#### IT 540 Operating Systems ECE519 Advanced Operating Systems

#### Prof. Dr. Hasan Hüseyin BALIK

(11<sup>th</sup> Week)

#### (Advanced) Operating Systems

#### 11. I/O Management and Disk Scheduling

#### 11. Outline

- I/O Devices
- Organization of the I/O Function
- Operating System Design Issues
- I/O Buffering
- Disk Scheduling
- RAID
- Disk Cache

# Categories of I/O Devices

External devices that engage in I/O with computer systems can be 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 available range of responses differs from one device to another

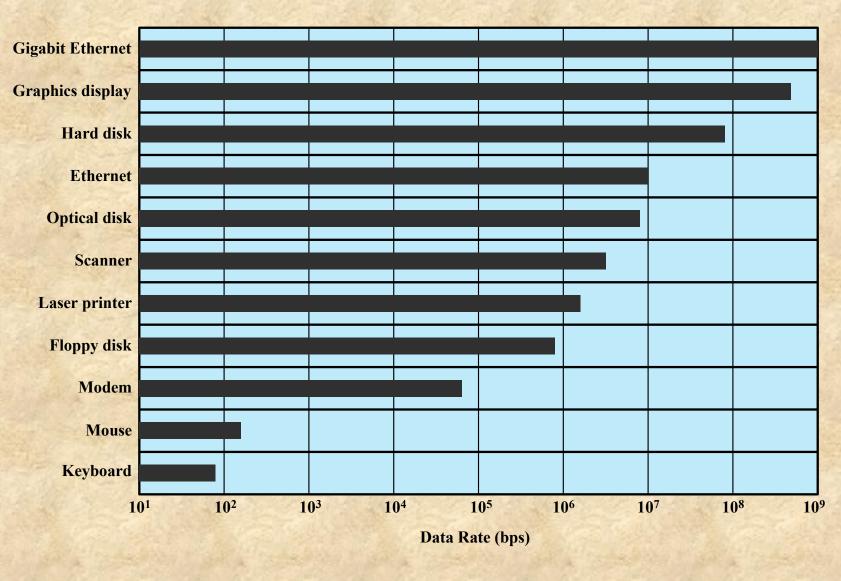


Figure 11.1 Typical I/O Device Data Rates

### 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

## I/O Techniques

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

### Evolution of the I/O Function

- Processor directly controls a peripheral device
- A controller or I/O module is added

2

3

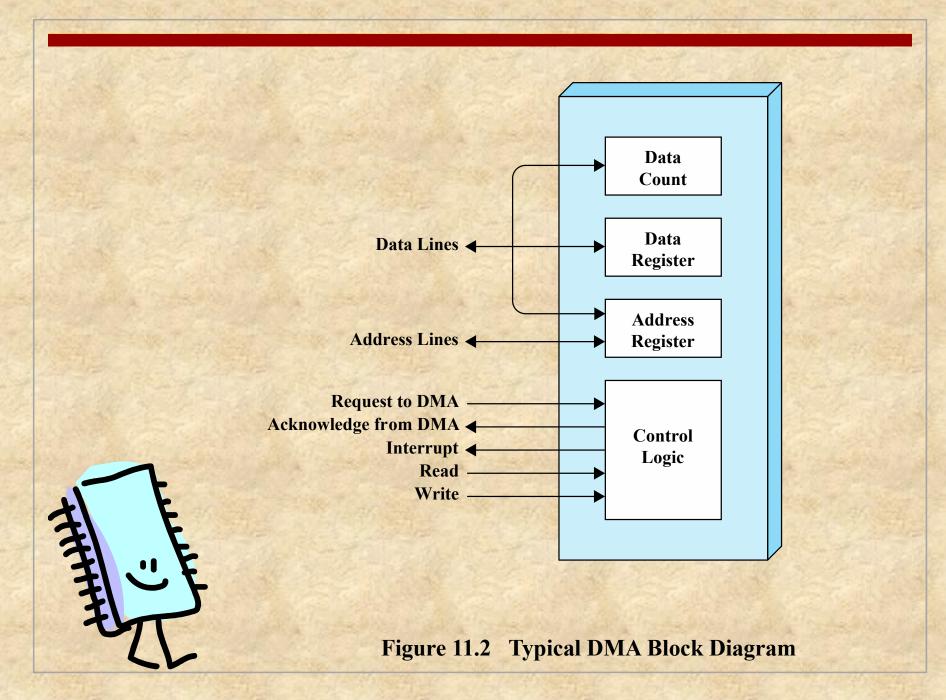
4

5

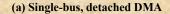
6

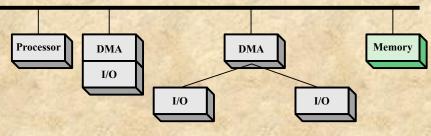
- Same configuration as step 2, but now interrupts are employed
- The I/O module is given direct control of memory via DMA
- The I/O module is enhanced to become a separate processor, with a specialized instruction set tailored for I/O

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









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

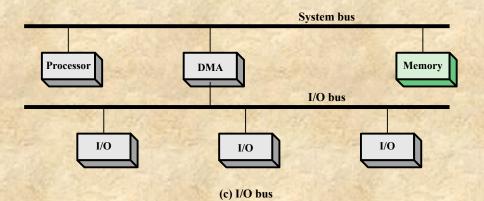


Figure 11.3 Alternative DMA Configurations

### **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

#### Hierarchical Design

- Functions of the operating system should be separated according to their complexity, their characteristic time scale, and their level of abstraction
- Leads to an organization of the operating system into a series of layers
- Each layer performs a related subset of the functions required of the operating system
- Layers should be defined so that changes in one layer do not require changes in other layers

## Buffering

 Perform input transfers in advance of requests being made and perform output transfers some time after the request is made

#### Block-oriented device

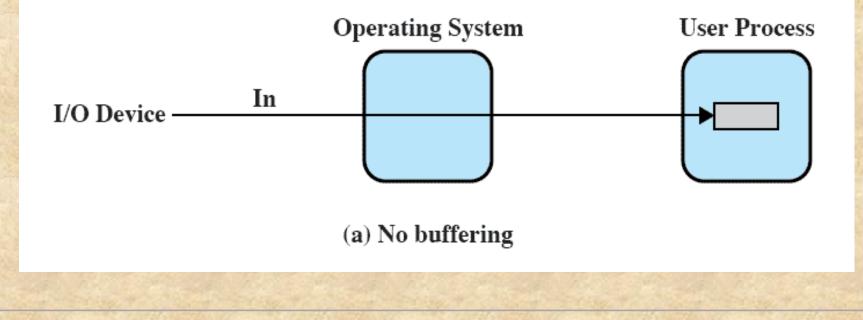
- stores information in blocks that are usually of fixed size
- transfers are made one block at a time
- possible to reference data by its block number
- disks and USB keys are examples

#### Stream-oriented device

- transfers data in and out as a stream of bytes
- no block structure
- terminals, printers, communications ports, and most other devices that are not secondary storage are examples

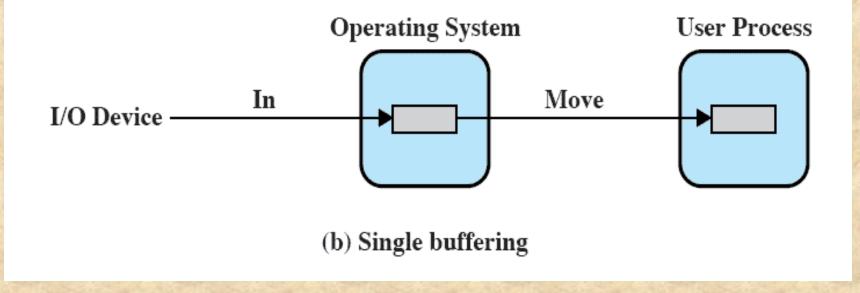
#### No Buffer

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



#### Single Buffer

 Operating system assigns a buffer in main memory for an I/O request



### Block-Oriented 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
- Generally provides a speedup compared to the lack of system buffering
- Disadvantages:
  - complicates the logic in the operating system
  - swapping logic is also affected

### Stream-Oriented Single Buffer

#### Line-at-a-time operation

230

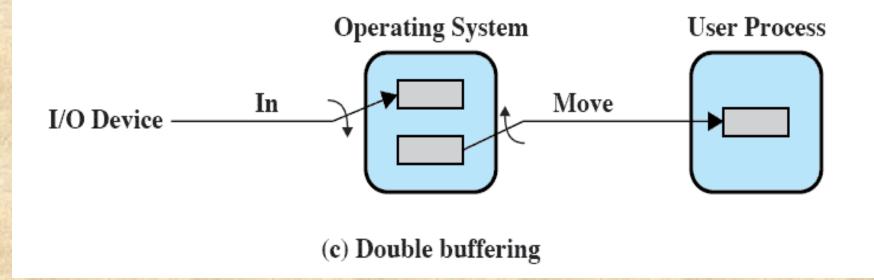
- appropriate for scroll-mode terminals (dumb terminals)
- user input is one line at a time with a carriage return signaling the end of a line
- output to the terminal is similarly one line at a time

- Byte-at-a-time operation
  - used on forms-mode terminals
  - when each keystroke is significant
  - other peripherals such as sensors and controllers



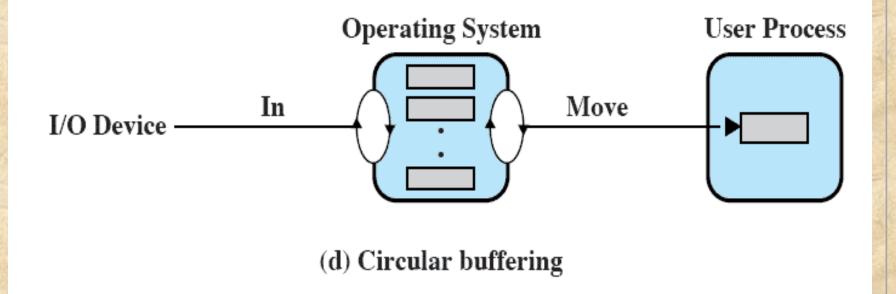
#### **Double Buffer**

- Use two system buffers instead of one
- A process can transfer data to or from one buffer while the operating system empties or fills the other buffer
- Also known as buffer swapping



#### **Circular Buffer**

- Two or more buffers are used
- Each individual buffer is one unit in a circular buffer
- Used when I/O operation must keep up with process



# The Utility of Buffering

- Technique that smoothes out peaks in I/O demand
  - with enough demand eventually all buffers become full and their advantage is lost
- When there is a variety of I/O and process activities to service, buffering can increase the efficiency of the OS and the performance of individual processes

## 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

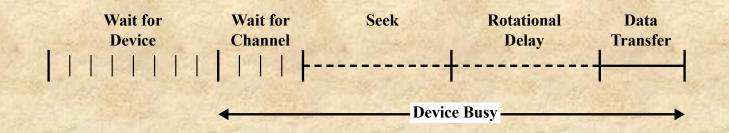


Figure 11.6 Timing of a Disk I/O Transfer

### Positioning the Read/Write Heads

- 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

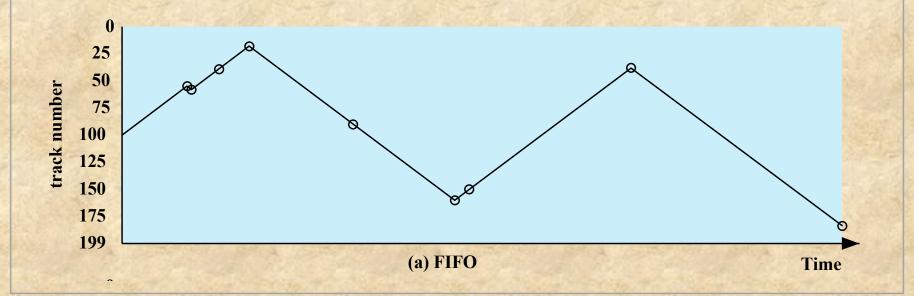
Name	Description	Remarks			
Selection according to requestor					
Random	Random scheduling	For analysis and simulation			
FIFO	First in first out	Fairest of them all			
PRI	Priority by process	Control outside of disk queu management			
LIFO	Last in first out	Maximize locality and resource utilization			
Selection according to requested item					
SSTF	Shortest service time first	High utilization, small queu			
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			

100) (starting lber Next acks track	o) SSTF g at track 100) Number of tracks	(starting a in the di increase num Next	CAN t track 100, rection of ing track hber) Number	(starting a in the di increasi	-SCAN t track 100, rection of ing track nber) Number
iber Next acks track	Number	in the dimincrease num	rection of ing track nber)	in the dim increase num	rection of ing track nber)
acks track		increasi nun Next	ing track nber)	increasi nun	ing track nber)
acks track		nun Next	nber)	nun	nber)
acks track		Next	,		,
acks track			Number	Next	Number 🖡
	of tracks		<b>A</b> . <b>T</b>	-	
ersed   accessed		track		track	of tracks
	d traversed	accessed	traversed	accessed	traversed
	10	1 5 0		1 = 0	
					50
3 58		160	10	160	10
19 55	3	184	24	184	24
21 39	16	90	94	18	166
72 38	1	58	32	38	20
70 18	20	55	3	39	1
10 150	132	39	16	55	16
12 160	10	38	1	58	3
46 184	24	18	20	90	32
55.3 Average	e 27.5	Average	27.8	Average	35.8
seek		seek		seek	
length		length		length	
	45 90   3 58   19 55   21 39   72 38   70 18   10 150   12 160   46 184   55.3 Average seek	ersedaccessedtraversed $45$ 9010358321955321391672381701820101501321216010461842455.3Average seek27.5	ersedaccessedtraversedaccessed459010150358321601955318421391690723815870182055101501323912160103846184241855.3Average seek27.5Average seek	ersedaccessedtraversedaccessedtraversed $45$ 90101505035832160101955318424213916909472381583270182055310150132391612160103814618424182055.3Average seek27.5Average seek27.8	acks ersedtrack accessedof tracks traversedtrack accessedof tracks traversedtrack accessed $45$ 90101505015035832160101601955318424184213916909418723815832387018205533910150132391655121601038158461842418209055.3Average seek27.5Average seek27.8Average seek

Table 11.2 Comparison of Disk Scheduling Algorithms

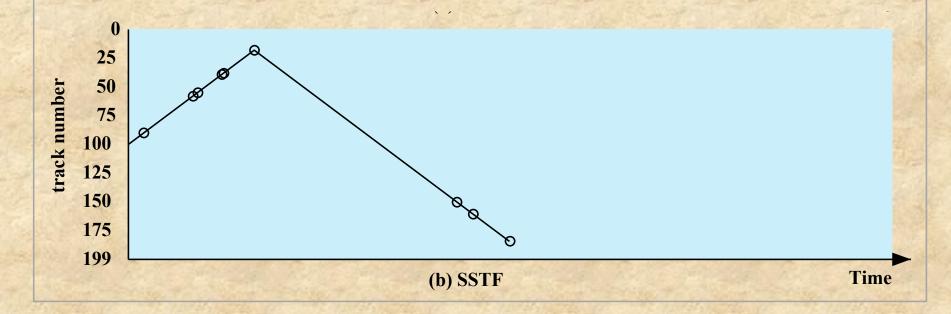
## First-In, First-Out (FIFO)

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



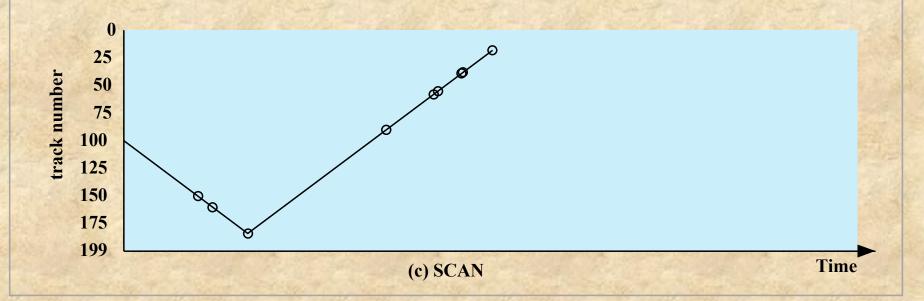
Shortest Service Time First (SSTF)  Select the disk I/O request that requires the least movement of the disk arm from its current position

 Always choose the minimum seek time



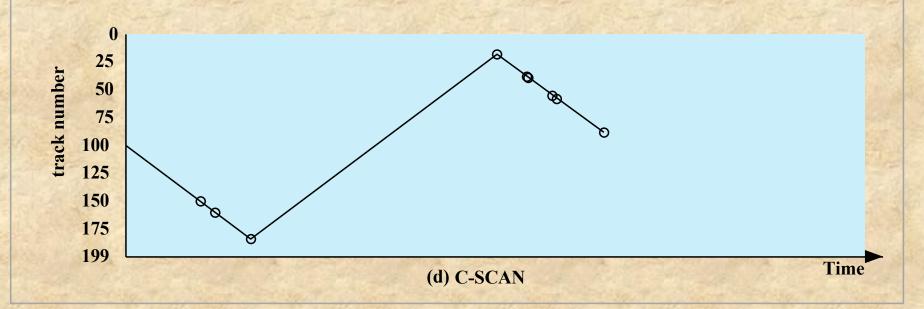
### SCAN

- 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



## 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



### 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

### **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

#### RAID

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

RAID is a set of physical disk drives viewed by the operating system as a single logical drive

> Design architectures share three characteristics:

redundant disk capacity is used to store parity information, which guarantees data recoverability in case of a disk failure

data are distributed across the physical drives of an array in a scheme known as striping



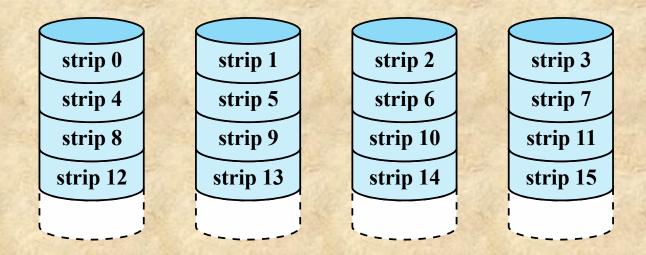
- The term was originally coined in a paper by a group of researchers at the University of California at Berkeley
  - the paper outlined various configurations and applications and introduced the definitions of the RAID levels
- Strategy employs multiple disk drives and distributes data in such a way 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

Category	Level	Description	Disks required	Data availability	Large I/O data transfer capacity	Small I/O request rate
Striping	0	Nonredundant	Ν	Lower than single disk	Very high	Very high for both read and write
Mirroring	1	Mirrored	2 <i>N</i>	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	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
access 3	3	Bit-interleaved parity	<i>N</i> + 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	<i>N</i> + 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	<i>N</i> +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

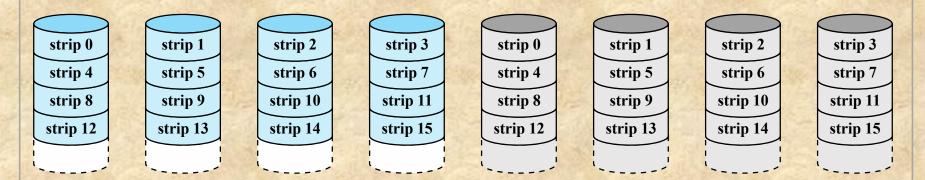
#### **RAID** Levels

- 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



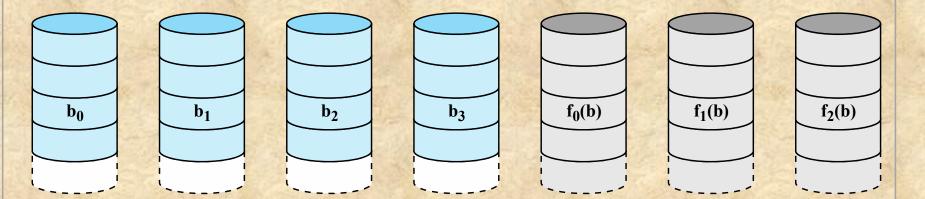
(a) RAID 0 (non-redundant)

- Redundancy is achieved by the simple expedient of duplicating all the data
- There is no "write penalty"
- When a drive fails the data may still be accessed from the second drive
- Principal disadvantage is the cost



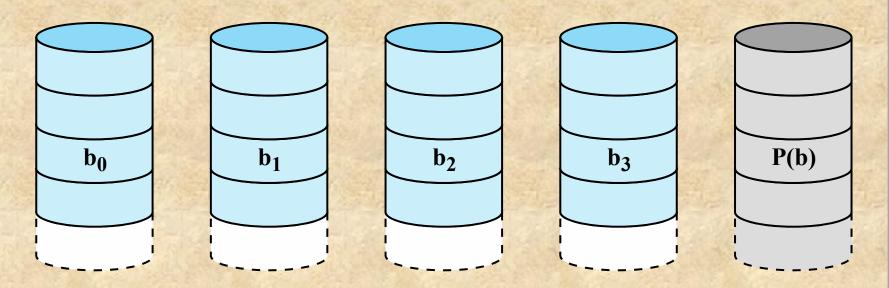
(b) RAID 1 (mirrored)

- 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



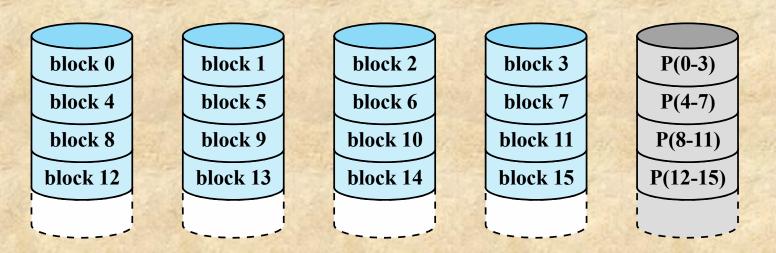
(c) RAID 2 (redundancy through Hamming code)

- 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



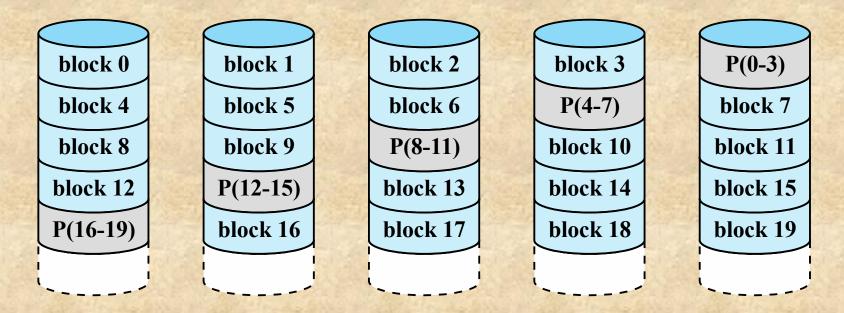
(d) RAID 3 (bit-interleaved parity)

- 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



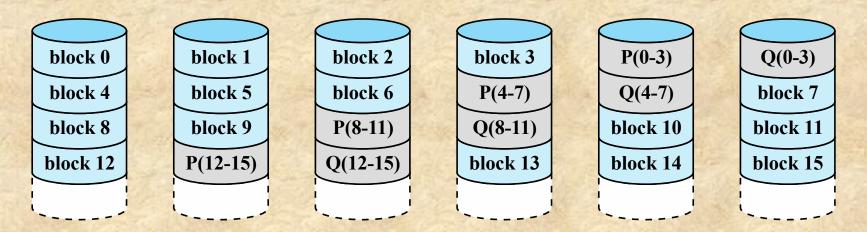
(e) RAID 4 (block-level parity)

- 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)

- 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

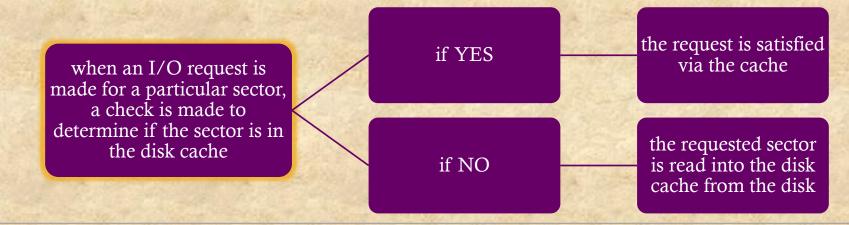


<sup>(</sup>g) RAID 6 (dual redundancy)

11

#### **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



## 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

