

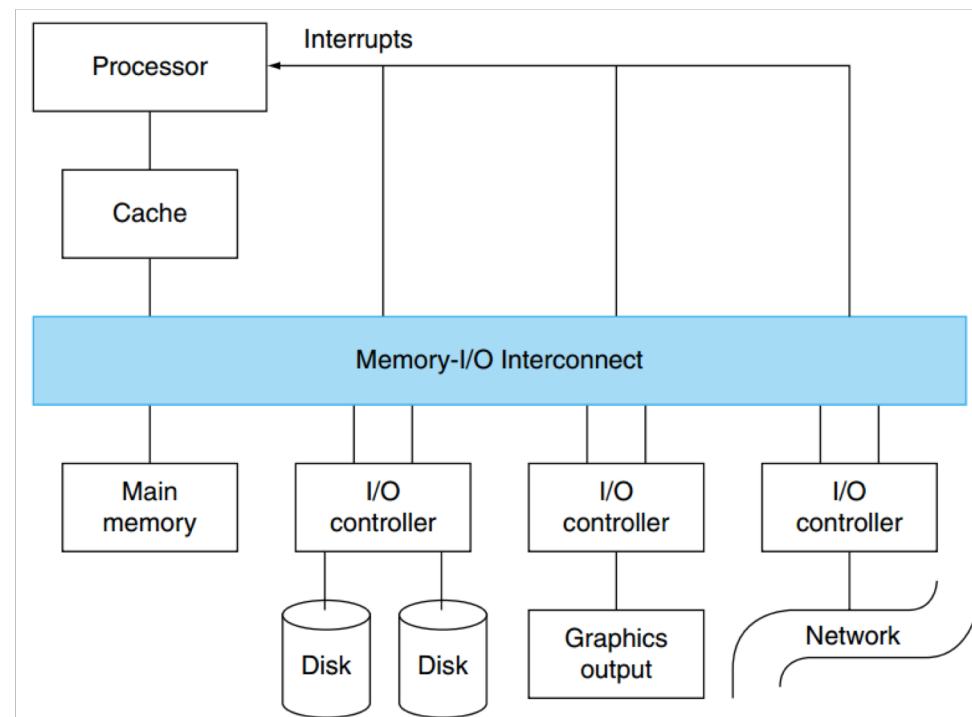


Chapter 6

Storage and Other I/O Topics

Introduction

- I/O devices can be characterized by
 - Behaviour: input, output, storage
 - Partner: human or machine
 - Data rate: bytes/sec, transfers/sec
- I/O bus connections



I/O System Characteristics

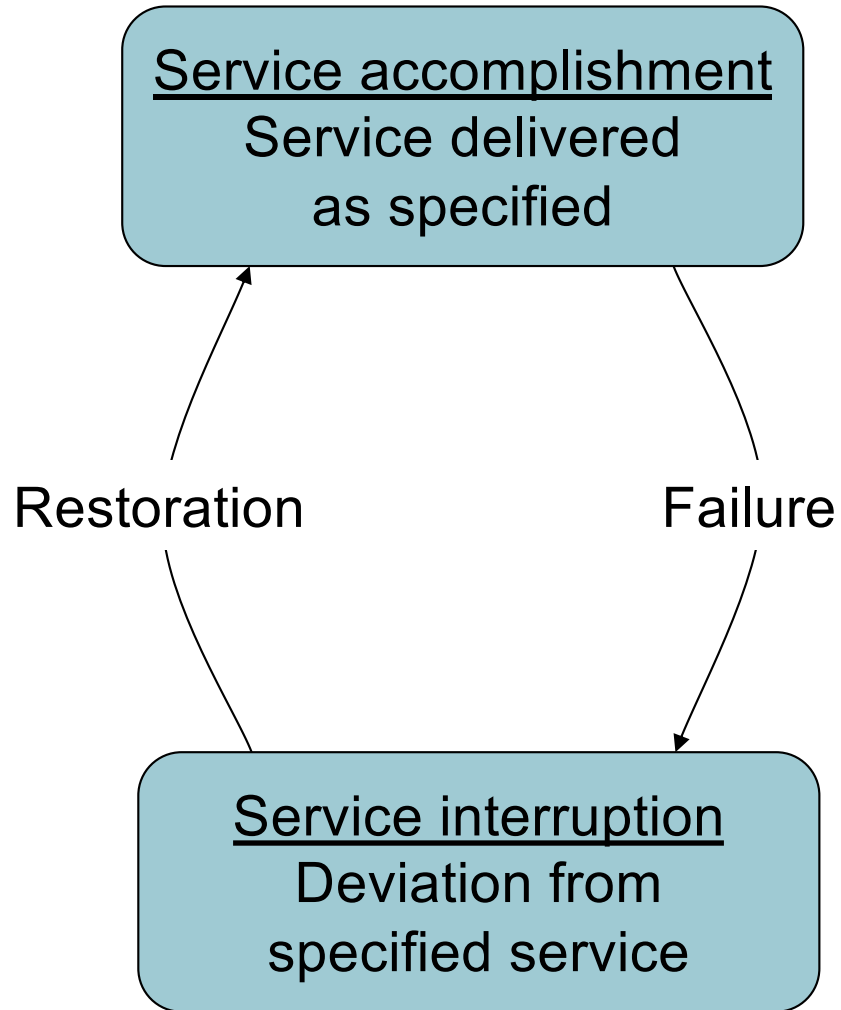
- Dependability is important
 - Particularly for storage devices
- Performance measures
 - Latency (response time)
 - Throughput (bandwidth)
 - Desktops & embedded systems
 - Mainly interested in response time & diversity of devices
 - Servers
 - Mainly interested in throughput & expandability of devices



DEPENDABILITY



Dependability



- Fault: failure of a component
 - May or may not lead to system failure

Dependability Measures

- Reliability: mean time to failure (MTTF)
- Service interruption: mean time to repair (MTTR)
- Mean time between failures
 - $MTBF = MTTF + MTTR$
- $Availability = MTTF / (MTTF + MTTR)$
- Improving Availability
 - Increase MTTF: fault avoidance, fault tolerance, fault forecasting
 - Reduce MTTR: improved tools and processes for diagnosis and repair

Dependability - Example

Mean Time Between Failures (MTBF), Mean Time To Replacement (MTTR), and Mean Time To Failure (MTTF) are useful metrics for evaluating the reliability and availability of a storage resource. Explore these concepts by answering the questions about devices with the following metrics.

	MTTF	MTTR
a.	3 Years	1 Day
b.	7 Years	3 Days

6.2.1 [5] <6.1, 6.2> Calculate the MTBF for each of the devices in the table.

6.2.2 [5] <6.1, 6.2> Calculate the availability for each of the devices in the table.

6.2.3 [5] <6.1, 6.2> What happens to availability as the MTTR approaches 0? Is this a realistic situation?

6.2.4 [5] <6.1, 6.2> What happens to availability as the MTTR gets very high, i.e., a device is difficult to repair? Does this imply the device has low availability?

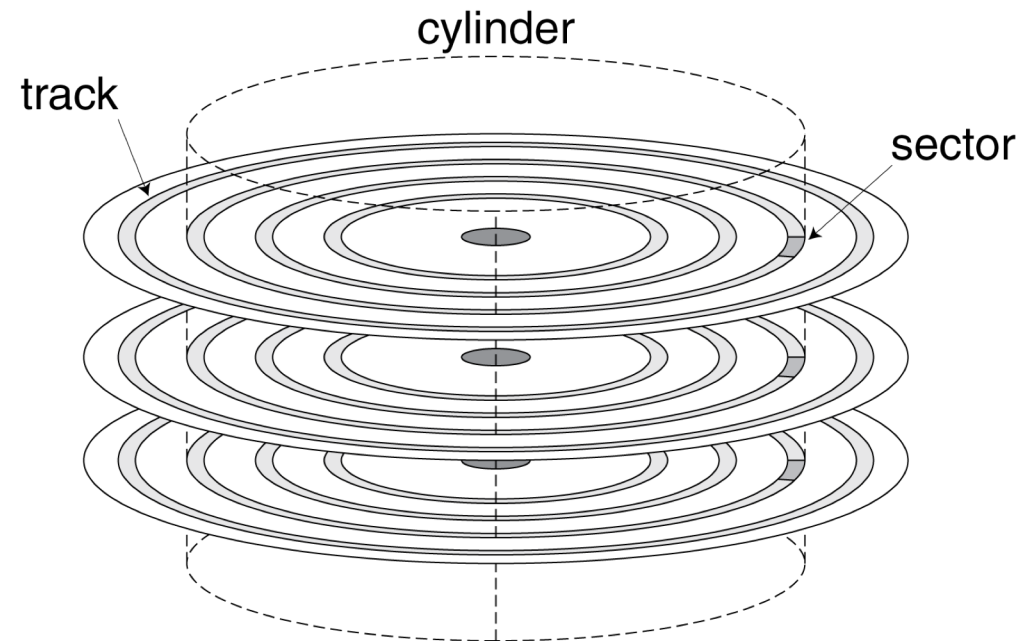


STORAGE



Disk Storage

- Nonvolatile, rotating magnetic storage



Disk Sectors and Access

- Each sector records
 - Sector ID
 - Data (512 bytes, 4096 bytes proposed)
 - Error correcting code (ECC)
 - Used to hide defects and recording errors
 - Gaps
- Access to a sector involves
 - Queuing delay if other accesses are pending
 - Seek: move the heads
 - Rotational latency
 - Data transfer
 - Controller overhead



Disk Access Example

- Given
 - 512B sector, 15,000rpm, 4ms average seek time, 100MB/s transfer rate, 0.2ms controller overhead, idle disk
- Average read time
 - 4ms seek time
 - + $\frac{1}{2} / (15,000/60) = 2\text{ms}$ rotational latency
 - + $512 / 100\text{MB/s} = 0.005\text{ms}$ transfer time
 - + 0.2ms controller delay
 - = 6.2ms
- If actual (vs. quoted) average seek time is 1ms
 - Average read time = 3.2ms



Disk Access - Exercise

Average and minimum times for reading and writing to storage devices are common measurements used to compare devices. Using techniques from [Chapter 6](#), calculate values related to read and write time for disks with the following characteristics.

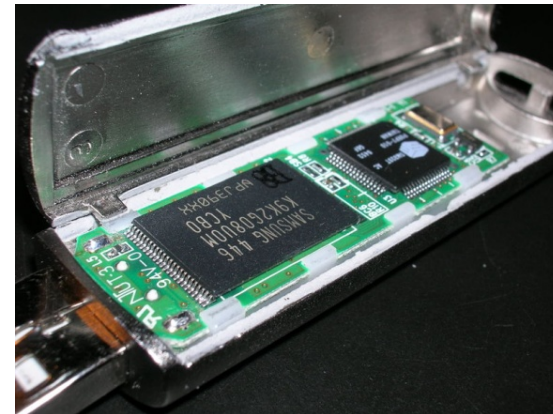
	Average Seek Time	RPM	Disk Transfer Rate	Controller Transfer Rate
a.	10 ms	7500	90 MB/s	100 MB/s
b.	7 ms	10,000	40 MB/s	200 MB/s

6.3.1 [10] <6.2, 6.3> Calculate the average time to read or write a 1024-byte sector for each disk listed in the table.

6.3.2 [10] <6.2, 6.3> Calculate the minimum time to read or write a 2048-byte sector for each disk listed in the table.

Flash Storage

- Nonvolatile semiconductor storage
 - 100× – 1000× faster than disk
 - Smaller, lower power, more robust
 - But more \$/GB (between disk and DRAM)



Flash Types

- NOR flash: bit cell like a NOR gate
 - Random read/write access
 - Used for instruction memory in embedded systems
- NAND flash: bit cell like a NAND gate
 - Denser (bits/area), but block-at-a-time access
 - Cheaper per GB
 - Used for USB keys, media storage, ...
- Flash bits wears out after 1000's of accesses
 - Not suitable for direct RAM or disk replacement
 - Wear leveling: remap data to less used blocks



RAID

- Redundant Array of Inexpensive (Independent) Disks
 - Use multiple smaller disks (c.f. one large disk)
 - Parallelism improves performance
 - Plus extra disk(s) for redundant data storage
- Provides fault tolerant storage system
 - Especially if failed disks can be “hot swapped”
- RAID 0
 - No redundancy (“AID”?)
 - Just stripe data over multiple disks
 - But it does improve performance



RAID 1 & 2

- RAID 1: Mirroring
 - $N + N$ disks, replicate data
 - Write data to both data disk and mirror disk
 - On disk failure, read from mirror
- RAID 2: Error correcting code (ECC)
 - $N + E$ disks (e.g., $10 + 4$)
 - Split data at bit level across N disks
 - Generate E -bit ECC
 - Too complex, not used in practice



RAID 3: Bit-Interleaved Parity

- $N + 1$ disks
 - Data striped across N disks at byte level
 - Redundant disk stores parity
 - Read access
 - Read all disks
 - Write access
 - Generate new parity and update all disks
 - On failure
 - Use parity to reconstruct missing data
- Not widely used



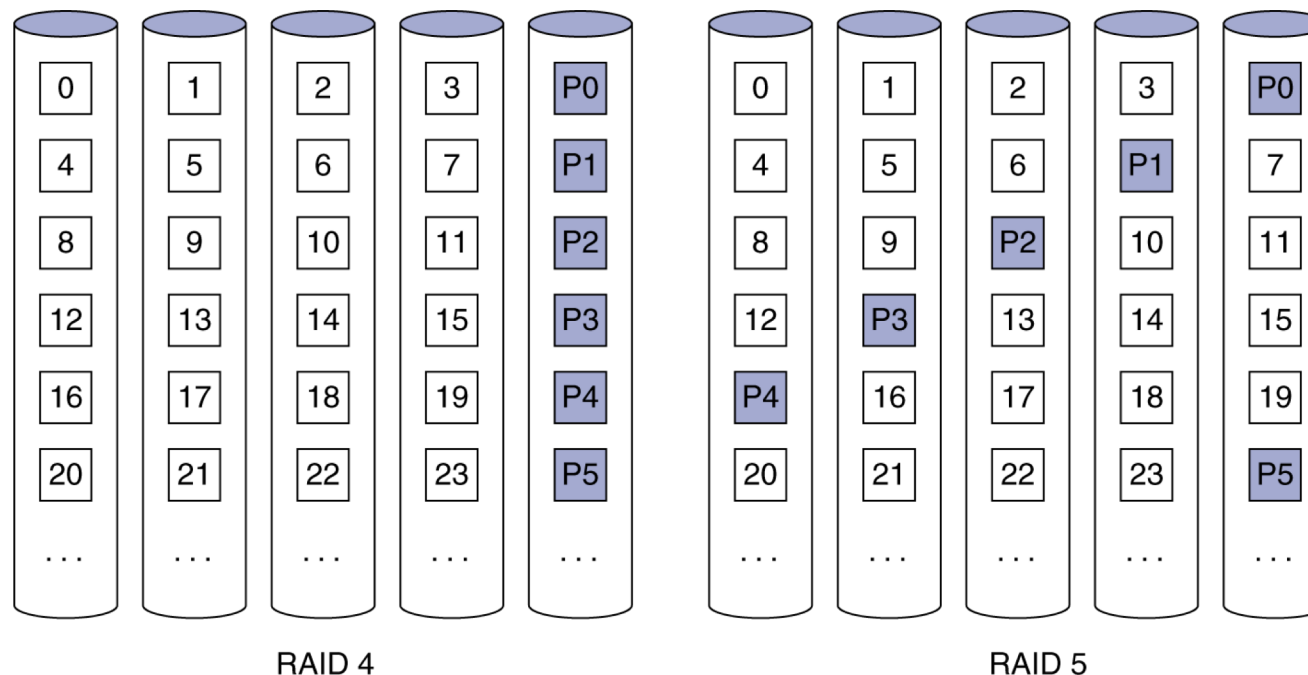
RAID 4: Block-Interleaved Parity

- N + 1 disks
 - Data striped across N disks at block level
 - Redundant disk stores parity for a group of blocks
 - Read access
 - Read only the disk holding the required block
 - Write access
 - Just read disk containing modified block, and parity disk
 - Calculate new parity, update data disk and parity disk
 - On failure
 - Use parity to reconstruct missing data
- Not widely used



RAID 5: Distributed Parity

- N + 1 disks
 - Like RAID 4, but parity blocks distributed across disks
 - Avoids parity disk being a bottleneck
- Widely used



RAID 6: P + Q Redundancy

- N + 2 disks
 - Like RAID 5, but two lots of parity
 - Greater fault tolerance through more redundancy
- Multiple RAID
 - More advanced systems give similar fault tolerance with better performance



RAID Summary

- RAID can improve performance and availability
 - High availability requires hot swapping
- Assumes independent disk failures
 - Too bad if the building burns down!
- See “Hard Disk Performance, Quality and Reliability”
 - <http://www.pcguide.com/ref/hdd/perf/index.htm>



INTERCONNECTS



Interconnecting Components

- Need interconnections between
 - CPU, memory, I/O controllers
- Bus: shared communication channel
 - Parallel set of wires for data and synchronization of data transfer
 - Can become a bottleneck
- Performance limited by physical factors
 - Wire length, number of connections
- More recent alternative: high-speed serial connections with switches
 - Like networks



Bus Types

- Processor-Memory buses
 - Short, high speed
 - Design is matched to memory organization
- I/O buses
 - Longer, allowing multiple connections
 - Specified by standards for interoperability
 - Connect to processor-memory bus through a bridge

Bus Signals and Synchronization

- Data lines
 - Carry address and data
 - Multiplexed or separate
- Control lines
 - Indicate data type, synchronize transactions
- Synchronous
 - Uses a bus clock
- Asynchronous
 - Uses request/acknowledge control lines for handshaking

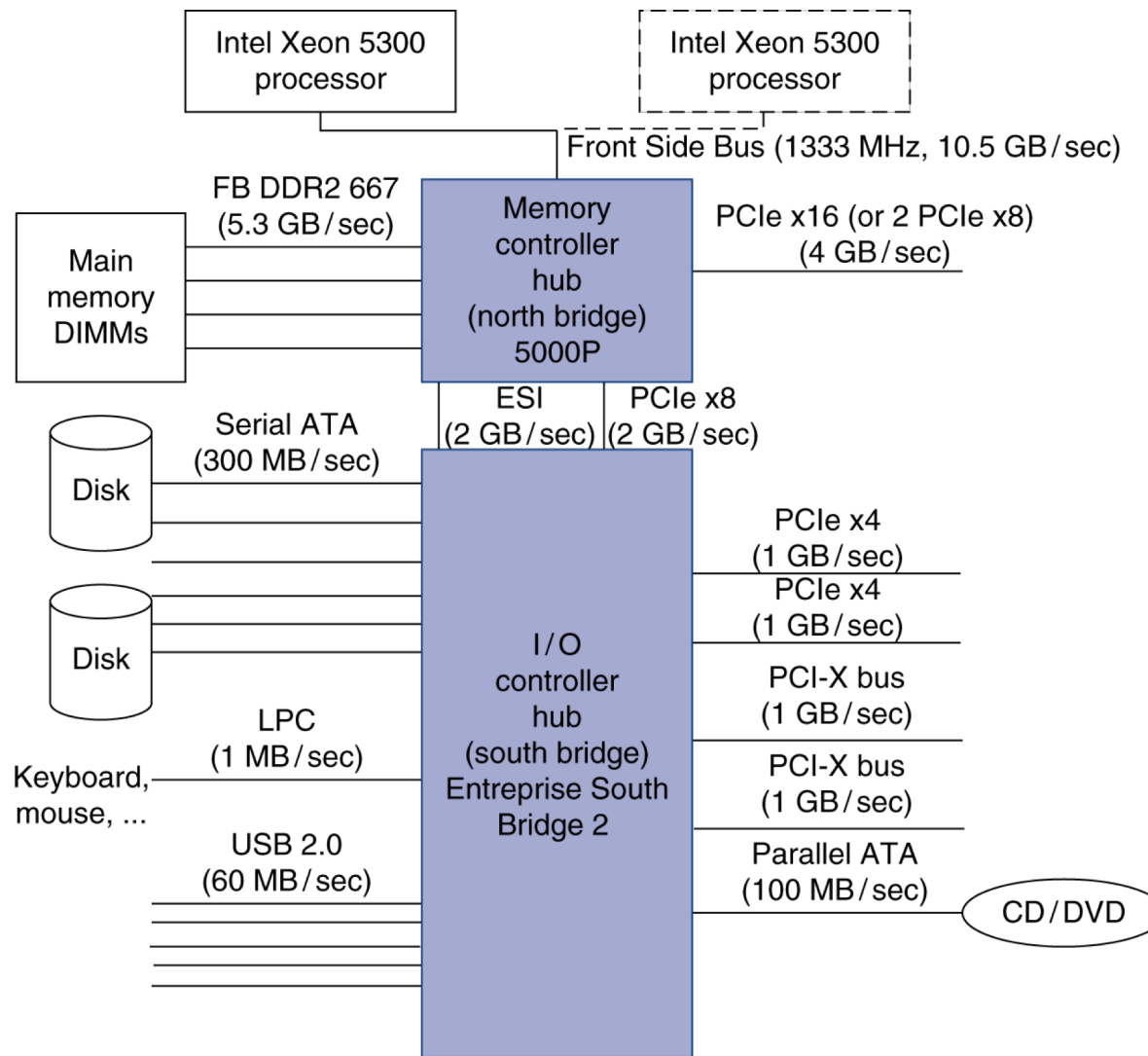


I/O Bus Examples

	Firewire	USB 2.0	PCI Express	Serial ATA	Serial Attached SCSI
Intended use	External	External	Internal	Internal	External
Devices per channel	63	127	1	1	4
Data width	4	2	2/lane	4	4
Peak bandwidth	50MB/s or 100MB/s	0.2MB/s, 1.5MB/s, or 60MB/s	250MB/s/lane 1×, 2×, 4×, 8×, 16×, 32×	300MB/s	300MB/s
Hot pluggable	Yes	Yes	Depends	Yes	Yes
Max length	4.5m	5m	0.5m	1m	8m
Standard	IEEE 1394	USB Implementers Forum	PCI-SIG	SATA-IO	INCITS TC T10



Typical x86 PC I/O System



I/O Management

- I/O is mediated by the OS
 - Multiple programs share I/O resources
 - Need protection and scheduling
 - I/O causes asynchronous interrupts
 - Same mechanism as exceptions
 - I/O programming is fiddly (detailed/complicated)
 - OS provides abstractions to programs



I/O Commands

- I/O devices are managed by I/O controller hardware
 - Transfers data to/from device
 - Synchronizes operations with software
- Command registers
 - Cause device to do something
- Status registers
 - Indicate what the device is doing and occurrence of errors
- Data registers
 - Write: transfer data to a device
 - Read: transfer data from a device



I/O Register Mapping

- Memory mapped I/O
 - Registers are addressed in same space as memory
 - Address decoder distinguishes between them
 - OS uses address translation mechanism to make them only accessible to kernel
- I/O instructions
 - Separate instructions to access I/O registers
 - Can only be executed in kernel mode
 - Example: x86



Polling

- Periodically check I/O status register
 - If device ready, do operation
 - If error, take action
- Common in small or low-performance real-time embedded systems
 - Predictable timing
 - Low hardware cost
- In other systems, wastes CPU time



Interrupts

- When a device is ready or error occurs
 - Controller interrupts CPU
- Interrupt is like an exception
 - But not synchronized to instruction execution
 - Can invoke handler between instructions
 - Cause information often identifies the interrupting device
- Priority interrupts
 - Devices needing more urgent attention get higher priority



I/O Data Transfer

- Polling and interrupt-driven I/O
 - CPU transfers data between memory and I/O data registers
 - Time consuming for high-speed devices
- Direct memory access (DMA)
 - OS provides starting address in memory
 - I/O controller transfers to/from memory autonomously
 - Controller interrupts on completion or error



Concluding Remarks

- I/O performance measures
 - Throughput, response time
 - Dependability and cost also important
- Buses used to connect CPU, memory, I/O controllers
 - Polling, interrupts, DMA
- RAID
 - Improves performance and dependability



Exercise

1. Given below are the parameters relevant to a hard disk access. Sector size 1024 Bytes, 25,000 rpm, 2 ms average seek time, 75 MB/s transfer rate, 0.4 ms controller overhead, and an idle disk. Compute the transfer time for the following cases.
 - a) Transfer of 512 KB data stored as a contiguous file
 - b) Transfer of 512 KB data stored in several 2 KB files/segments
2. For a hard disk, the manufacturer provides an average seek time. Describe the mechanisms used by the manufacturer to calculate the average seek time.
3. Input/output management is handled by the operating system of computers and not by the individual applications. State your stand on statement with reasons.



Exercise

- A. RAID 1 mirrors data among several disks. Assuming that inexpensive disks have lower Mean Time Between Failure (MTBF) than expensive disks, state how can redundancy using inexpensive disks result in a system with a higher MTBF? Use the definition of MTBF to explain your answer.
- B. For each of the activities listed below, will RAID 1 help better achieve their goals?
 - a. Online video services
 - b. High Performance Mathematical Computations
- C. In case of writing to disk, will RAID 3 be more efficient than RAID 4. Justify your answer.
- D. RAID 0 is stripping data across multiple disks. State the importance of RAID 0 in case of applications given below. In your answer include the impact of Input/Output efficiency compare to that of processing power of the computing system.
 - A. High performance mathematical computation
 - B. Video services

