

I/O System Characteristics(特性)

- Dependability(可靠性) is important
 - Particularly for storage devices
- Performance measures
 - Latency (response time, 响应时间)
 - Throughput (bandwidth, 带宽)
 - Desktops & embedded systems
 - Mainly interested in response time & diversity(多样性) of devices
 - Servers
 - Mainly interested in throughput & expandability(扩展性) of devices



Chapter 6

Storage and Other I/O Topics

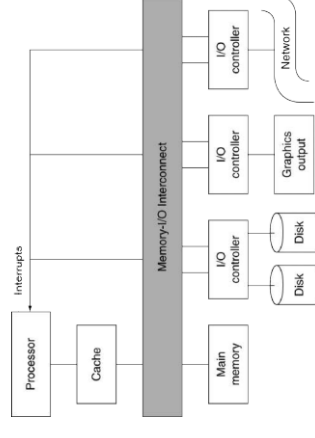
Dependability(可信度) Measures

- Reliability(可靠性): 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(可用度)
 - Increase MTTF: fault avoidance(避免), fault tolerance(容忍), fault forecasting(预测)
 - Reduce MTTR: improved tools and processes for diagnosis(诊断) and repair



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



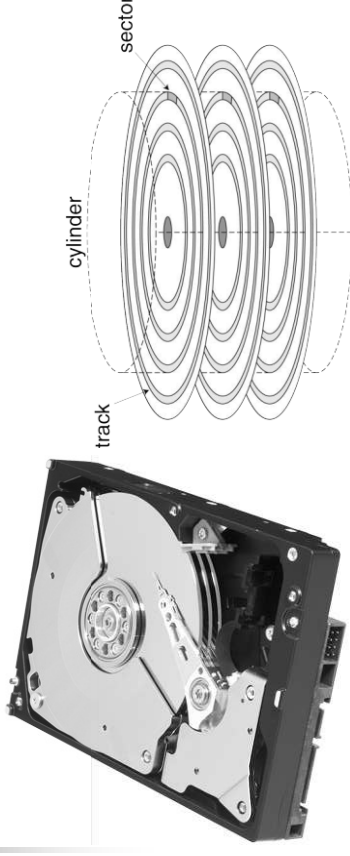
Disk Access Example

- 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
 - $+ 0.2\text{ms}$ controller delay
 - $= 6.2\text{ms}$
- If actual average seek time is 1ms
 - Average read time = 3.2ms



Disk Storage

- Nonvolatile(非易失性), rotating(旋转式) magnetic storage



Flash Storage

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



Disk Sectors and Access

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



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



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



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



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



I/O Management

- I/O is mediated by the OS
 - Multiple programs share I/O resources
 - Need protection and scheduling
 - I/O causes asynchronous(异步的) interrupts
 - Same mechanism as exceptions
 - I/O programming is fiddly(麻烦的, 精巧的)
 - OS provides abstractions to programs



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

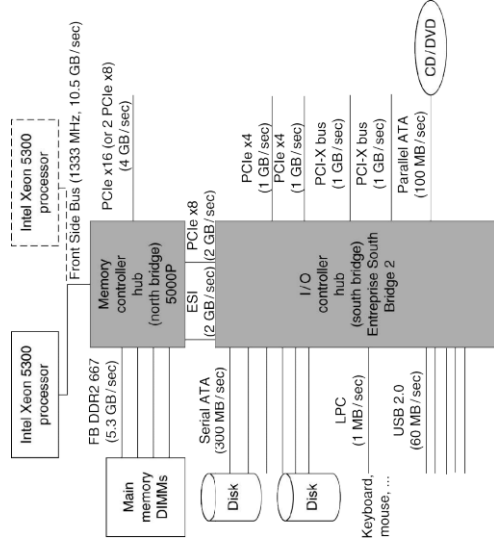


I/O Commands

- 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



Typical x86 PC I/O System



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



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
 - Can interrupt handler for a lower priority interrupt



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



Polling(轮询)

- 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



I/O vs. CPU Performance

- Amdahl's Law
 - Don't neglect I/O performance as parallelism increases compute performance
- Example
 - Benchmark takes 90s CPU time, 10s I/O time
 - Double the number of CPUs/2 years
 - I/O unchanged

Year	CPU time	I/O time	Elapsed time	% I/O time
now	90s	10s	100s	10%
+2	45s	10s	55s	18%
+4	23s	10s	33s	31%
+6	11s	10s	21s	47%



DMA/Cache Interaction

- If DMA writes to a memory block that is cached
 - Cached copy becomes stale
- If write-back cache has dirty block, and DMA reads memory block
 - Reads stale data
- Need to ensure cache coherence
 - Flush blocks from cache if they will be used for DMA
 - Or use non-cacheable memory locations for I/O



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
- I/O benchmarks
 - TPC, SPECIFS, SPECWeb



§ 6.7 I/O Performance Measures: ...

Measuring I/O Performance

- I/O performance depends on
 - Hardware: CPU, memory, controllers, buses
 - Software: operating system, database management system, application
 - Workload: request rates and patterns
- I/O system design can trade-off between response time and throughput
 - Measurements(测量, 量度) of throughput often done with constrained response-time



§ 6.9 Parallelism and I/O: RAID