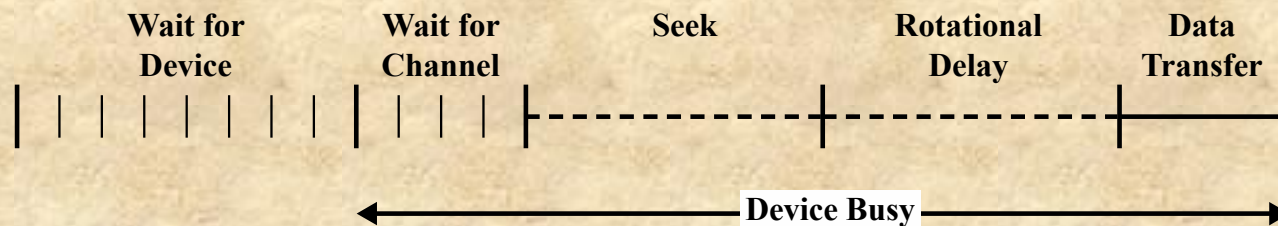


# Disk Scheduling

- Over the last 40 years, the increase in the speed of processors and main memory has far outpaced that for disk access.
- The result is disks are currently at least four orders of magnitude slower than main memory.
- The performance of disk storage subsystem is of vital concern, and much research has gone into schemes for improving that performance.

# Disk Performance Parameters

- The actual details of disk I/O operation depend on the:
  - Computer system
  - Operating system
  - Nature of the I/O channel and disk controller hardware



**Figure 11.6 Timing of a Disk I/O Transfer**

# Disk Performance Parameters

- When the disk drive is operating, the disk is rotating at constant speed
- To read or write the head must be positioned at the desired track and at the beginning of the desired sector on that track
- Track selection involves moving the head in a movable-head system or electronically selecting one head on a fixed-head system
- On a movable-head system the time it takes to position the head at the track is known as **seek time**
- The time it takes for the beginning of the sector to reach the head is known as **rotational delay**
- The sum of the seek time and the rotational delay equals the **access time**



# Seek Time

- The time required to move the disk arm to the required track
- Consists of two key components:
  - The initial startup time
  - The time taken to traverse the tracks that have to be crossed once the access arm is up to speed
- Settling time
  - Time after positioning the head over the target track until track identification is confirmed
- Much improvement comes from smaller and lighter disk components
- A typical average seek time on contemporary hard disks is under 10ms

# Disk Performance

## ■ Rotational delay

- The time required for the addressed area of the disk to rotate into a position where it is accessible by the read/write head
- Disks rotate at speeds ranging from 3,600 rpm (for handheld devices such as digital cameras) up to 15,000 rpm

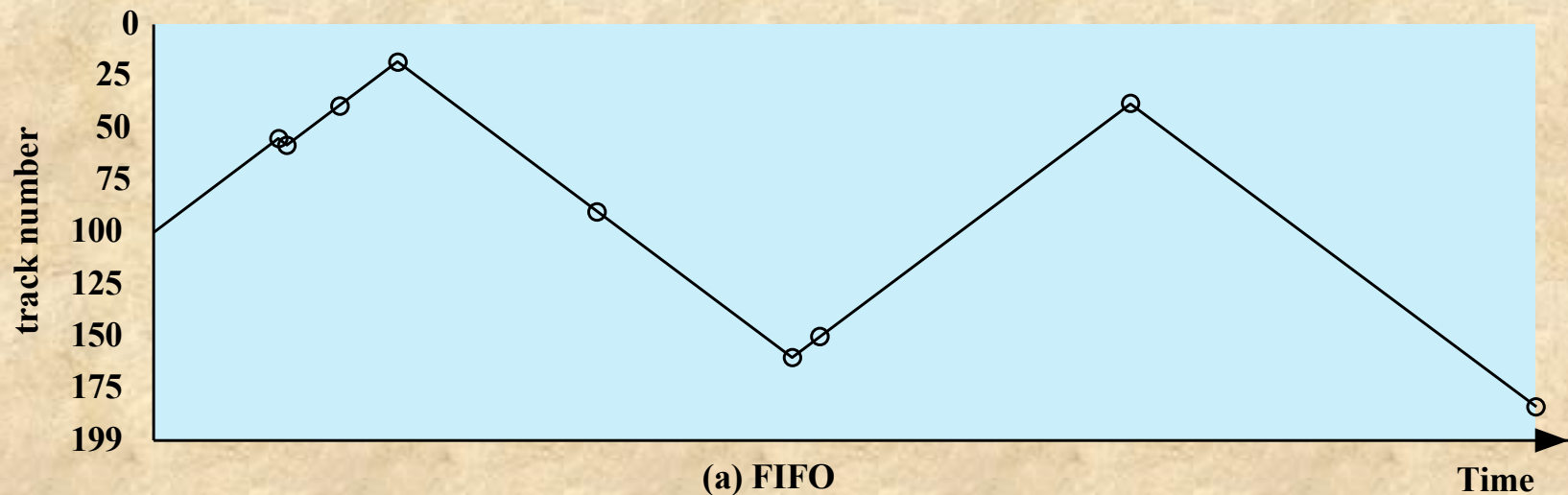
# Disk scheduling policies

- Let's try to do some calculations.
- $T_s = 4ms, r = 7500rpm, 512 \text{ byte sectors}$
- *File size of 1.28MB with 2500 sectors*
  - 1. Compacted
  - 2. Randomly put



# First-In, First-Out (FIFO)

- Processes in sequential order
- Fair to all processes
- Approximates random scheduling in performance if there are many processes competing for the disk

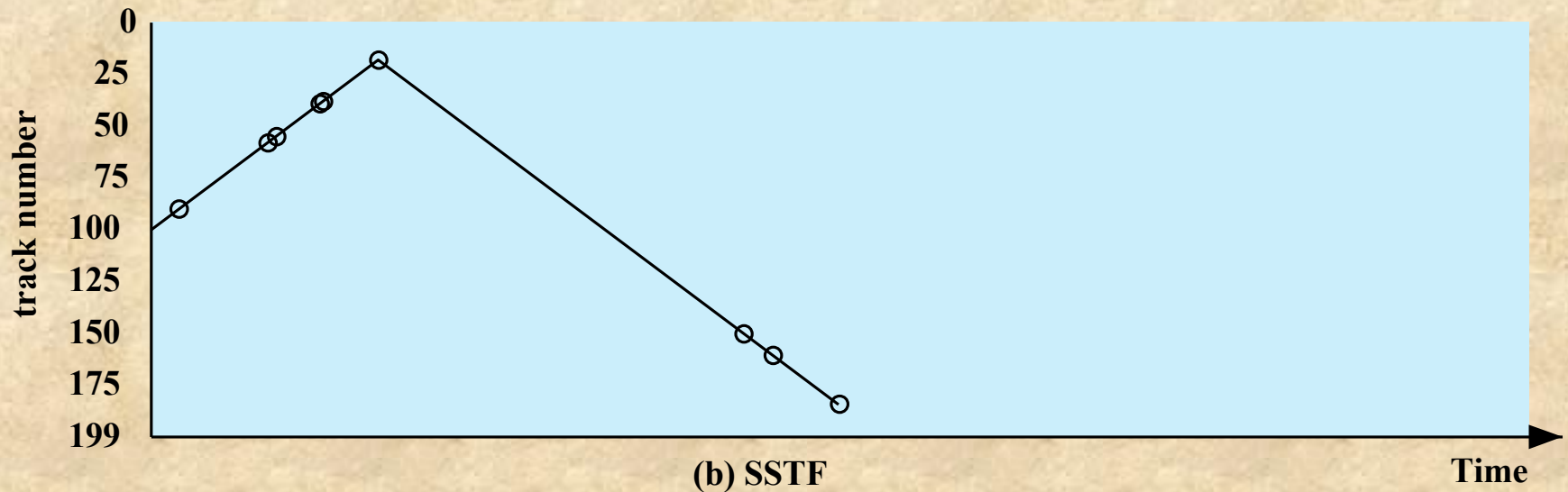


# Priority (PRI)

- Control of the scheduling is outside the control of disk management software
- Goal is not to optimize disk utilization but to meet other objectives
- Short batch jobs and interactive jobs are given higher priority
- Provides good interactive response time
- Longer jobs may have to wait an excessively long time
- A poor policy for database systems

# Shortest Service Time First (SSTF)

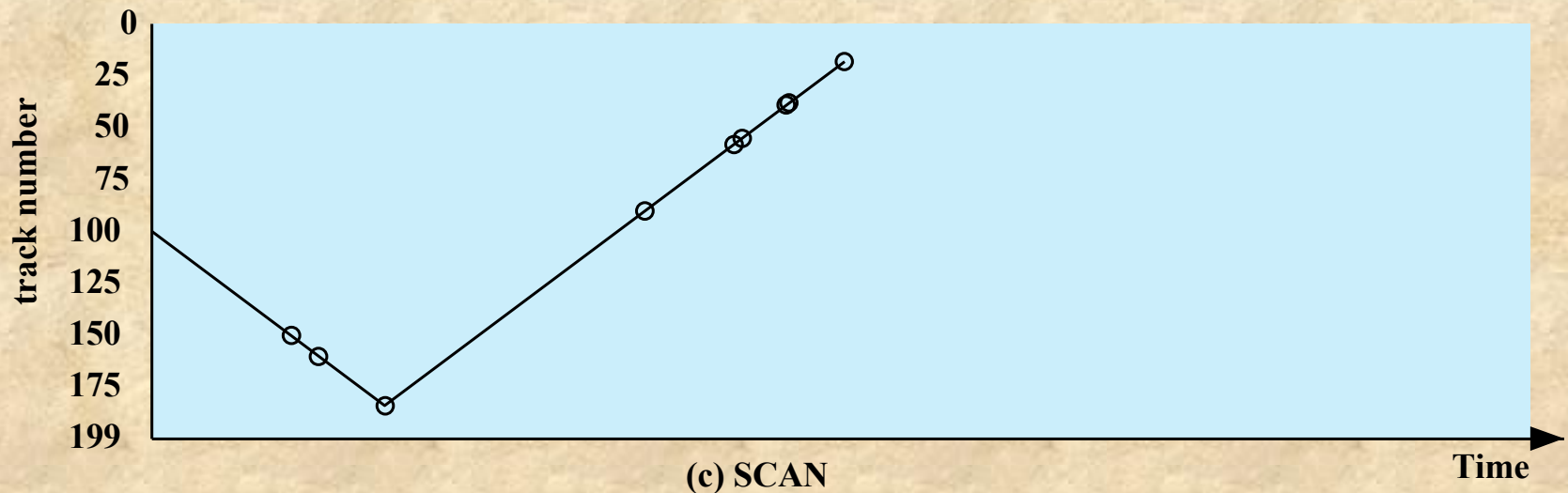
- Select the disk I/O request that requires the least movement of the disk arm from its current position
- Always choose the minimum seek time





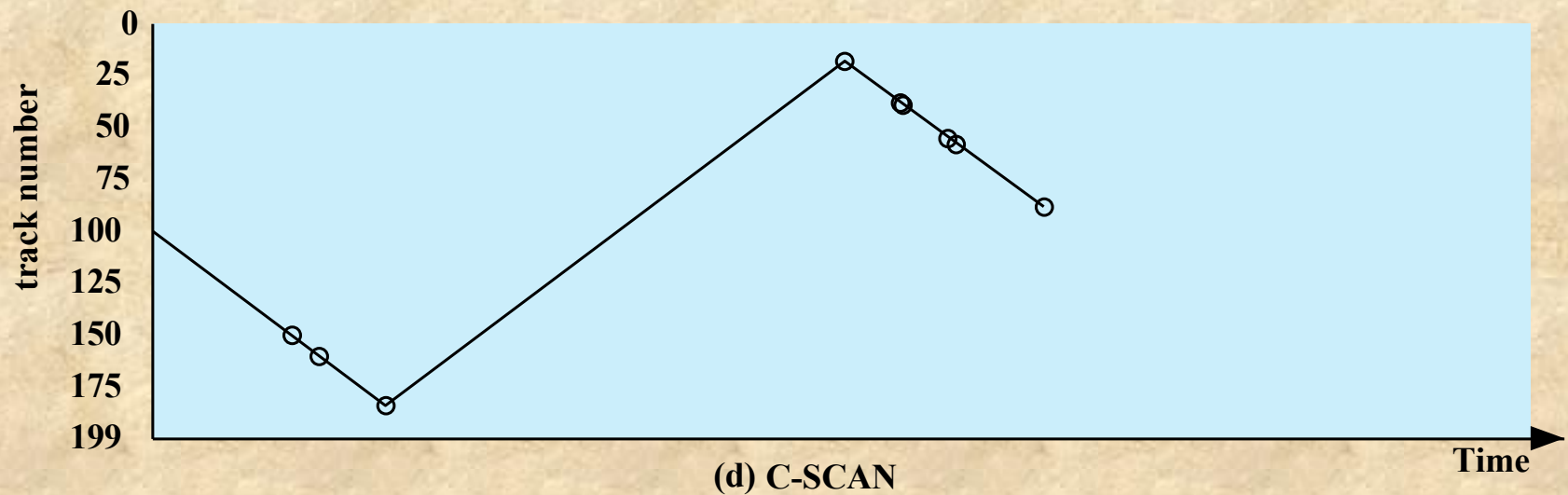
# SCAN

- Also known as the elevator algorithm
- Arm moves in one direction only
  - Satisfies all outstanding requests until it reaches the last track in that direction then the direction is reversed
- Favors jobs whose requests are for tracks nearest to both innermost and outermost tracks and favors the latest-arriving jobs



# C-SCAN (Circular SCAN)

- Restricts scanning to one direction only
- When the last track has been visited in one direction, the arm is returned to the opposite end of the disk and the scan begins again



(a) FIFO (starting at track 100)		(b) SSTF (starting at track 100)		(c) SCAN (starting at track 100, in the direction of increasing track number)		(d) C-SCAN (starting at track 100, in the direction of increasing track number)	
Next track accessed	Number of tracks traversed	Next track accessed	Number of tracks traversed	Next track accessed	Number of tracks traversed	Next track accessed	Number of tracks traversed
55	45	90	10	150	50	150	50
58	3	58	32	160	10	160	10
39	19	55	3	184	24	184	24
18	21	39	16	90	94	18	166
90	72	38	1	58	32	38	20
160	70	18	20	55	3	39	1
150	10	150	132	39	16	55	16
38	112	160	10	38	1	58	3
184	146	184	24	18	20	90	32
<b>Average seek length</b>	55.3	<b>Average seek length</b>	27.5	<b>Average seek length</b>	27.8	<b>Average seek length</b>	35.8

**Table 11.2 Comparison of Disk Scheduling Algorithms**



# N-Step-SCAN

- Segments the disk request queue into subqueues of length  $N$
- Subqueues are processed one at a time, using SCAN
- While a queue is being processed new requests must be added to some other queue
- If fewer than  $N$  requests are available at the end of a scan, all of them are processed with the next scan

# FSCAN

- Uses two subqueues
- When a scan begins, all of the requests are in one of the queues, with the other empty
- During scan, all new requests are put into the other queue
- Service of new requests is deferred until all of the old requests have been processed

Name	Description	Remarks
<b>Selection according to requestor</b>		
Random	Random scheduling	For analysis and simulation
FIFO	First in first out	Fairest of them all
PRI	Priority by process	Control outside of disk queue management
LIFO	Last in first out	Maximize locality and resource utilization
<b>Selection according to requested item</b>		
SSTF	Shortest service time first	High utilization, small queues
SCAN	Back and forth over disk	Better service distribution
C-SCAN	One way with fast return	Lower service variability
N-step-SCAN	SCAN of $N$ records at a time	Service guarantee
FSCAN	N-step-SCAN with $N$ = queue size at beginning of SCAN cycle	Load sensitive

**Table 11.3 Disk Scheduling Algorithms**



# RAID

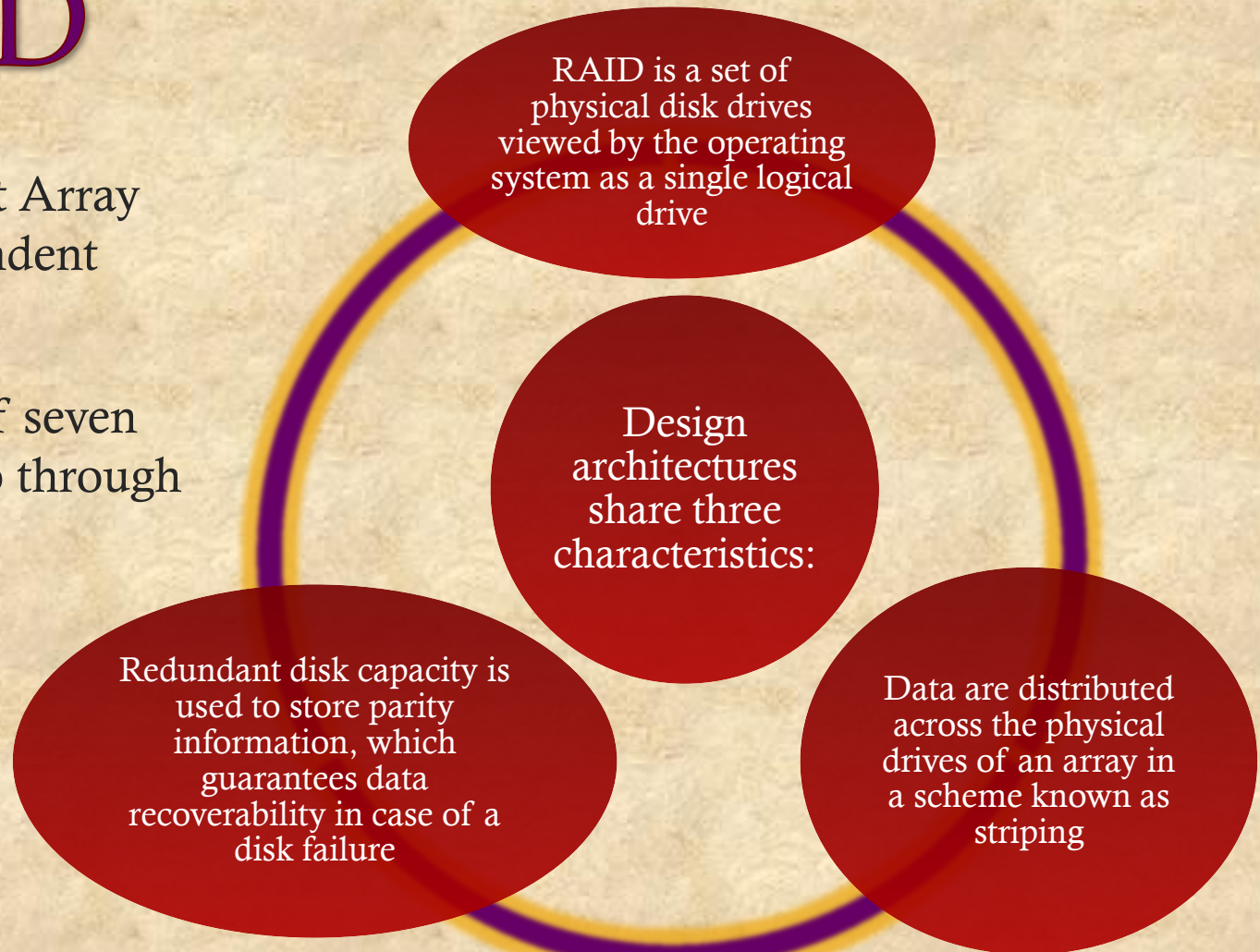
- Rate of improvement of secondary storage device performance has been staggeringly low compared to the speed of main memory and processor.
- Disk storage system has been the focus of concern in improving the overall computer system.
- Disk storage designers recognized that if one component can only be pushed so far then they should look at gains by using multiple parallel components.
- Same idea with multi-core processors.
- In case of disk storage this leads to the development of arrays of disks that operate independently and in parallel.

# RAID

- Key points:
  - With multiple disks, separate I/O requests can be handled in parallel, as long as the data required reside on separate disks.
  - Further, a single I/O request can be executed in parallel if the block of data to be accessed is distributed across multiple disks.

# RAID

- Redundant Array of Independent Disks
- Consists of seven levels, zero through six



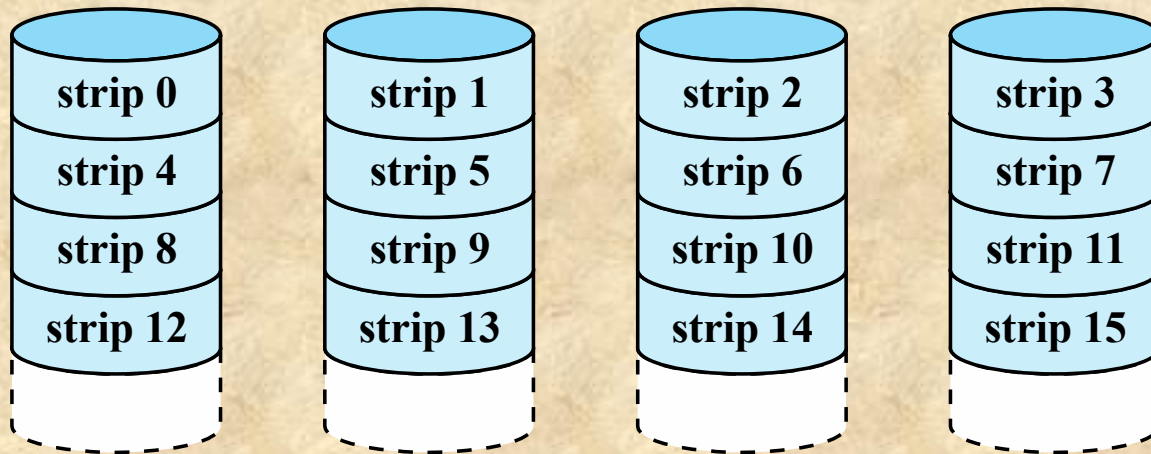


# RAID

- The term was originally coined in a paper by a group of researchers at the University of California at Berkeley
  - The paper outlined various configurations and applications and introduced the definitions of the RAID levels
- Strategy employs multiple disk drives and distributes data in such a way as to enable simultaneous access to data from multiple drives
  - Improves I/O performance and allows easier incremental increases in capacity
- The unique contribution is to address effectively the need for redundancy
- Makes use of stored parity information that enables the recovery of data lost due to a disk failure

# RAID Level 0

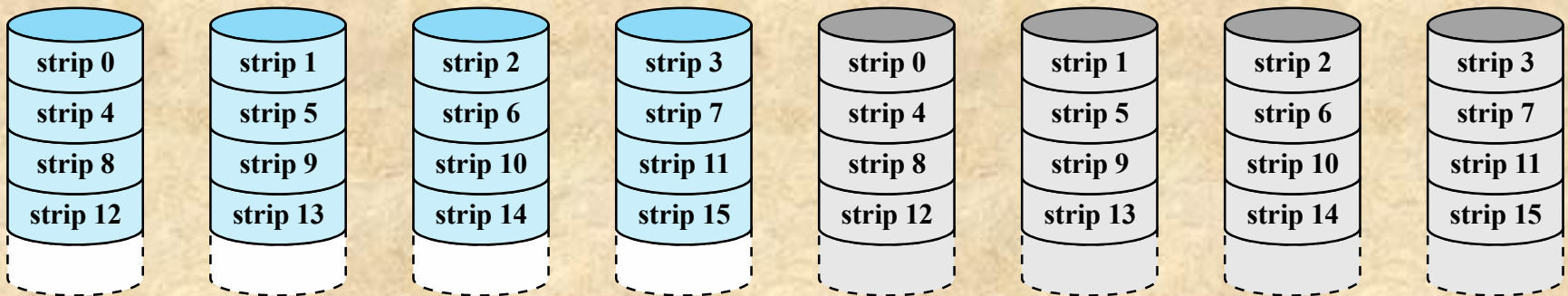
- Not a true RAID because it does not include redundancy to improve performance or provide data protection
- User and system data are distributed across all of the disks in the array
- Logical disk is divided into strips



**(a) RAID 0 (non-redundant)**

# RAID Level 1

- Redundancy is achieved by the simple expedient of duplicating all the data
- There is no “write penalty” (NO PARITY BITS)
- When a drive fails the data may still be accessed from the second drive
- Principal disadvantage is the cost

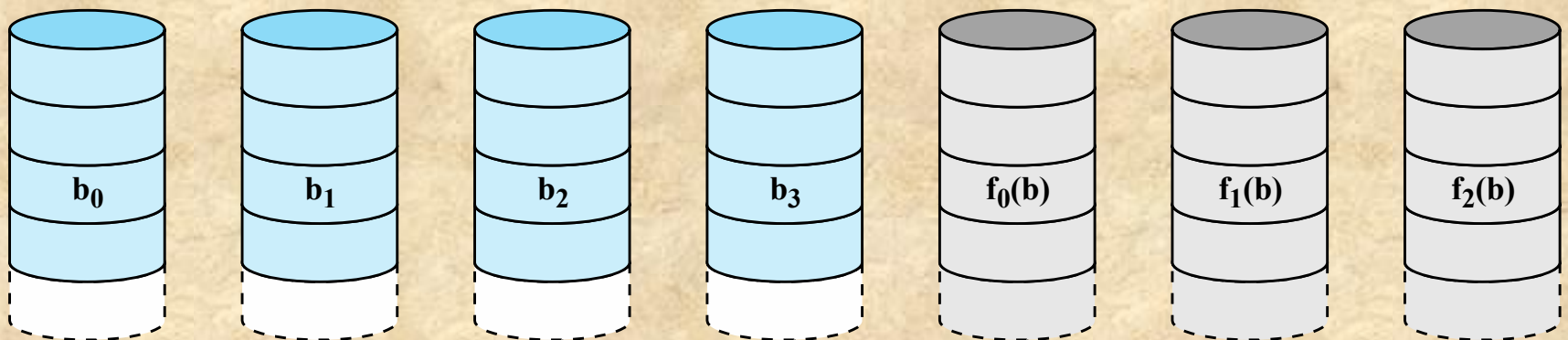


(b) RAID 1 (mirrored)



# RAID Level 2

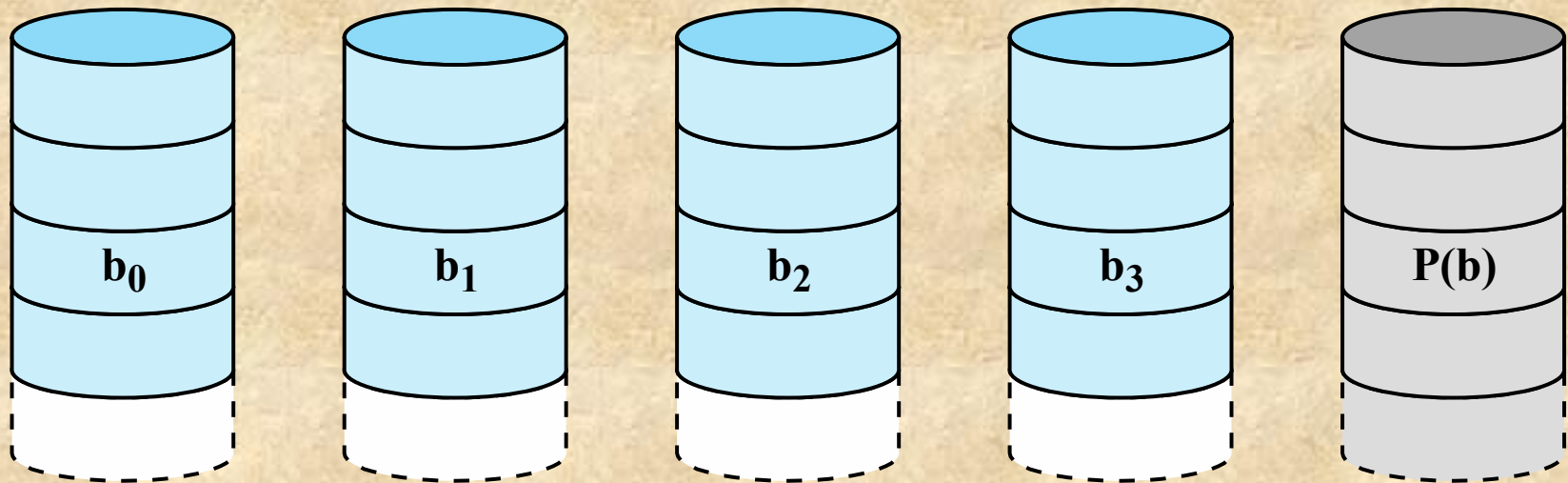
- Makes use of a parallel access technique
- Data striping is used
- Typically, a Hamming code is used
- Effective choice in an environment in which many disk errors occur
- $N+m$  ( $m$  =parity disks)



(c) RAID 2 (redundancy through Hamming code)

# RAID Level 3

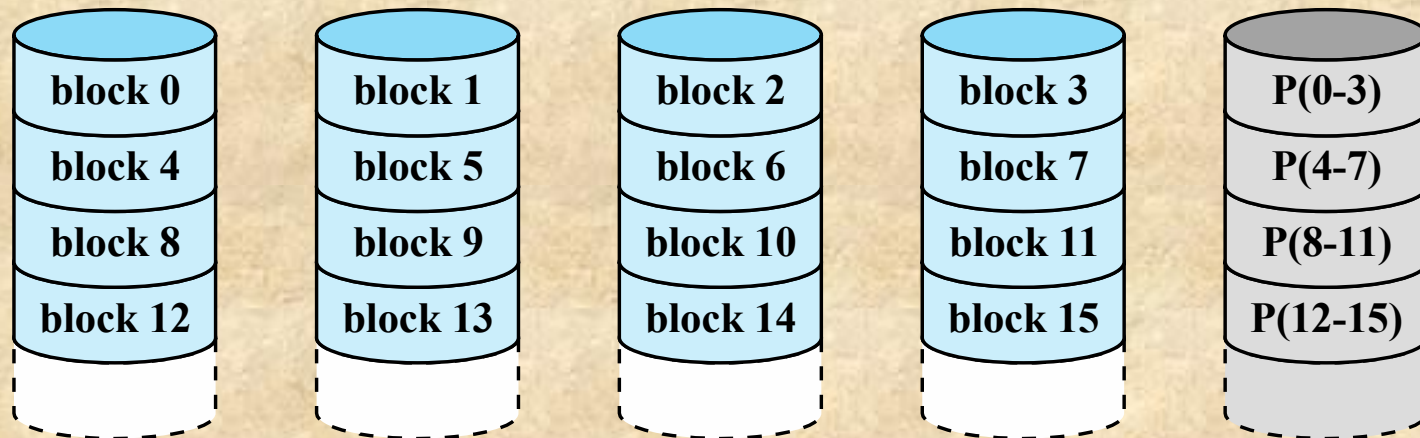
- Requires only a single redundant disk, no matter how large the disk array
- Employs parallel access, with data distributed in small strips
- Can achieve very high data transfer rates
- $N + 1$



**(d) RAID 3 (bit-interleaved parity)**

# RAID Level 4

- Makes use of an independent access technique
- A bit-by-bit parity strip is calculated across corresponding strips on each data disk, and the parity bits are stored in the corresponding strip on the parity disk
- Involves a write penalty when an I/O write request of small size is performed

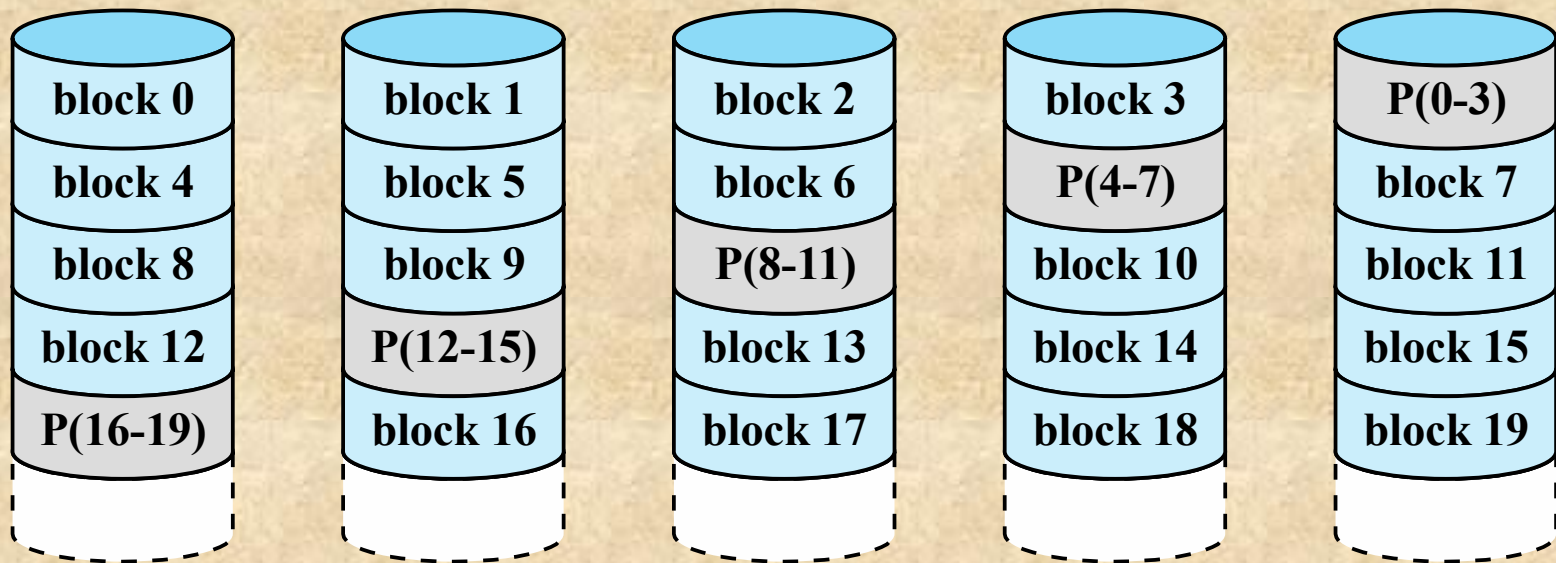


(e) RAID 4 (block-level parity)



# RAID Level 5

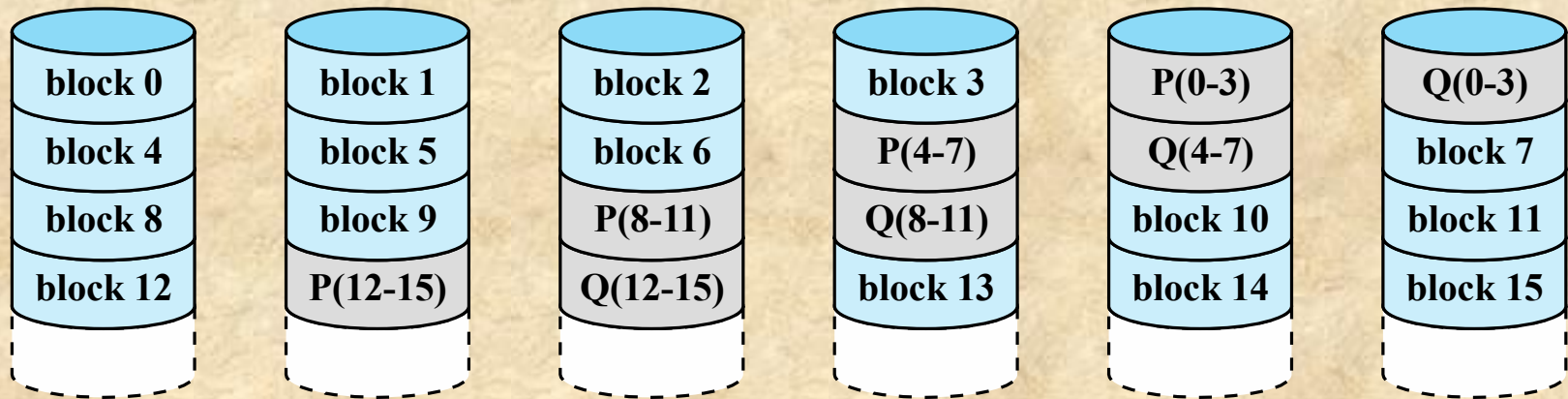
- Similar to RAID-4 but distributes the parity bits across all disks
- Typical allocation is a round-robin scheme
- Has the characteristic that the loss of any one disk does not result in data loss



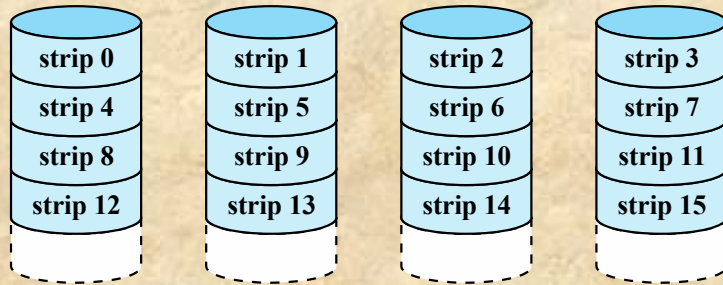
(f) RAID 5 (block-level distributed parity)

# RAID Level 6

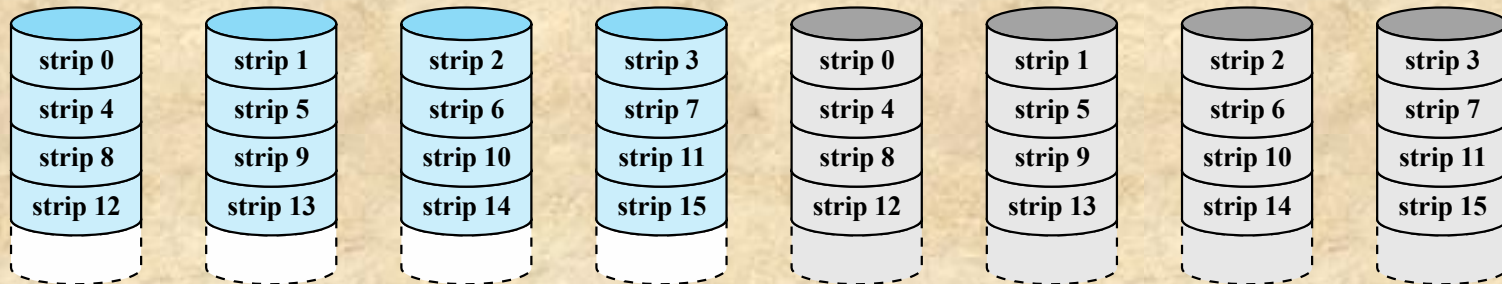
- Two different parity calculations are carried out and stored in separate blocks on different disks
- Provides extremely high data availability
- Incurs a substantial write penalty because each write affects two parity blocks
- $N + 2$



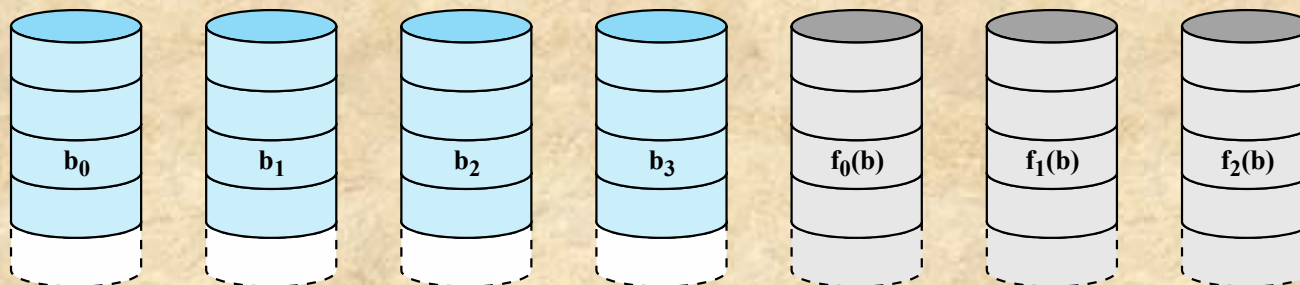
(g) RAID 6 (dual redundancy)



(a) RAID 0 (non-redundant)



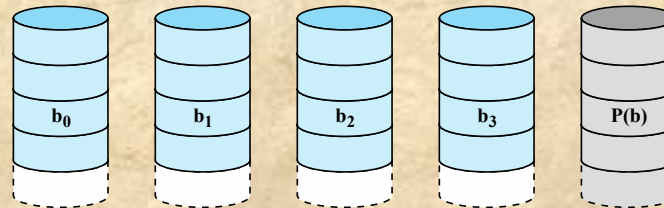
(b) RAID 1 (mirrored)



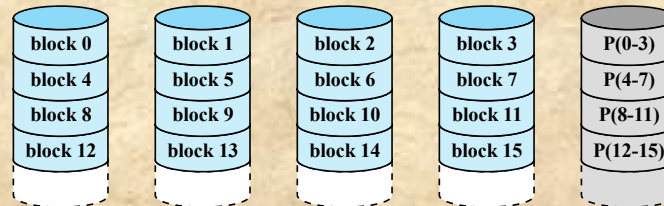
(c) RAID 2 (redundancy through Hamming code)

Figure 11.8 RAID Levels (page 1 of 2)

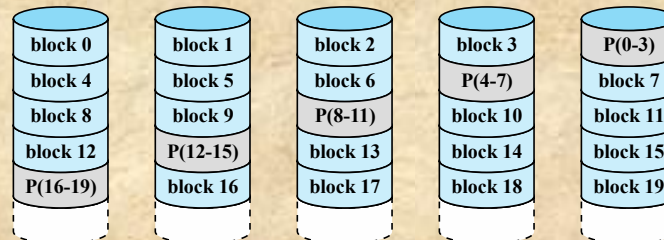




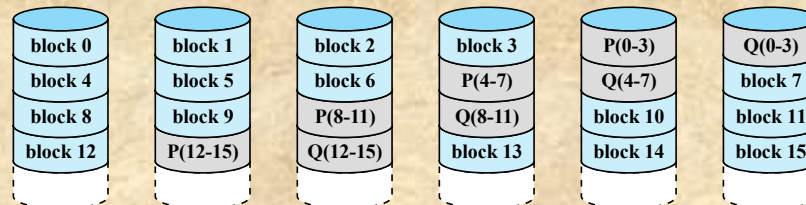
(d) RAID 3 (bit-interleaved parity)



(e) RAID 4 (block-level parity)



(f) RAID 5 (block-level distributed parity)



(g) RAID 6 (dual redundancy)

**Figure 11.8 RAID Levels** (page 2 of 2)

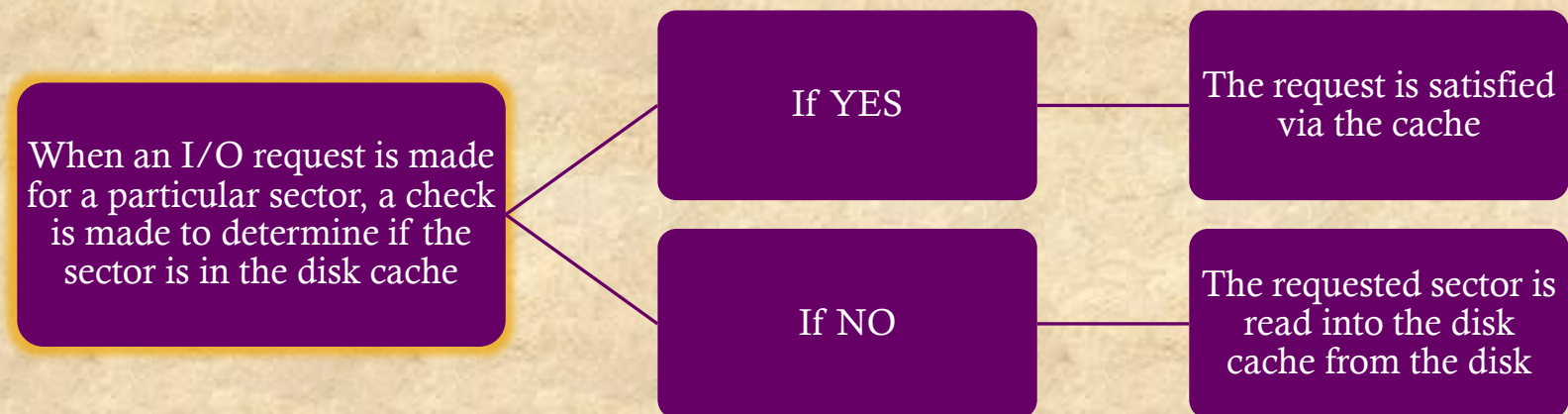
Category	Level	Description	Disks required	Data availability	Large I/O data transfer capacity	Small I/O request rate
Striping	0	Nonredundant	$N$	Lower than single disk	Very high	Very high for both read and write
Mirroring	1	Mirrored	$2N$	Higher than RAID 2, 3, 4, or 5; lower than RAID 6	Higher than single disk for read; similar to single disk for write	Up to twice that of a single disk for read; similar to single disk for write
Parallel access	2	Redundant via Hamming code	$N + m$	Much higher than single disk; comparable to RAID 3, 4, or 5	Highest of all listed alternatives	Approximately twice that of a single disk
	3	Bit-interleaved parity	$N + 1$	Much higher than single disk; comparable to RAID 2, 4, or 5	Highest of all listed alternatives	Approximately twice that of a single disk
Independent access	4	Block-interleaved parity	$N + 1$	Much higher than single disk; comparable to RAID 2, 3, or 5	Similar to RAID 0 for read; significantly lower than single disk for write	Similar to RAID 0 for read; significantly lower than single disk for write
	5	Block-interleaved distributed parity	$N + 1$	Much higher than single disk; comparable to RAID 2, 3, or 4	Similar to RAID 0 for read; lower than single disk for write	Similar to RAID 0 for read; generally lower than single disk for write
	6	Block-interleaved dual distributed parity	$N + 2$	Highest of all listed alternatives	Similar to RAID 0 for read; lower than RAID 5 for write	Similar to RAID 0 for read; significantly lower than RAID 5 for write

$N$  = number of data disks;  $m$  proportional to  $\log N$

**Table 11.4 RAID Levels** (Page 498 in textbook)

# Disk Cache

- *Cache memory* is used to apply to a memory that is smaller and faster than main memory and that is interposed between main memory and the processor
- Reduces average memory access time by exploiting the principle of locality
- *Disk cache* is a buffer **in main memory** for disk sectors
- Contains a copy of some of the sectors on the disk





# Disk Cache

## ■ Issues

- How can we deliver the data in the cache to the process?
  - Bring the data from the cache to the memory assigned to the user process
- What is the replacement strategy?
  - LRU (Least Recently Used)
  - LFU (Least Frequently Used)

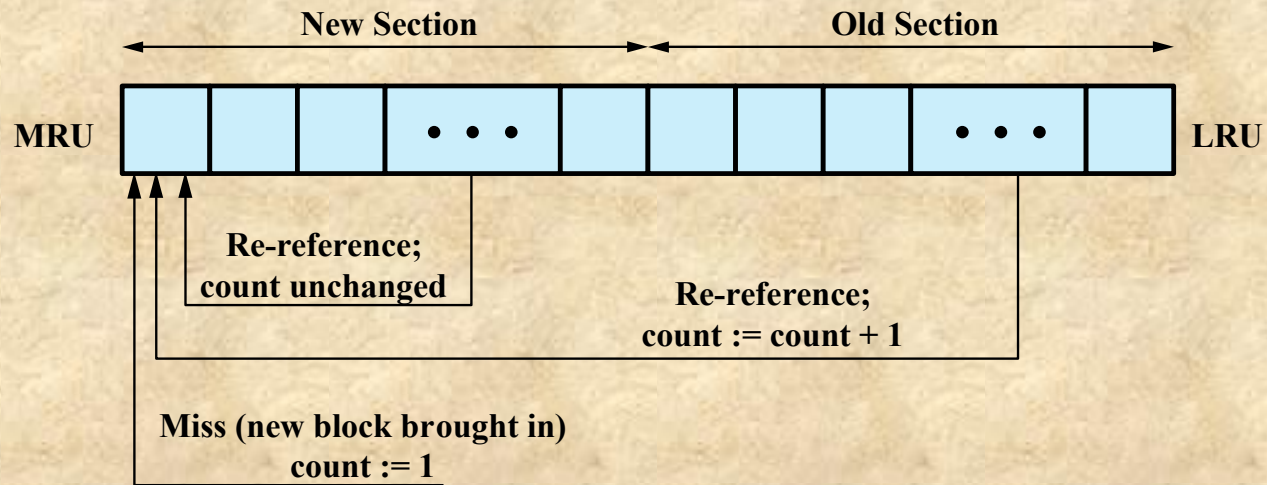
# Least Recently Used (LRU)

- Most commonly used algorithm that deals with the design issue of replacement strategy
- The block that has been in the cache the longest with no reference to it is replaced
- A stack of pointers reference the cache
  - Most recently referenced block is on the top of the stack
  - When a block is referenced or brought into the cache, it is placed on the top of the stack

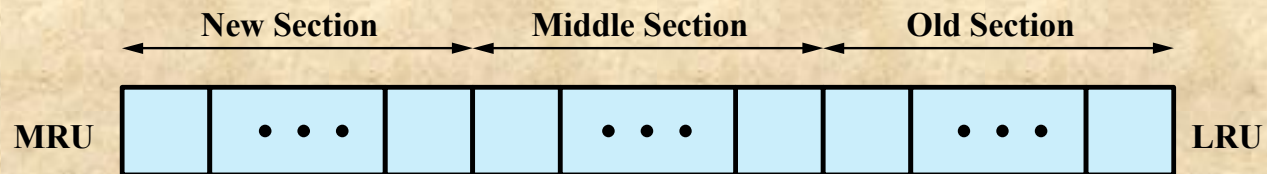
# Least Frequently Used (LFU)

- The block that has experienced the fewest references is replaced
- A counter is associated with each block
- Counter is incremented each time block is accessed
- When replacement is required, the block with the smallest count is selected





(a) FIFO



(b) Use of three sections

**Figure 11.9 Frequency-Based Replacement**