

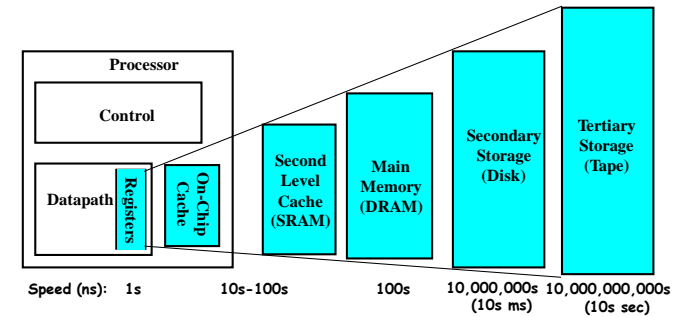
CS162
Operating Systems and
Systems Programming
Lecture 16

I/O Systems

March 31, 2008
Prof. Anthony D. Joseph
<http://inst.eecs.berkeley.edu/~cs162>

Review: Memory Hierarchy of a Modern Computer System

- Take advantage of the principle of locality to:
 - Present as much memory as in the cheapest technology
 - Provide access at speed offered by the fastest technology



3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.2

Goals for Today

- I/O Systems
 - Hardware Access
 - Device Drivers
- Queuing Theory

Note: Some slides and/or pictures in the following are adapted from slides ©2005 Silberschatz, Galvin, and Gagne. Many slides generated from my lecture notes by Kubiawicz.

3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.3

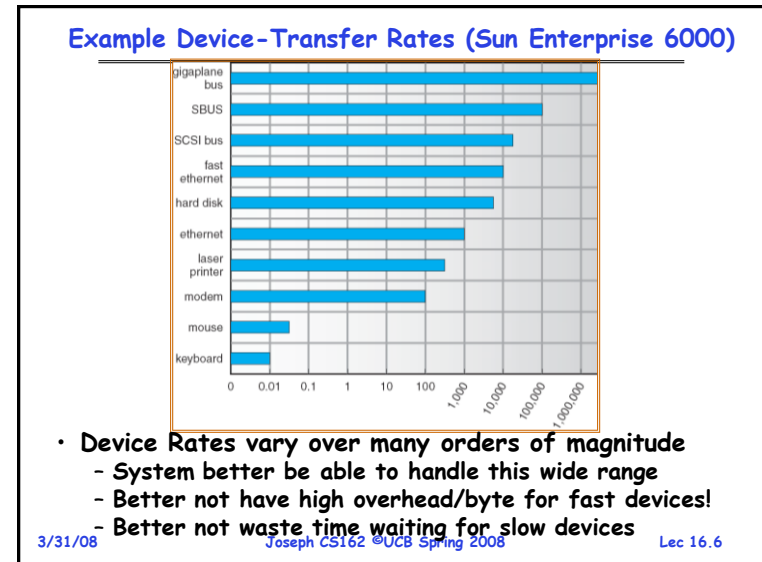
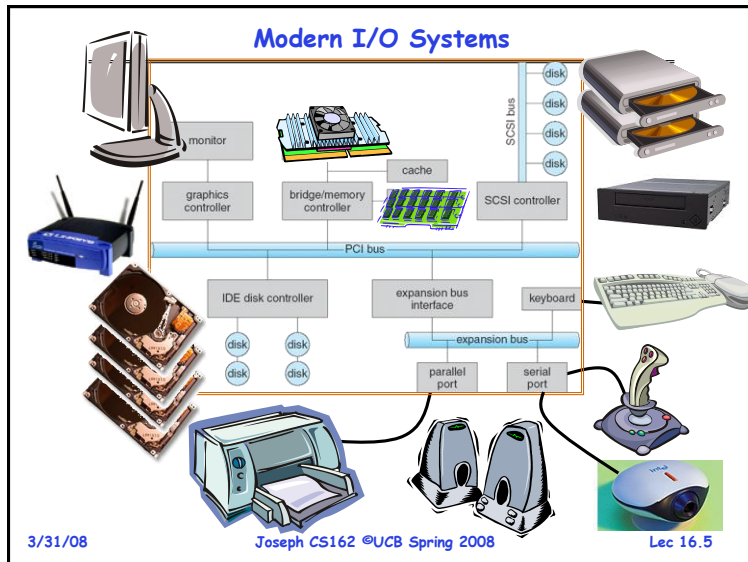
The Requirements of I/O

- So far in this course:
 - We have learned how to manage CPU, memory
- What about I/O?
 - Without I/O, computers are useless (disembodied brains?)
 - But... thousands of devices, each slightly different
 - » How can we standardize the interfaces to these devices?
 - Devices unreliable: media failures and transmission errors
 - » How can we make them reliable???
 - Devices unpredictable and/or slow
 - » How can we manage them if we don't know what they will do or how they will perform?
- Some operational parameters:
 - Byte/Block
 - » Some devices provide single byte at a time (e.g. keyboard)
 - » Others provide whole blocks (e.g. disks, networks, etc)
 - Sequential/Random
 - » Some devices must be accessed sequentially (e.g. tape)
 - » Others can be accessed randomly (e.g. disk, cd, etc.)
 - Polling/Interrupts
 - » Some devices require continual monitoring
 - » Others generate interrupts when they need service

3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.4



The Goal of the I/O Subsystem

- Provide Uniform Interfaces, Despite Wide Range of Different Devices
 - This code works on many different devices:


```
FILE fd = fopen("/dev/something", "rw");
for (int i = 0; i < 10; i++) {
    fprintf(fd, "Count %d\n", i);
}
close(fd);
```
 - Why? Because code that controls devices ("device driver") implements standard interface.
- We will try to get a flavor for what is involved in actually controlling devices in rest of lecture
 - Can only scratch surface!

3/31/08 Joseph CS162 ©UCB Spring 2008 Lec 16.7

Administrivia

- Would you like an extra 5% for your course grade?
 - Attend lectures and sections! 5% of grade is participation
 - Midterm 1 was only 15%
- Project #3 design doc due next Monday (4/7) at 11:59pm
- Midterm #2 is in two weeks (Wed 4/16) 6-7:30pm in 10 Evans

3/31/08 Joseph CS162 ©UCB Spring 2008 Lec 16.8

Want Standard Interfaces to Devices

- **Block Devices:** *e.g.* disk drives, tape drives, DVD-ROM
 - Access blocks of data
 - Commands include `open()`, `read()`, `write()`, `seek()`
 - Raw I/O or file-system access
 - Memory-mapped file access possible
- **Character Devices:** *e.g.* keyboards, mice, serial ports, some USB devices
 - Single characters at a time
 - Commands include `get()`, `put()`
 - Libraries layered on top allow line editing
- **Network Devices:** *e.g.* Ethernet, Wireless, Bluetooth
 - Different enough from block/character to have own interface
 - Unix and Windows include **socket** interface
 - » Separates network protocol from network operation
 - » Includes `select()` functionality
 - Usage: pipes, FIFOs, streams, queues, mailboxes

3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.9

How Does User Deal with Timing?

- **Blocking Interface:** "Wait"
 - When request data (*e.g.* `read()` system call), put process to sleep until data is ready
 - When write data (*e.g.* `write()` system call), put process to sleep until device is ready for data
- **Non-blocking Interface:** "Don't Wait"
 - Returns quickly from read or write request with count of bytes successfully transferred
 - Read may return nothing, write may write nothing
- **Asynchronous Interface:** "Tell Me Later"
 - When request data, take pointer to user's buffer, return immediately; later kernel fills buffer and notifies user
 - When send data, take pointer to user's buffer, return immediately; later kernel takes data and notifies user

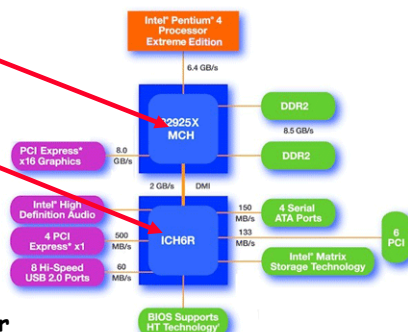
3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.10

Main components of Intel Chipset: Pentium 4

- **Northbridge:**
 - Handles memory
 - Graphics
- **Southbridge: I/O**
 - PCI bus
 - Disk controllers
 - USB controllers
 - Audio
 - Serial I/O
 - Interrupt controller
 - Timers

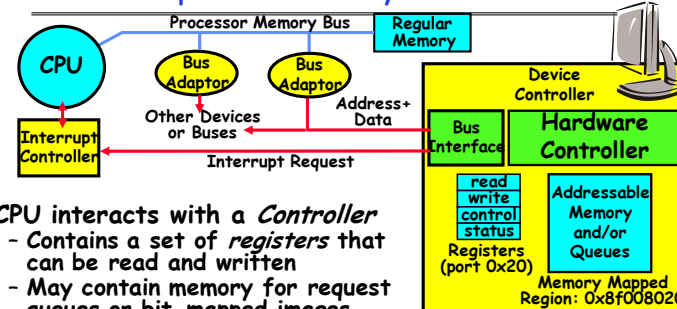


3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.11

How does the processor actually talk to the device?



- **CPU interacts with a Controller**
 - Contains a set of *registers* that can be read and written
 - May contain memory for request queues or bit-mapped images
- **Regardless of the complexity of the connections and buses, processor accesses registers in two ways:**
 - **I/O instructions:** in/out instructions
 - » Example from the Intel architecture: `out 0x21, AL`
 - **Memory mapped I/O:** load/store instructions
 - » Registers/memory appear in physical address space
 - » I/O accomplished with load and store instructions

3/31/08

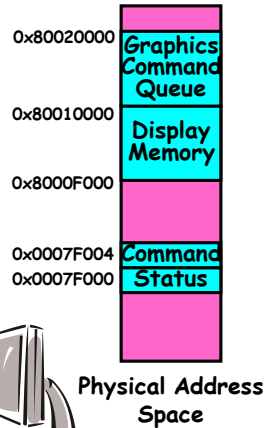
Joseph CS162 ©UCB Spring 2008

Lec 16.12

Example: Memory-Mapped Display Controller

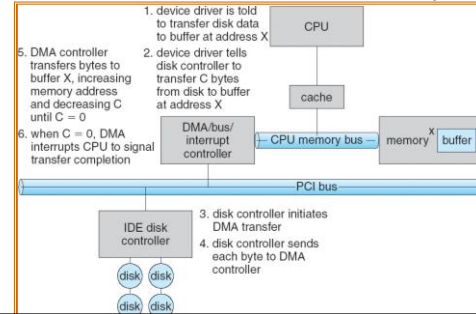
- **Memory-Mapped:**

- Hardware maps control registers and display memory into physical address space
 - » Addresses set by hardware jumpers or programming at boot time
- Simply writing to display memory (also called the "frame buffer") changes image on screen
 - » Addr: 0x8000F000–0x8000FFFF
- Writing graphics description to command-queue area
 - » Say enter a set of triangles that describe some scene
 - » Addr: 0x80010000–0x8001FFFF
- Writing to the command register may cause on-board graphics hardware to do something
 - » Say render the above scene
 - » Addr: 0x0007F004

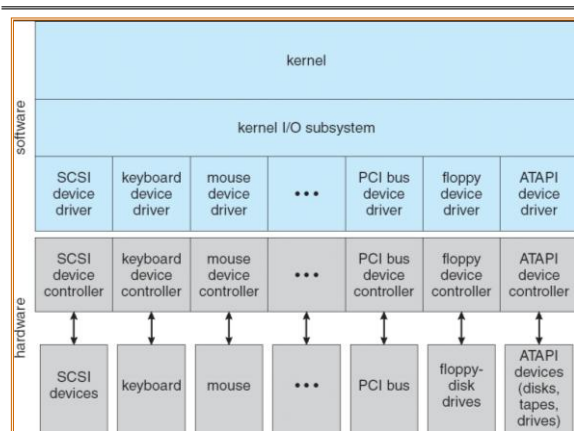


Transferring Data To/From Controller

- **Programmed I/O:**
 - Each byte transferred via processor in/out or load/store
 - Pro: Simple hardware, easy to program
 - Con: Consumes processor cycles proportional to data size
- **Direct Memory Access:**
 - Give controller access to memory bus
 - Ask it to transfer data to/from memory directly
- Sample interaction with DMA controller (from book):



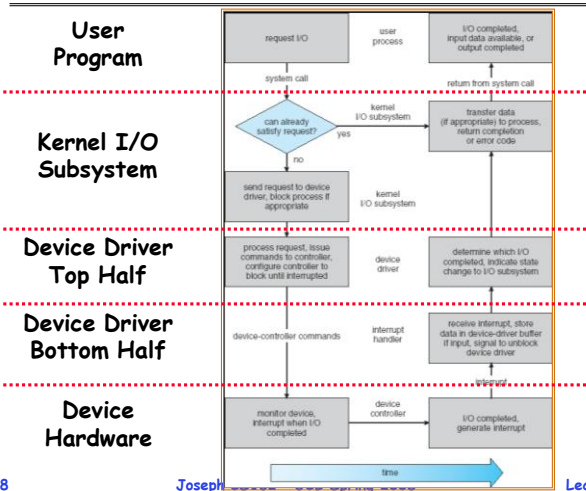
A Kernel I/O Structure



Device Drivers

- **Device Driver:** Device-specific code in the kernel that interacts directly with the device hardware
 - Supports a standard, internal interface
 - Same kernel I/O system can interact easily with different device drivers
 - Special device-specific configuration supported with the `ioctl()` system call
- Device Drivers typically divided into two pieces:
 - Top half: accessed in call path from system calls
 - » Implements a set of **standard, cross-device calls** like `open()`, `close()`, `read()`, `write()`, `ioctl()`, `strategy()`
 - » This is the kernel's interface to the device driver
 - » Top half will *start I/O* to device, may put thread to sleep until finished
 - Bottom half: run as interrupt routine
 - » Gets input or transfers next block of output
 - » May wake sleeping threads if I/O now complete

Life Cycle of An I/O Request



I/O Device Notifying the OS

- The OS needs to know when:
 - The I/O device has completed an operation
 - The I/O operation has encountered an error
- **I/O Interrupt:**
 - Device generates an interrupt whenever it needs service
 - Handled in bottom half of device driver
 - » Often run on special kernel-level stack
 - Pro: handles unpredictable events well
 - Con: interrupts relatively high overhead
- **Polling:**
 - OS periodically checks a device-specific status register
 - » I/O device puts completion information in status register
 - » Could use timer to invoke lower half of drivers occasionally
 - Pro: low overhead
 - Con: may waste many cycles on polling if infrequent or unpredictable I/O operations
- Actual devices combine both polling and interrupts
 - For instance: High-bandwidth network device:
 - » Interrupt for first incoming packet
 - » Poll for following packets until hardware empty

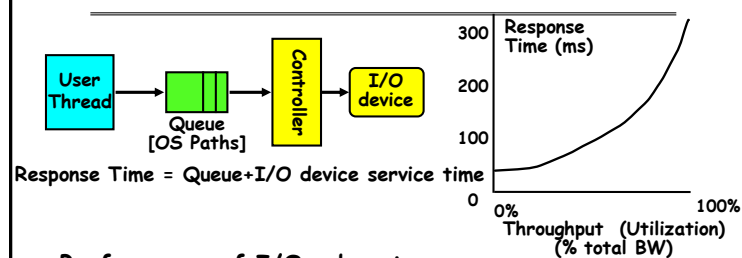
3/31/08

Joseph CS162 @UCB Spring 2008

Lec 16.18

BREAK

I/O Performance



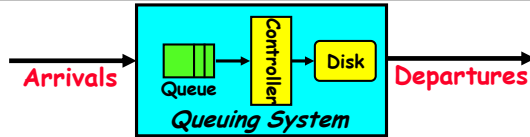
- Performance of I/O subsystem
 - Metrics: Response Time, Throughput
 - Contributing factors to latency:
 - » Software paths (can be loosely modeled by a queue)
 - » Hardware controller
 - » I/O device service time
- Queuing behavior:
 - Can lead to big increases of latency as utilization approaches 100%

3/31/08

Joseph CS162 @UCB Spring 2008

Lec 16.20

Introduction to Queuing Theory



- What about queuing time??
 - Let's apply some queuing theory
 - Queuing Theory applies to long term, steady state behavior \Rightarrow Arrival rate = Departure rate
- Little's Law:
 - Mean # tasks in system = arrival rate \times mean response time**
 - Observed by many, Little was first to prove
 - Simple interpretation: you should see the same number of tasks in queue when entering as when leaving.
- Applies to any system in equilibrium, as long as nothing in black box is creating or destroying tasks
 - **Typical queuing theory doesn't deal with transient behavior, only steady-state behavior**

3/31/08

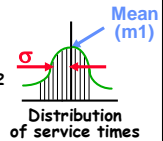
Joseph CS162 ©UCB Spring 2008

Lec 16.21

Background: Use of random distributions

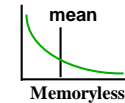
- Server spends variable time with customers

- Mean (Average) $m1 = \Sigma p(T) \times T$
- Variance $\sigma^2 = \Sigma p(T) \times (T - m1)^2 = \Sigma p(T) \times T^2 - m1^2$
- Squared coefficient of variance: $C = \sigma^2 / m1^2$



- Important values of C:

- No variance or deterministic $\Rightarrow C=0$
- "memoryless" or exponential $\Rightarrow C=1$
 - » Past tells nothing about future
 - » Many complex systems (or aggregates) well described as memoryless
- Disk response times $C \approx 1.5$ (wider variance \Rightarrow long tail)



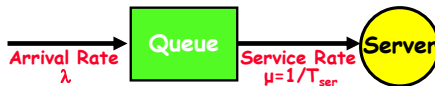
3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.22

A Little Queuing Theory: Some Results

- Assumptions:
 - System in equilibrium; No limit to the queue
 - Time between successive arrivals is random and memoryless



- Parameters that describe our system:
 - λ : mean number of arriving customers/second
 - T_{ser} : mean time to service a customer ("m1")
 - C: squared coefficient of variance = $\sigma^2 / m1^2$
 - μ : service rate = $1/T_{ser}$
 - u: server utilization ($0 \leq u \leq 1$): $u = \lambda / \mu = \lambda \times T_{ser}$
- Parameters we wish to compute:
 - T_q : Time spent in queue
 - L_q : Length of queue = $\lambda \times T_q$ (by Little's law)
- Results:
 - Memoryless service distribution ($C = 1$):
 - » Called M/M/1 queue: $T_q = T_{ser} \times u / (1 - u)$
 - General service distribution (no restrictions), 1 server:
 - » Called M/G/1 queue: $T_q = T_{ser} \times \frac{1}{2}(1 + C) \times u / (1 - u)$

3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.23

A Little Queuing Theory: An Example

- Example Usage Statistics:
 - User requests $10 \times 8\text{KB}$ disk I/Os per second
 - Requests & service exponentially distributed ($C=1.0$)
 - Avg. service = 20 ms (From controller+seek+rot+trans)

- Questions:
 - How utilized is the disk?
 - » Ans: server utilization, $u = \lambda T_{ser}$
 - What is the average time spent in the queue?
 - » Ans: T_q
 - What is the number of requests in the queue?
 - » Ans: $L_q = \lambda T_q$ (Little's law)
 - What is the avg response time for disk request?
 - » Ans: $T_{sys} = T_q + T_{ser}$

- Computation:

$$\begin{aligned} \lambda & \text{ (avg \# arriving customers/s)} = 10/s \\ T_{ser} & \text{ (avg time to service customer)} = 20 \text{ ms (0.02s)} \\ u & \text{ (server utilization)} = \lambda \times T_{ser} = 10/s \times 0.02s = 0.2 \\ T_q & \text{ (avg time/customer in queue)} = T_{ser} \times u / (1 - u) \\ & = 20 \times 0.2 / (1 - 0.2) = 20 \times 0.25 = 5 \text{ ms (0.005s)} \\ L_q & \text{ (avg length of queue)} = \lambda \times T_q = 10/s \times 0.005s = 0.05 \\ T_{sys} & \text{ (avg time/customer in system)} = T_q + T_{ser} = 25 \text{ ms} \end{aligned}$$

3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.24

Summary

- **Working Set:** Set of pgs touched by a process recently
- **Thrashing:** a process is busy swapping pages in and out
 - Process will thrash if working set doesn't fit in memory
 - Need to swap out a process
- **I/O Devices Types:**
 - Many different speeds (0.1 bytes/sec to GBytes/sec)
 - Different Access Patterns: block, char, net devices
 - Different Access Timing: Non-/Blocking, Asynchronous
- **I/O Controllers:** Hardware that controls actual device
 - CPU accesses thru I/O insts, ld/st to special phy memory
 - Report results thru interrupts or a status register polling
- **Device Driver:** Device-specific code in kernel
- **Queuing Latency:**
 - M/M/1 and M/G/1 queues: simplest to analyze
 - As utilization approaches 100%, latency $\rightarrow \infty$

$$T_q = T_{ser} \times \frac{1}{2}(1+C) \times u/(1-u)$$

3/31/08

Joseph CS162 ©UCB Spring 2008

Lec 16.25