

# Disk Scheduling and Input/ Output Management

Chapter 14

# Principles of I/O Hardware

- » I/O can be roughly divided into two categories
  - » Block devices - Stores information in a fixed size block each with its own address.  
Transfers are in units of blocks
  - » Character devices - Delivers or accepts a stream of characters without regard to a block structure.

# Principles of I/O Hardware

- » Classification scheme is not perfect.
- » Clocks, not block addressable, don't generate or accept streams.
- » Memory mapped screens
- » Touch screens

# Typical Device Data Rates

Device	Data rate
Keyboard	10 bytes/sec
Mouse	100 bytes/sec
56K modem	7 KB/sec
Scanner at 300 dpi	1 MB/sec
Digital camcorder	3.5 MB/sec
4x Blu-ray disc	18 MB/sec
802.11n Wireless	37.5 MB/sec
USB 2.0	60 MB/sec
FireWire 800	100 MB/sec
Gigabit Ethernet	125 MB/sec
SATA 3 disk drive	600 MB/sec
USB 3.0	625 MB/sec
SCSI Ultra 5 bus	640 MB/sec
Single-lane PCIe 3.0 bus	985 MB/sec
Thunderbolt 2 bus	2.5 GB/sec
SONET OC-768 network	5 GB/sec

# Device Controllers

- » I/O often consists of a mechanical component and an electronic component.
- » Generally split into two modular portions
- » Electric component is the device controller or device adapter.
- » Mechanical component is the device itself

# Device Controllers

- » If the interface between the controller and device is a standard interface companies can make controllers or devices that fit that interface.
- » ANSI, IEEE, ISO or de facto.

# Standards

- » de facto standard - One vendor comes up with a good idea and other vendors follow the lead. Or some organizations come together and agree on a standard.
  - QWERTY keyboard layout
  - Microsoft Word DOC
- » de jure standard - Technically have the force of law behind them. Set by professional, national, or international organizations such as IEEE, ANSI and ISO.
  - PDF and HTML ( started de facto eventually made de jury )
  - 802.11

# HOW STANDARDS PROLIFERATE:

(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC)

SITUATION:  
THERE ARE  
14 COMPETING  
STANDARDS.

14?! RIDICULOUS!  
WE NEED TO DEVELOP  
ONE UNIVERSAL STANDARD  
THAT COVERS EVERYONE'S  
USE CASES.



SOON:

SITUATION:  
THERE ARE  
15 COMPETING  
STANDARDS.



# Device Controllers

- » Interface is often very low level.
- » Disk formatted with 2,000,000 sectors.
- » What comes off the disk is a serial bit stream.
- » Preamble
- » 4096 bit sector
- » Checksum or Error Correcting Code

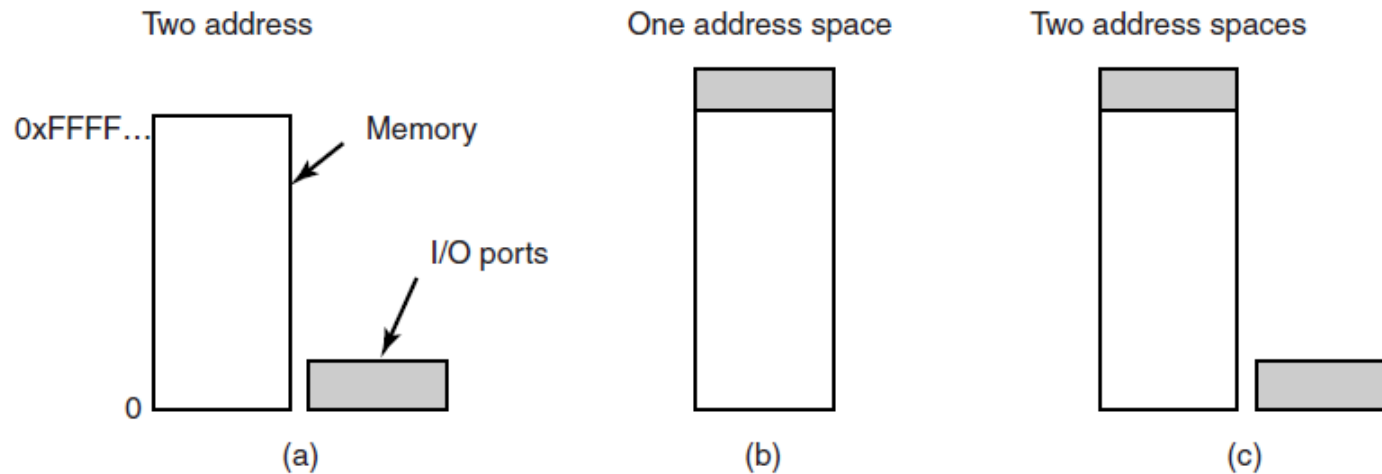
# Device Controllers

- » Controller converts the serial bit stream into a block of bytes and performs any error correction.
- » Once error free it can be copied to main memory.

# Memory-Mapped I/O

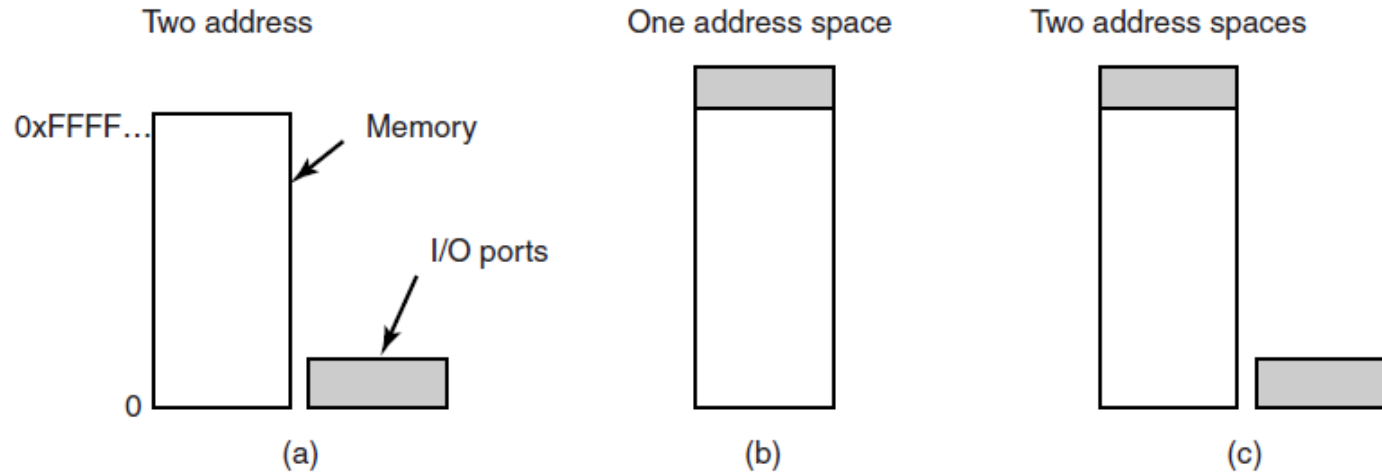
- » Each controller has a couple registers used for communicating with the CPU.
- » Command the device to deliver data
- » Switch it on/off
- » Other actions
- » Device also maintains data buffer.
  - » Video RAM is effectively a data buffer

# Device Controllers



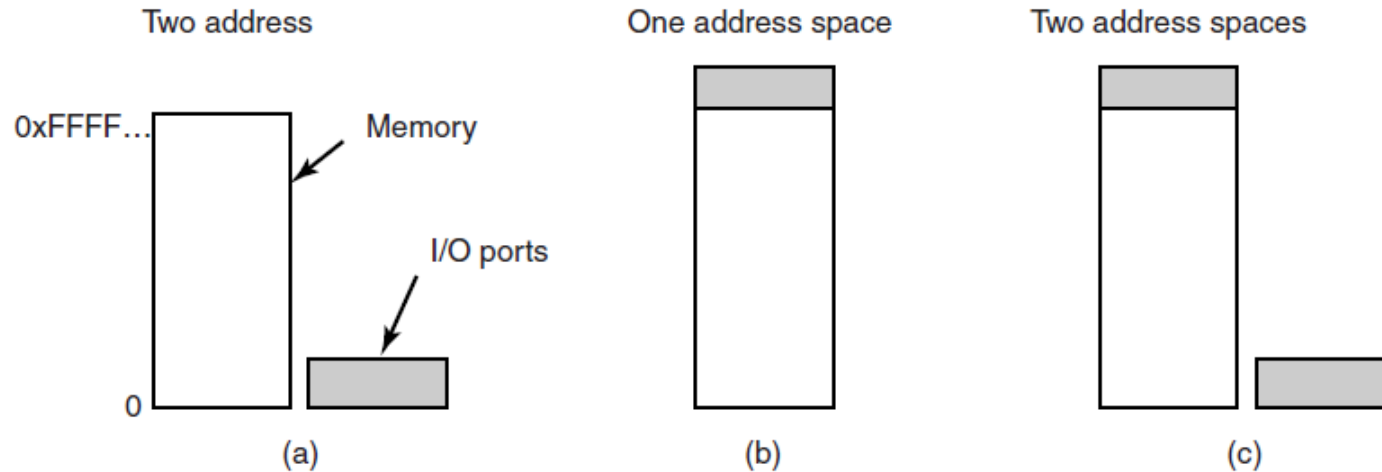
- » Two methods for CPU to communicate
  - » Each control register is assigned an I/O port number ( 8 or 16 bit number )
  - » Set of all ports is the I/O port space.
  - » Ordinary user programs can not access
  - » IN REG PORT or OUT PORT REG

# Device Controllers



- » Other method maps all control registers into memory space. ( Memory Mapped I/O )
- » Each controller is assigned a unique memory address to which no memory is assigned.
- » Usually near the top of address space.

# Device Controllers



- » Hybrid scheme (x86) - Memory mapped I/O data buffers and separate I/O ports for control registers.
- » 640k to 1M-1 reserved for device data buffers. I/O ports 0 to 64K-1.

# Device Controllers

- » Different strengths and weaknesses
  - » Memory mapped I/O
    - » If special I/O instructions needed then need to access the control registers with assembly
    - » Memory mapped allow devices drivers to be written in C.
    - » No memory protection needed to keep user processes from performing I/O.
    - » O/S just doesn't map I/O address range in process address space

# Device Controllers

- » Different strengths and weaknesses
  - » Memory mapped I/O
    - » Multiple devices can be given their own page of memory. Access can be given to users over some devices and not others
  - » Different address spaces keep drivers from conflicting with each other.

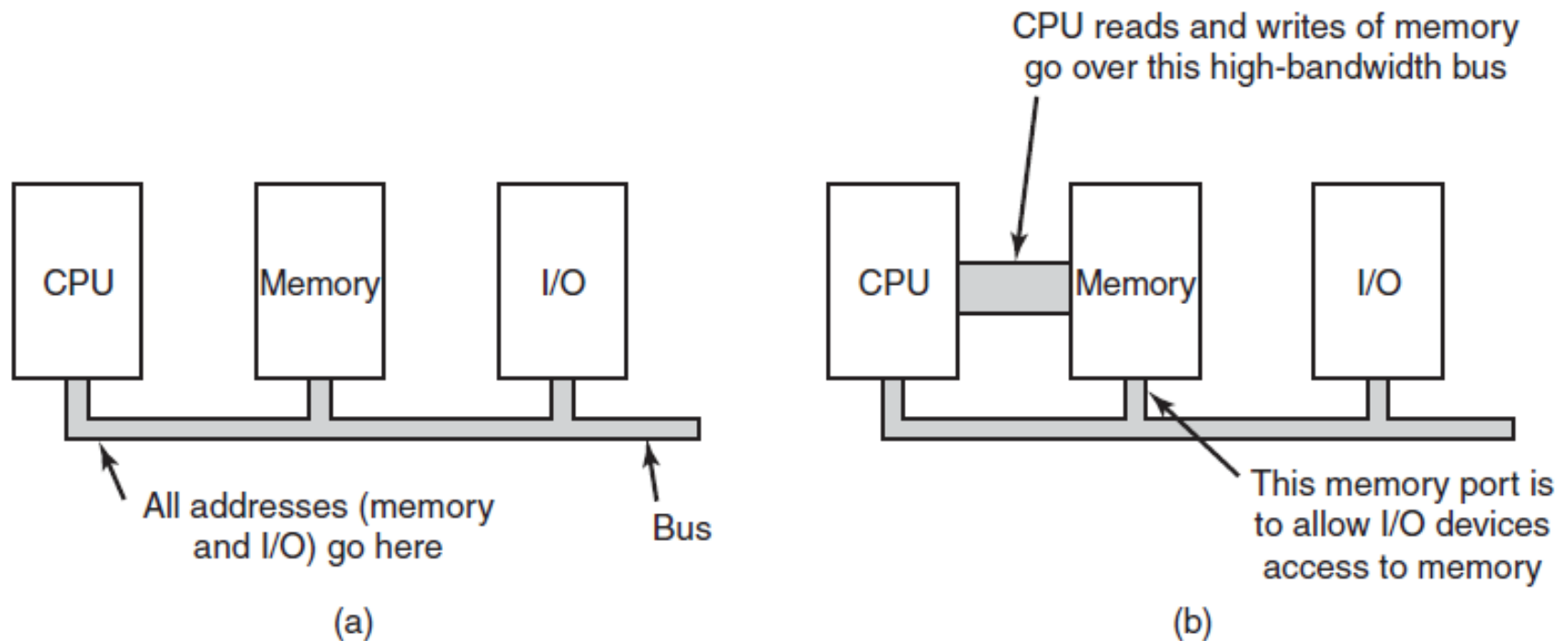


# Device Controllers

- » Different strengths and weaknesses
  - » Memory mapped I/O weaknesses
    - » Caching a control register would be disastrous.
    - » Hardware has to disable caching on a per page basis.
    - » If there is only one address space then all memory modules and all I/O devices must examine all memory references to see which to respond to.

# Device Controllers

- » Different strengths and weaknesses
- » Memory mapped I/O weaknesses

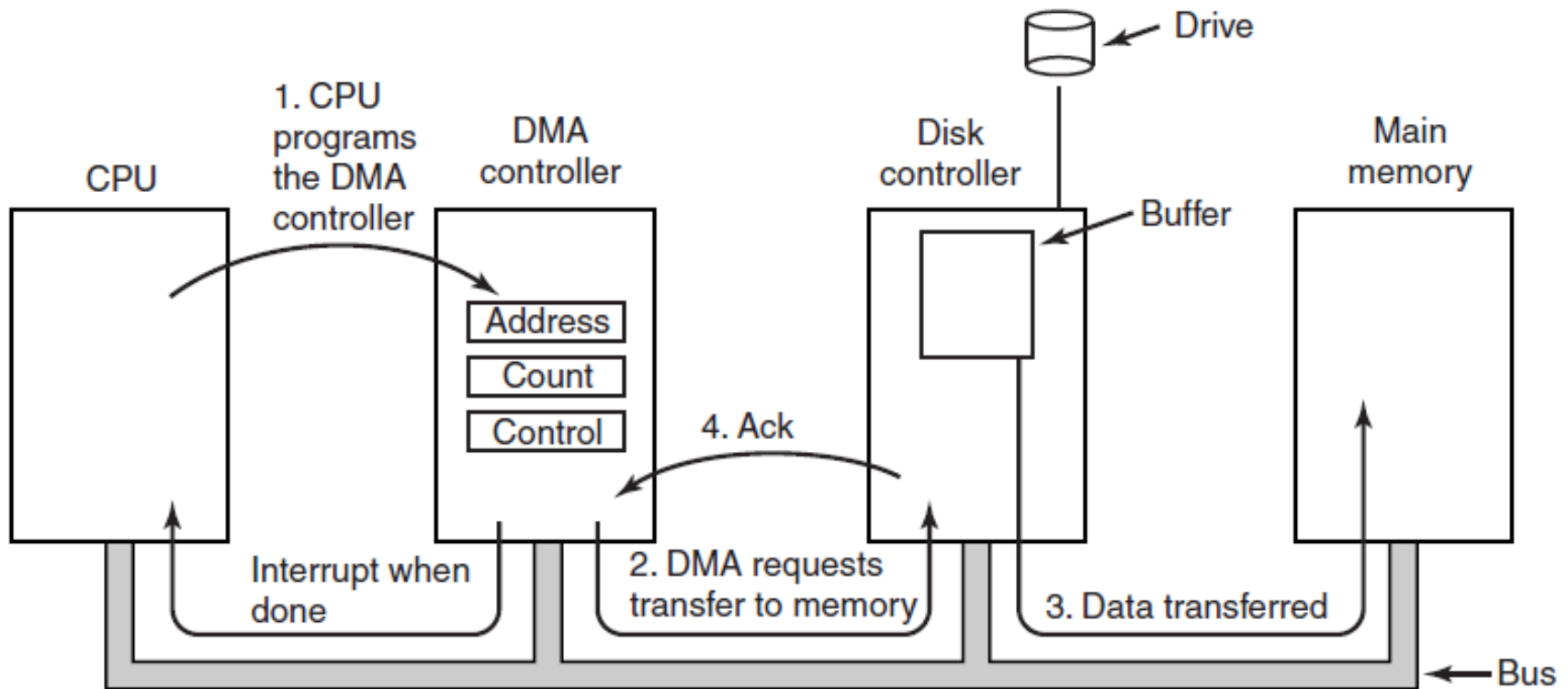


- (a) A single-bus architecture.  
(b) A dual-bus memory architecture.

# Direct Memory Access

- » CPU can request data from the controller one byte at a time, but wastes the CPU's time.
- » Direct Memory Access (DMA)
  - » Must have a DMA controller
  - » DMA controller has access to the bus independent of the CPU.
  - » CPU can read/write to the DMA controller control registers.
    - » Address, Count, Control

# Direct Memory Access



# Direct Memory Access

- » Normal disk read:
  - » Disk controlled read the block from the driver bit by bit until the entire block is in the controller's internal buffer
  - » Checksum computed
  - » Controller causes interrupt.
  - » O/S reads the disk block from the controller's buffer one byte or a word at a time.
  - » Loop until done

# Direct Memory Access

- » DMA:
  - » CPU programs the DMA controller (Step 1)
  - » CPU commands the disk controller to read the data from disk and verify the checksum
  - » DMA controller initiates the transfer by issuing a read request over the bus to the disk controller.
  - » Disk controller doesn't know if it's the CPU or DMA controller issuing it.
  - » Loop until done

# Direct Memory Access

- » DMA:
  - » More sophisticated DMA can handle multiple transfers at a time.
  - » Round robin the requests or priority requests
  - » System busses can operate in two modes: word-at-a-time and block more.
  - » Word-at-a-time mode allows the device controller to sneak in and steal an occasional bus cycle from the CPU
    - » Cycle Stealing

# Direct Memory Access

- » DMA:
  - » Block mode - the DMA controller tells the device to acquire the bus, issue the transfers then release the bus.
  - » Burst mode
  - » More efficient than cycle stealing
  - » Blocks the CPU and other devices



# Direct Memory Access

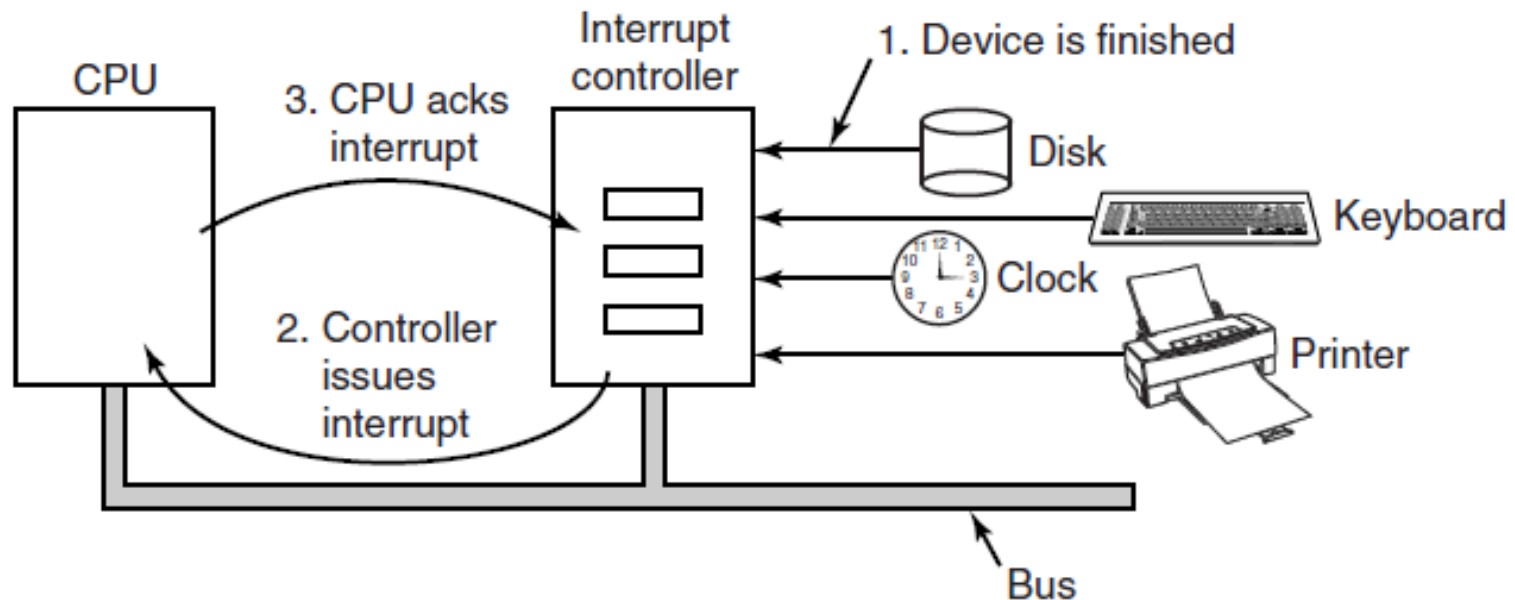
- » DMA:
  - » Fly-by mode - DMA tells the device controller to write directly to main memory
  - » Alternate mode allows the device to send to the DMA controller which then issues a second bus request to write to a destination
    - » Allows device to device DMA
- » Some computers don't use DMA
  - » Main CPU is often faster than DMA controller
  - » Idling CPU waiting for slower DMA is pointless
  - » Less hardware means cheaper cost

# Interrupts Revisited

- » Refresher:
  - » When an I/O device is done with its work it issues an interrupt.
    - » Asserts a signal on a bus line that it has been assigned.
    - » Interrupt controlled chip on motherboard decides what to do
      - » If no interrupt pending, the interrupt controller handles the interrupt immediately.
      - » If other interrupts in progress it is ignored until the interrupt controller is free.
  - » Interrupt handler puts a number on the address line specifying which device needs attention and asserts a signal to interrupt the CPU.

# Interrupts Revisited

- » Refresher:
  - » The number on the address line is used as an index into the interrupt vector table that allows the CPU to fetch a new program counter.



# Interrupts Revisited

- » When executing the interrupt, what does the CPU save?
  - » Program counter is bare minimum
  - » All visible registers and most internal registers at the other end
- » Where do you put the data?
  - » Internal registers, but then can't acknowledge interrupt handled until all registers read back out
    - » Takes time, leaves dead space
  - » Stack
    - » Can't be user stack. Stack pointer may not be legal
      - » Might be on the end of a page. If you page fault where do you save the state to handle the fault?

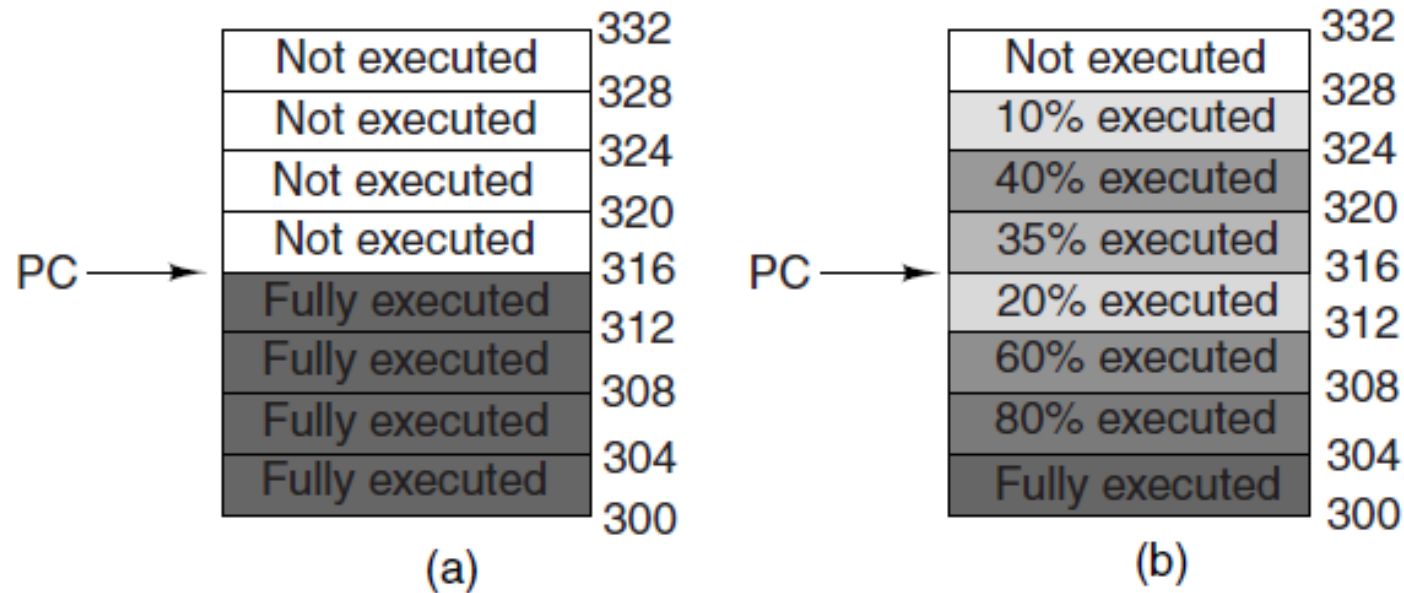
# Interrupts Revisited

- » Kernel stack?
  - » May require switching into kernel mode
    - » Change MMU context
      - » Invalidate the cache and TLB

# Interrupts Revisited

- » Kernel stack?
  - » May require switching into kernel mode
    - » Change MMU context
      - » Invalidate the cache and TLB
- » Modern CPUs are heavily pipelined, often superscalar.
  - » Can't assume if an interrupt occurs after an instruction that all instructions leading up to and including that instruction have been executed completely.
    - » Many partially executed instructions.

# Interrupts Revisited



(a) A precise interrupt. (b) An imprecise interrupt.

# Interrupts Revisited

- » Precise Interrupt - Interrupt that leaves the machine in a well-defined state.
- » Four properties of a precise interrupt:
  1. The PC saved in a known place.
  2. All instructions before that pointed to by PC have fully executed.
  3. No instruction beyond that pointed to by PC has been executed.
  4. Execution state of instruction pointed to by PC is known.



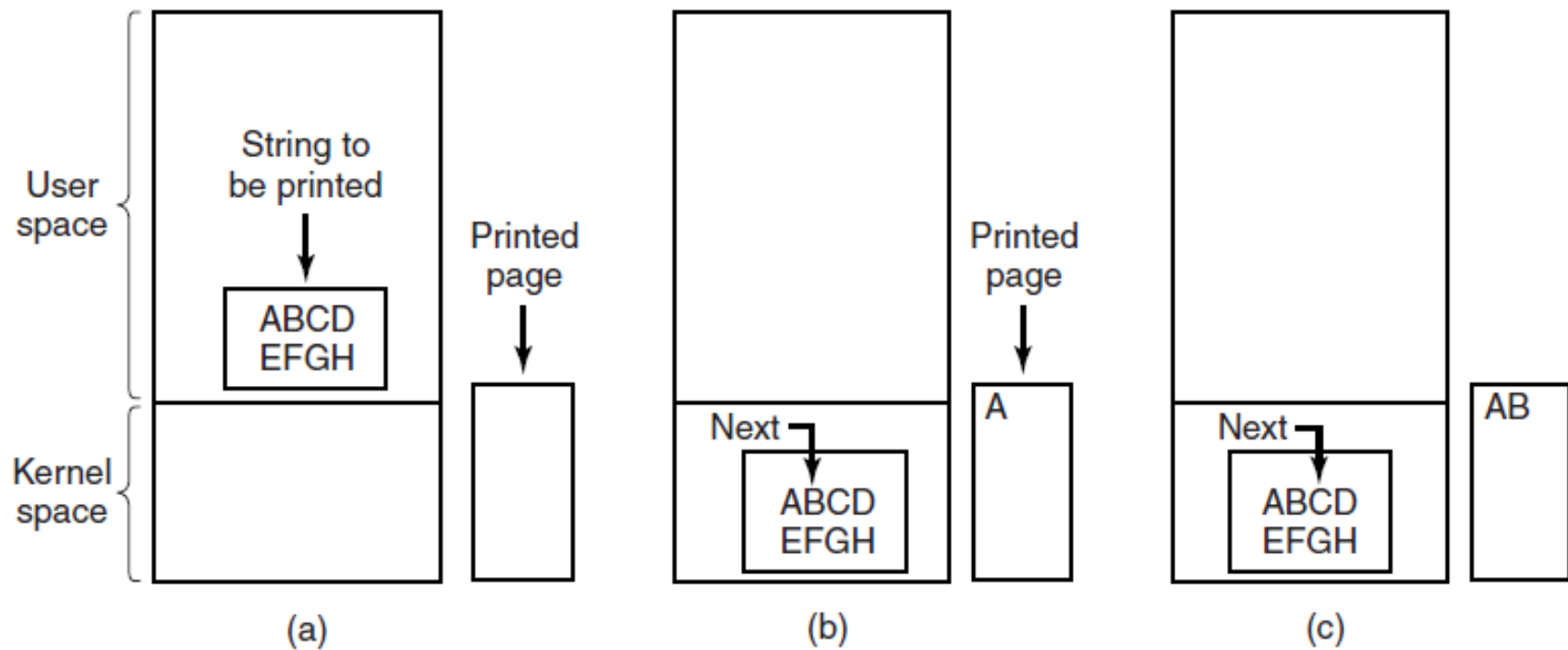
# Interrupts Revisited

- » Imprecise Interrupt - Interrupt that does not meet the four requirements
  - » Unpleasant for OS developer.
  - » Machines with imprecise interrupts usually vomit a large amount of internal state onto the stack to give the operating system the possibility of figuring out what was going on.
    - » Large amount of work to restart
    - » Can make very fast superscalar CPUs unsuitable for real-time work due to slow interrupts
- » x86 has very complex logic in the CPU to have precise interrupts

# Programmed I/O

- » Three fundamentally different ways I/O can be performed.
  - » Programmed I/O
  - » Interrupt Driven I/O
  - » I/O using DMA
- » Simplest is programmed I/O, also known as letting the CPU do all the work.

# Programmed I/O



Steps in printing a string.

# Programmed I/O

- » Essential aspect of programmed I/O:
  - » The CPU continuously polls the device to see if it is ready to accept more data
  - » Busy waiting / Polling

```
copy_from_user(buffer, p, count);  
for (i = 0; i < count; i++) {  
    while (*printer_status_reg != READY) ;  
    *printer_data_register = p[i];  
}  
return_to_user();
```

/\* p is the kernel buffer \*/  
/\* loop on every character \*/  
/\* loop until ready \*/  
/\* output one character \*/

# Programmed I/O

- » Essential aspect of programmed I/O:
  - » Busy waiting is inefficient.
  - » Ok for small embedded devices where the CPU has nothing else to do, but larger systems need a better solution

# Programmed I/O

## » Interrupt-driven I/O

```
copy_from_user(buffer, p, count);  
enable_interrupts();  
while (*printer_status_reg != READY) ;  
*printer_data_register = p[0];  
scheduler();
```

(a)

```
if (count == 0) {  
    unblock_user();  
} else {  
    *printer_data_register = p[i];  
    count = count - 1;  
    i = i + 1;  
}  
acknowledge_interrupt();  
return_from_interrupt();
```

(b)

Writing a string to the printer using interrupt-driven I/O.

(a) Code executed at the time the print system call is made. (b) Interrupt service procedure for the printer.

# Programmed I/O

- » Problem with interrupt driven I/O:
  - » Interrupt every character
- » Use DMA for I/O

```
copy_from_user(buffer, p, count);  
set_up_DMA_controller();  
scheduler();
```

(a)

```
acknowledge_interrupt();  
unblock_user();  
return_from_interrupt();
```

(b)

Printing a string using DMA. (a) Code executed when the print system call is made. (b) Interrupt service procedure.

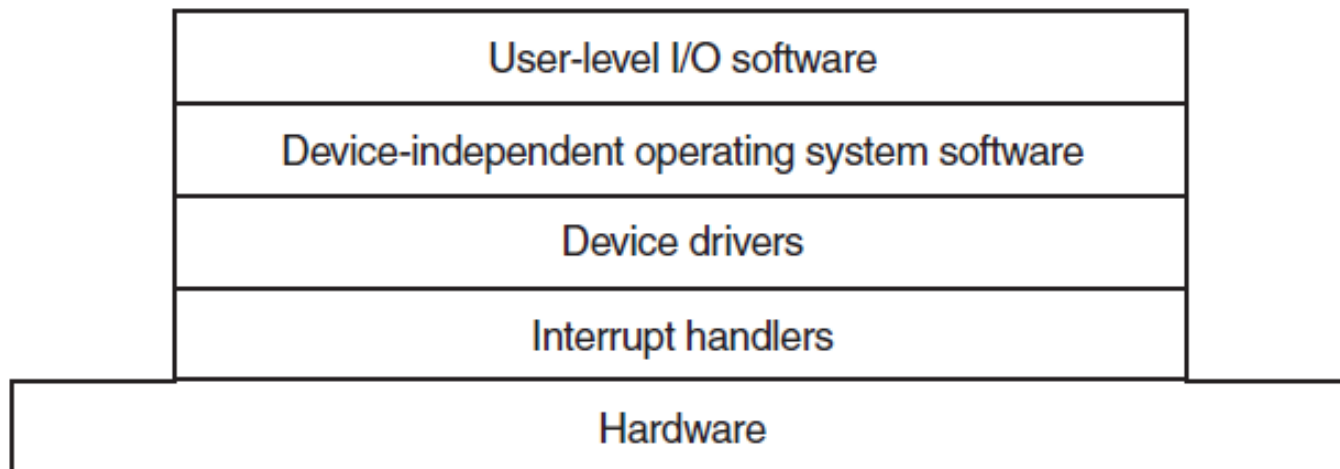
# Programmed I/O

- » Big win for DMA is reducing the number of interrupts from one per character to one per buffer.
- » DMA controller is slower than CPU.
- » If DMA can't drive the device at full speed then CPU driven I/O or interrupt driven I/O may be better.



# Programmed I/O

» I/O software typically organized in four layers.



# Device Drivers

- » Device registers and nature of commands can vary dramatically from device to device.
- » Each I/O device attached to a computer needs some device specific code for controlling it, a device driver.
  - » Usually written by the device manufacturer
- » While wildly different many devices based on the same underlying technology
  - » USB

# Device Drivers

» USB devices are stacked.

1. Link Layer (Serial IO) - deals with differential line transitions and signaling, and decode the stream to binary data, very often in hardware
2. USB Packet Layer - deals with structure of USB data packets
3. USB Required Functionality - enumeration, buffers, endpoints
4. USB higher level APIs - Audio, HID, etc, that have their own restrictions and needs.

# Device Drivers

- » Generally device drivers must be part of the kernel
- » Can write user space device drivers, e.g. MINIX
  - » Most run in kernel space
- » Need a well defined mode of what a driver does and how it interacts with the system so third party developers can write drivers.
  - » Block devices
  - » Character devices
- » Standard interface all block drivers must support and a second standard interface for all character devices.

# Device Drivers

- » For many years the norm was to compile the drivers into the kernel in one single binary program.
- » Advent of personal computers and myriad of devices changed.
  - » Home users don't want to recompile their kernel.

# Principle of I/O software

» Key concepts:

- » Device independence - Write a program that can access any I/O without having to specify the device in advance
- » open should work on any device: drive, flash, cd
- » Uniform Naming - should not depend on device in any way.

# Principle of I/O software

- » Key concepts:
  - » Error handling: errors should be handled as close to the hardware as possible.
  - » If the controller discovers a read error it should correct it .
- » Synchronous / Asynchronous
  - » Most physical I/O is asynchronous
  - » Blocking programs easier to write

# Principle of I/O software

- » Key concepts:
  - » Buffering: Usually can't store data directly off the device into its final destination
    - » Sometime must decouple rate at which the data is filled from the rate at which it is emptied.
  - » Shareable / Dedicated devices
    - » Drives can have multiple users sharing
    - » Printers can only have one user at a time



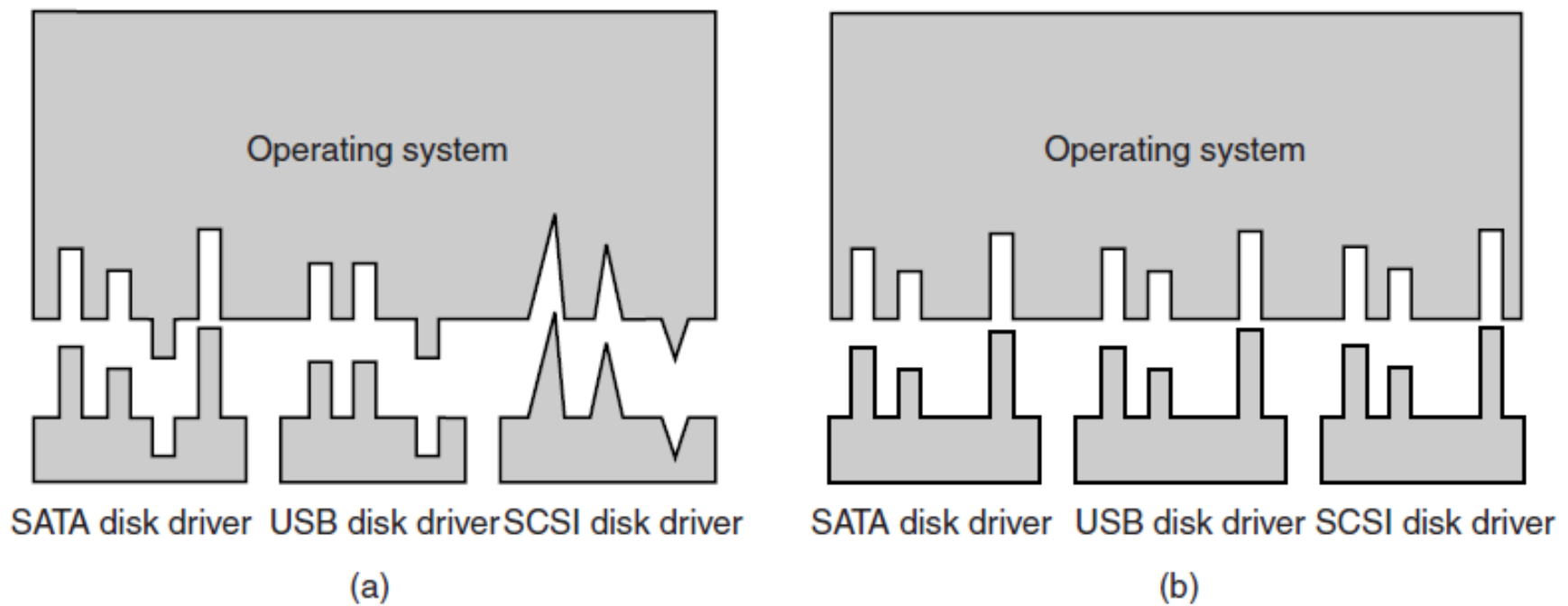
# Device Independent I/O Software

- » While some I/O is device specific, other parts are device independent.
- » Lines varies
  - » Efficiencies

Uniform interfacing for device drivers
Buffering
Error reporting
Allocating and releasing dedicated devices
Providing a device-independent block size

# Uniform Interfacing

» How do we make all the drivers and I/O devices look the same?



(a) Without a standard driver interface.

(b) With a standard driver interface.

# Uniform Interfacing

- » Device driver writers know what is expected in their drivers
- » While all I/O devices are not identical usually there are a small number of device types
- » Each class defines a set of operations a driver must supply

# Linux Device Drivers

```
# ls -l /dev/hda[1-3]  
brw-rw---- 1 root disk 3, 1 Jul 5 2000 /dev/hda1  
brw-rw---- 1 root disk 3, 2 Jul 5 2000 /dev/hda2  
brw-rw---- 1 root disk 3, 3 Jul 5 2000 /dev/hda3
```

- » Major and Minor Numbers
  - » Major number tells you which driver is used to access the hardware.
    - » Each driver is assigned a unique major number; all device files with the same major number are controlled by the same driver.
  - » The minor number is used by the driver to distinguish between the various hardware it controls.

# Linux Device Drivers

- » Adding a driver to your system means registering it with the kernel.
- » Assigning it a major number during the module's initialization. You do this by using the `register_chrdev` function, defined by `linux/fs.h`.

# Linux Device Drivers

- » How do you get a major number without hijacking one that's already in use?
- » Look through Documentation/devices.txt and pick an unused one.
- » Ask the kernel to assign you a dynamic major number.

# Linux Device Drivers

- » How do you get a major number without hijacking one that's already in use?
- » Look through Documentation/devices.txt and pick an unused one.
- » Ask the kernel to assign you a dynamic major number.

# Linux Device Drivers

- » Each device is represented in the kernel by a file structure, which is defined in `linux/fs.h`.
- » The “file” is a kernel level structure and never appears in a user space program.
- » It's not the same thing as a `FILE`, which is defined by `glibc` and would never appear in a kernel space function.
- » Also, its name is a bit misleading; it represents an abstract open ‘file’, not a file on a disk, which is represented by a structure named `inode`.
- » VFS





# Linux Device Drivers

```
struct file_operations fops = {  
    read: device_read,  
    write: device_write,  
    open: device_open,  
    release: device_release  
};
```

- » Populate the `file_operations` struct with the functions your device will support.

# Linux Device Drivers

```
/*
 * This function is called when the module is loaded
 */
int init_module(void)
{
    Major = register_chrdev(0, DEVICE_NAME, &fops);

    if (Major < 0) {
        printk(KERN_ALERT "Registering char device failed with %d\n", Major);
        return Major;
    }

    printk(KERN_INFO "I was assigned major number %d. To talk to\n", Major);
    printk(KERN_INFO "the driver, create a dev file with\n");
    printk(KERN_INFO "'mknod /dev/%s c %d 0'.\n", DEVICE_NAME, Major);
    printk(KERN_INFO "Try various minor numbers. Try to cat and echo to\n");
    printk(KERN_INFO "the device file.\n");
    printk(KERN_INFO "Remove the device file and module when done.\n");

    return SUCCESS;
}
```

» Create `init__module` function to register the device.

# Linux Device Drivers

```
/*
 * Called when a process tries to open the device file, like
 * "cat /dev/mycharfile"
 */
static int device_open(struct inode *inode, struct file *file)
{
    static int counter = 0;

    if (Device_Open)
        return -EBUSY;

    Device_Open++;
    sprintf(msg, "I already told you %d times Hello world!\n", counter++);
    msg_Ptr = msg;
    try_module_get(THIS_MODULE);

    return SUCCESS;
}
```

» For each function you've registered in the `file_operations` structure define your function.

# Linux Device Drivers

```
/*
 * Called when a process tries to open the device file, like
 * "cat /dev/mycharfile"
 */
static int device_open(struct inode *inode, struct file *file)
{
    static int counter = 0;

    if (Device_Open)
        return -EBUSY;

    Device_Open++;
    sprintf(msg, "I already told you %d times Hello world!\n", counter++);
    msg_Ptr = msg;
    try_module_get(THIS_MODULE);

    return SUCCESS;
}
```

» For each function you've registered in the `file_operations` structure define your function.

# Linux Device Drivers

```
/*
 * Called when a process, which already opened the dev file, attempts to
 * read from it.
 */
static ssize_t device_read(struct file *filp, /* see include/linux/fs.h */
                           char *buffer,    /* buffer to fill with data */
                           size_t length,    /* length of the buffer */
                           loff_t * offset)
{
    /*
     * Number of bytes actually written to the buffer
     */
    int bytes_read = 0;

    /*
     * If we're at the end of the message,
     * return 0 signifying end of file
     */
    if (*msg_Ptr == 0)
        return 0;

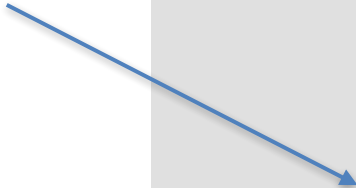
    /*
     * Actually put the data into the buffer
     */
    while (length && *msg_Ptr) {

        /*
         * The buffer is in the user data segment, not the kernel
         * segment so "*" assignment won't work. We have to use
         * put_user which copies data from the kernel data segment to
         * the user data segment.
         */
        put_user(*(msg_Ptr++), buffer++);

        length--;
        bytes_read++;
    }

    /*
     * Most read functions return the number of bytes put into the buffer
     */
    return bytes_read;
}
```

Copy from  
kernel to user

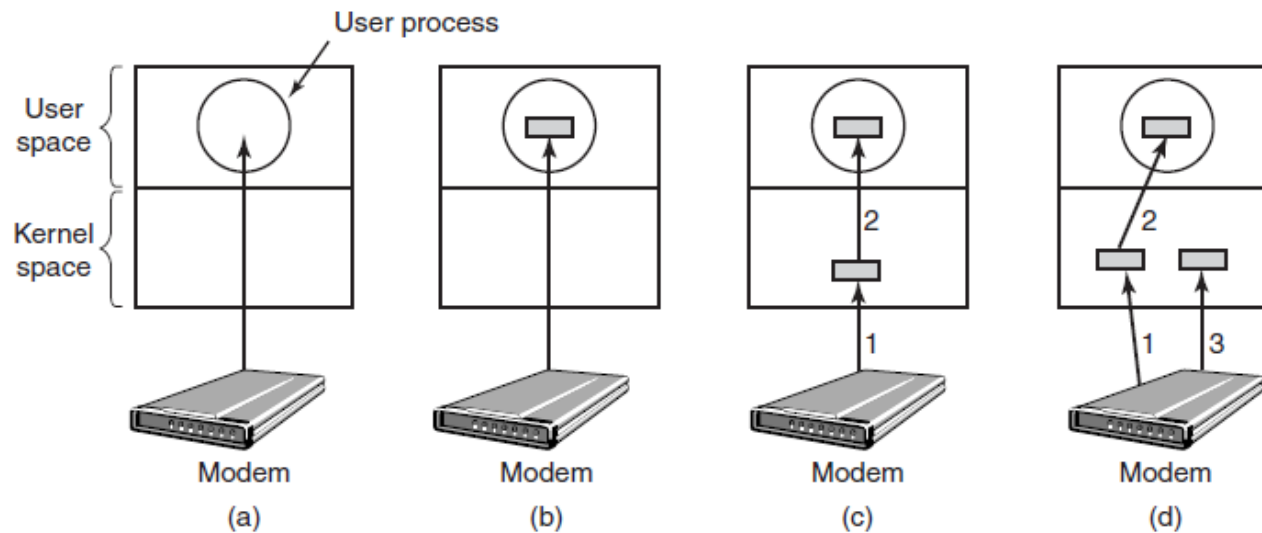


# Uniform Interfacing

- » Uniform naming.
  - » /dev/sda
  - » /dev/hda
- » Uniform protections and permissions

# Buffering

» Both block and character devices must deal with buffering.



(a) Unbuffered input. (b) Buffering in user space. (c) Buffering in the kernel followed by copying to user space. (d) Double buffering in the kernel.



# Buffering

- » (a) Unbuffered input.
  - » read, block, interrupt, repeat
  - » inefficient
- » (b) Buffering in user space.
  - » Process provides n-byte buffer
  - » When buffer full the process is awakened
    - » What if buffer is paged out?

# Buffering

- » (c) Buffering in the kernel followed by copying to user space.
  - » Copy in one operation
- » (d) Double buffering in the kernel.
  - » Copy one while accumulating in another
- » Circular buffer

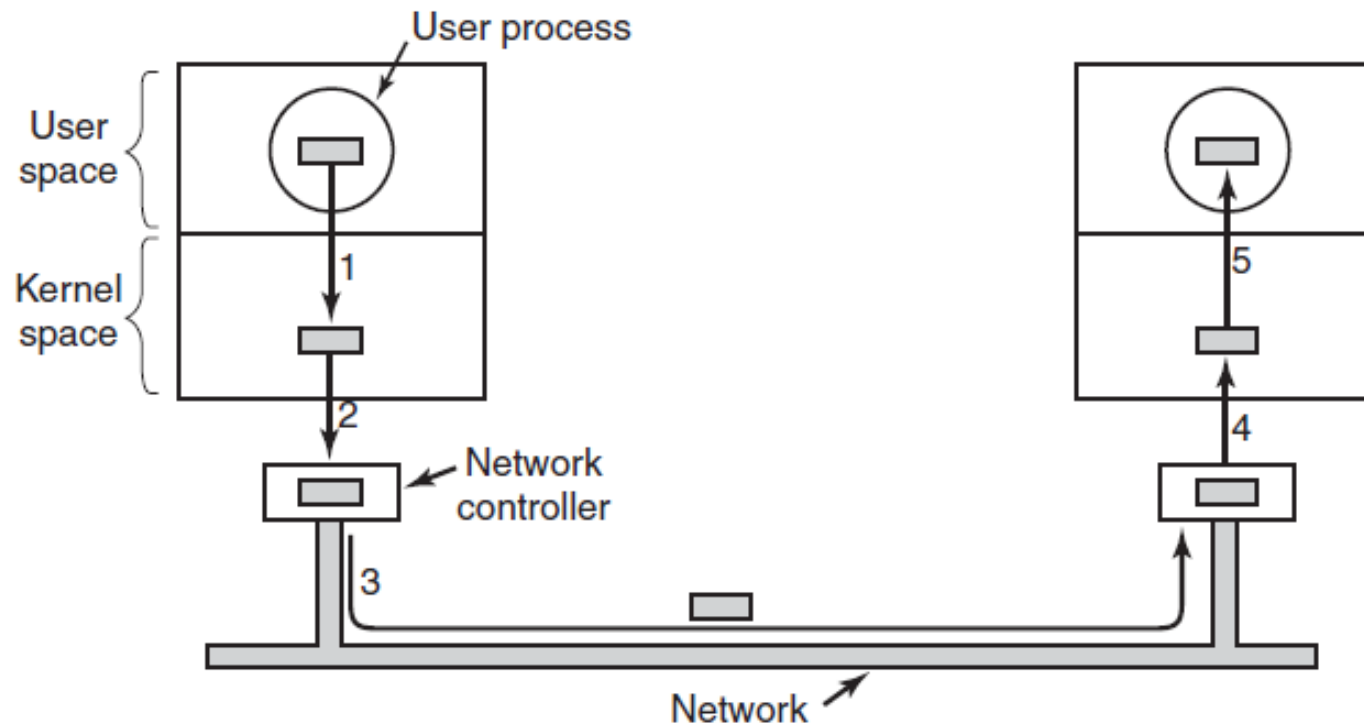
# Buffering

- » Also important on output
- » On write:
  - » Block the user until all characters written.
  - » Slow
  - » Release the user and process I/O while user process continues
  - » How do you know when it's finished?

# Buffering

- » System could generate signal or interrupts
  - » Prone to race conditions
- » Better to copy to kernel buffer and in (c)
- » Too many buffers on the path cause performance problems.

# Buffering



# Error Reporting

- » Errors are far more common in I/O context than in any other.
- » Device specific and handled by driver
- » Framework for error handling is device independent

# Classes of Errors

## 1. Programmer error

- Writing to an input
- Just report error

## 2. I/O error

- Writing to bad disk block
- Driver needs to determine what to do
- May pass error up to device independent software

# Dedicated Devices

- » Some device, such as printers, can only be used by a single process at a time
  - Operating system adjudicates requests
    - Require processes to call `open()`
  - Special mechanism for requesting and releasing devices
  - Processes block and wait in queue



# Device Independent Block Size

- » Different disks may have different block sizes.
  - Device independent software hides this and provides a uniform block size to higher layers.
  - Treat several sectors as a single logical block.

# RAID Structure

- Reliability measured in mean time between failure
- RAID - redundant array of inexpensive disks
- RAID – multiple disk drives provides reliability via redundancy
- Frequently combined with NVRAM to improve write performance
- RAID is arranged into six different levels

# RAID (Cont.)

- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively
- Disk **striping** uses a group of disks as one storage unit
- RAID schemes improve performance and improve the reliability of the storage system by storing redundant data
  - **Mirroring** or **shadowing** (**RAID 1**) keeps duplicate of each disk
  - Striped mirrors (**RAID 1+0**) or mirrored stripes (**RAID 0+1**) provides high performance and high reliability
  - **Block interleaved parity** (**RAID 4, 5, 6**) uses much less redundancy

# RAID (Cont.)

- RAID within a storage array can still fail if the array fails, so automatic **replication** of the data between arrays is common
- Frequently, a small number of **hot-spare** disks are left unallocated, automatically replacing a failed disk and having data rebuilt onto them

# RAID Levels



(a) RAID 0: non-redundant striping.



(b) RAID 1: mirrored disks.



(c) RAID 2: memory-style error-correcting codes.



(d) RAID 3: bit-interleaved parity.



(e) RAID 4: block-interleaved parity.



(f) RAID 5: block-interleaved distributed parity.



(g) RAID 6: P + Q redundancy.

- RAID 0 - **Striped disk array** without parity.
- Data spread across multiple disk drives but no data redundancy
- Improves performance because multiple reads and writes can be carried out at the same time.
- Works best with large requests.
- Works worst with OS that ask for data one sector at a time.
- No parallelism

## RAID 0

# RAID 0

- RAID 0 - **Striped disk array** without parity.
- Does not increase fault tolerance.
- Actually decreases reliability since if one drive fails the all data is lost since the OS treats the array as a single drive.
- N drives means configuration is N times as likely to fail.

# RAID 1

- RAID 1 - [Mirroring](#)
  - Duplicate set of drives
  - When data is written to one it's also written to its duplicate.
  - When one fails, swap in new and copy the data over.



# RAID 1

- RAID 1 - [Mirroring](#)
  - Assuming drive and its mirror can be read at the same time it provides twice the read transaction rate
  - Write transaction is unaffected.
    - No performance gain
  - Most expensive raid configuration

# RAID 2 and RAID 3

- RAID 2 - Error-correcting coding
- RAID 3 - Bit interleaved parity
- Prohibitively expensive and inferior to other RAID levels.
- Drive spinning must be synchronized

# RAID 4

- RAID 4 - Dedicated parity drive
- Block level striping like RAID 0 with a parity disk
- If a data disk fails the parity data is used to create a replacement disk.
- Every time a block is written the parity block must also be read, recalculated and rewritten. I/O bottleneck.

# RAID 4

- Same reliability as RAID 1 but if a drive fails the performance hit is worse.
- Heavy load on the parity drive
  - Bottleneck

# RAID 5

- RAID 5 - Block interleaved distribution parity
- Like RAID 4 but instead of parity on the same drive the parity block is assigned to the drives in a round robin fashion.
- Removes excessive use of the parity drive.

# RAID 6

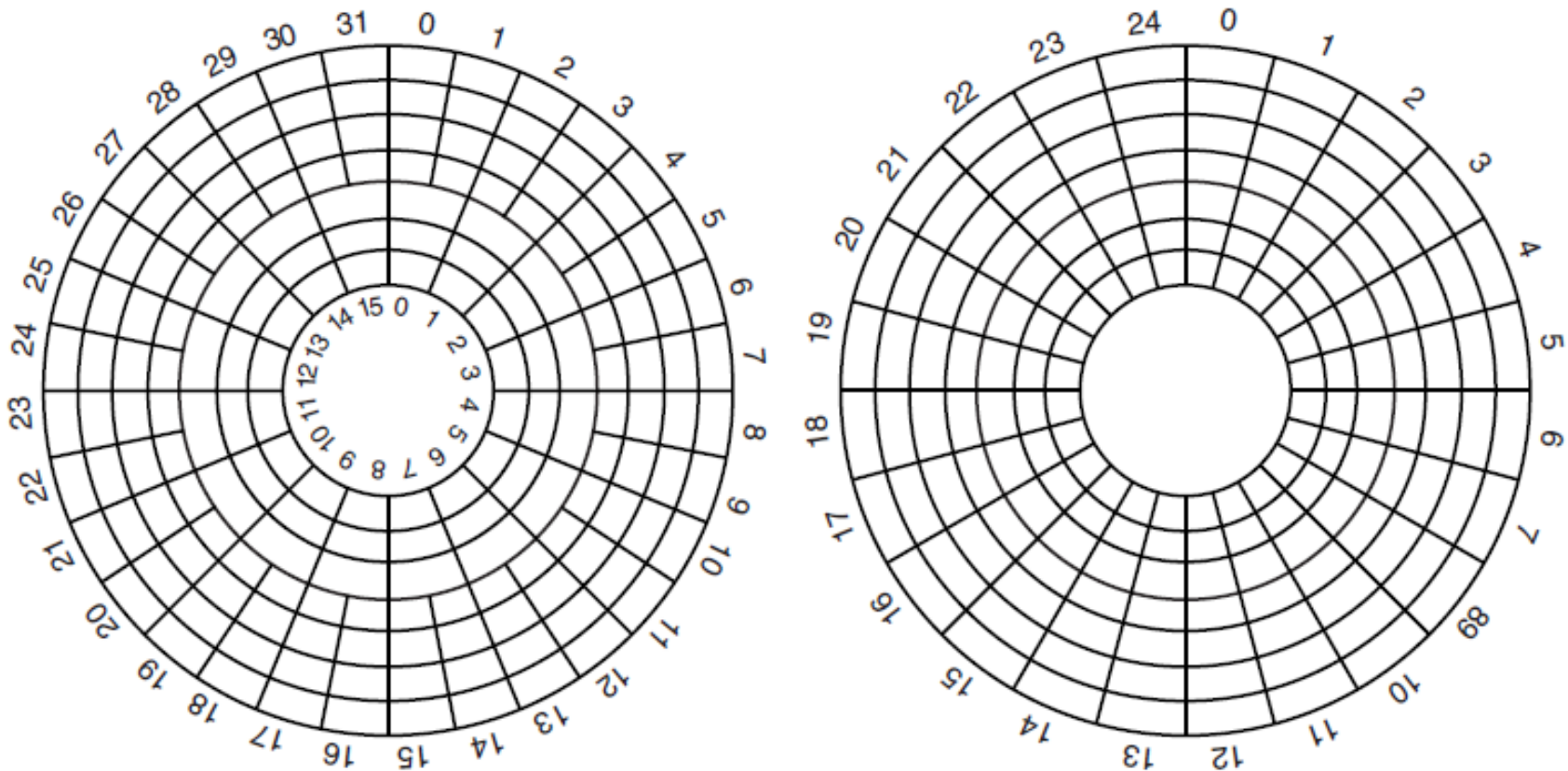
- RAID 6 - Independent data disks with double parity
- Block-level striping with parity across the disks like RAID 5 but instead of a simple parity scheme it calculates parity using two different algorithms.
- Requires an extra disk drive of RAID 5 but will tolerate the loss of two drives at the same time.

# Disks

- » Magnetic disks are organized into cylinders, each having as many tracks as there are heads stack vertically.
  - Tracks are divided into sectors
- » Early disks delivered serial bit stream
- » IDE and SATA have a microcontroller that does considerable work and allows the real controller to issue higher commands.
  - Overlapped seeks

# Disks

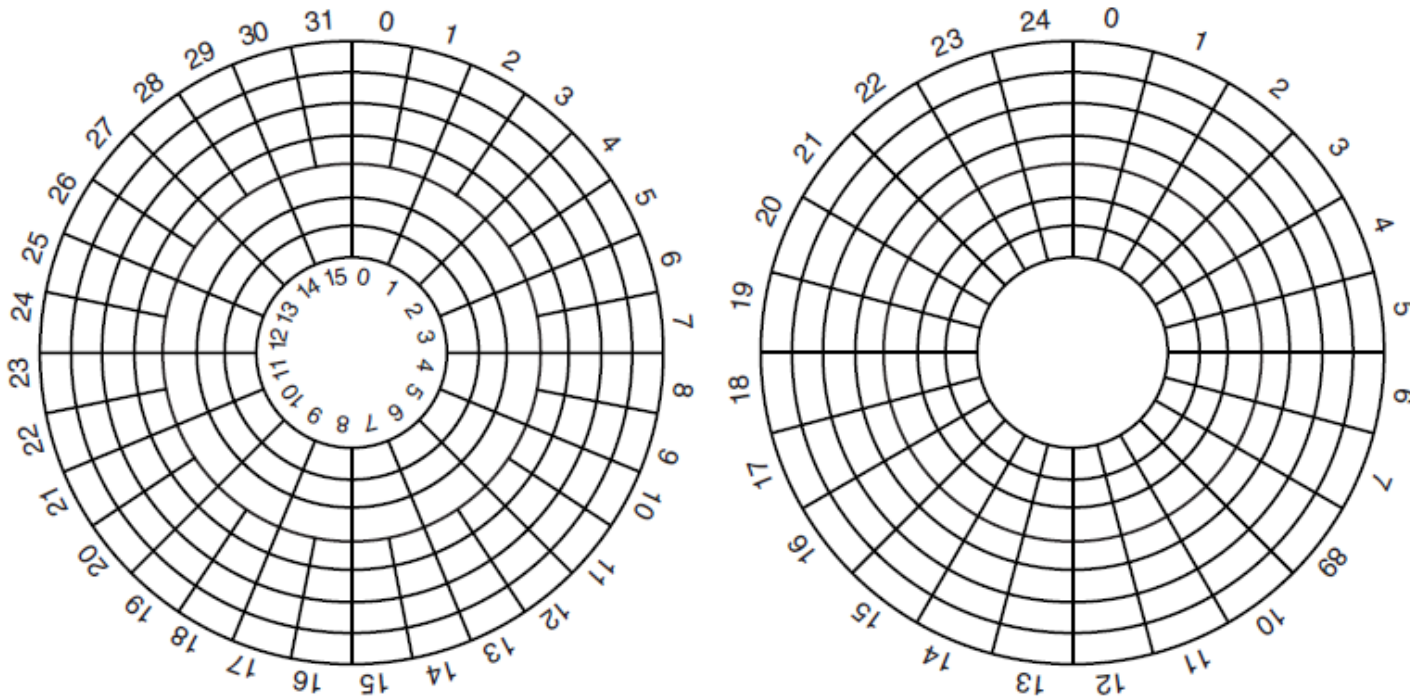
- » Geometry specified and used by the driver software is almost always different from the physical format.





# Disks

- » Logical Block Addressing - disk sectors are numbered consecutively starting at 0 with no regard for the disk geometry,



# Disk Formatting

- » Hard disk consists of a stack of aluminum or glass platters.
- » On each is deposited a thin magnetizable metal oxide.
- » No information
- » Each platter must receive a low-level format.
- » Series of sectors with short gaps between



# Disk Formatting

- » Preamble contains bit pattern hardware can recognize the sector.
- » Cylinder number
- » Sector number
- » Usually 512 -byte sectors.
- » ECC size varies by manufacturers but usually 16-bytes.

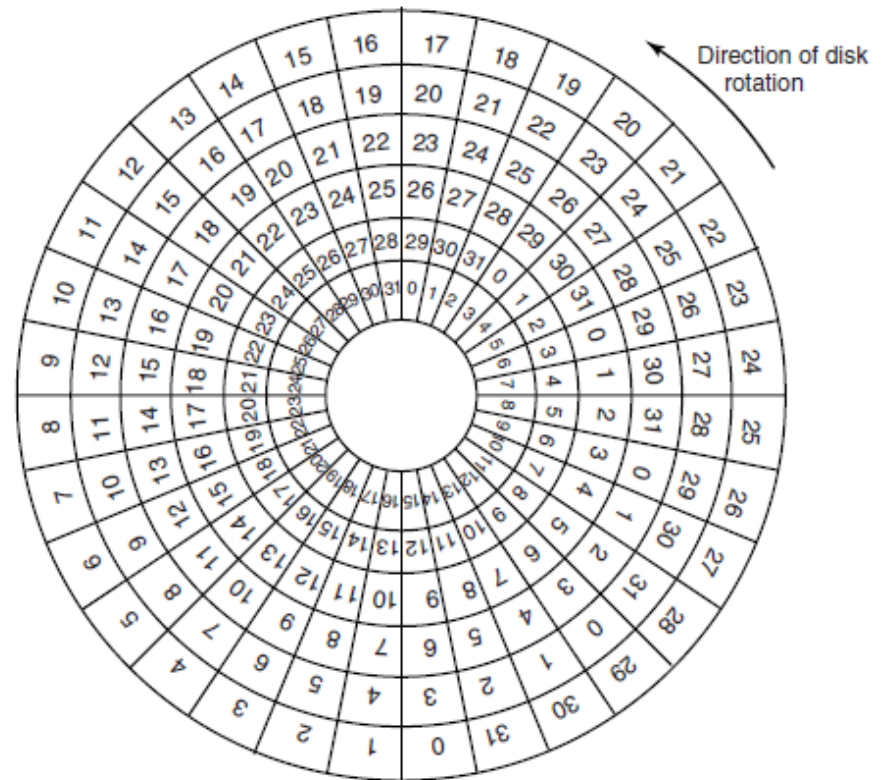
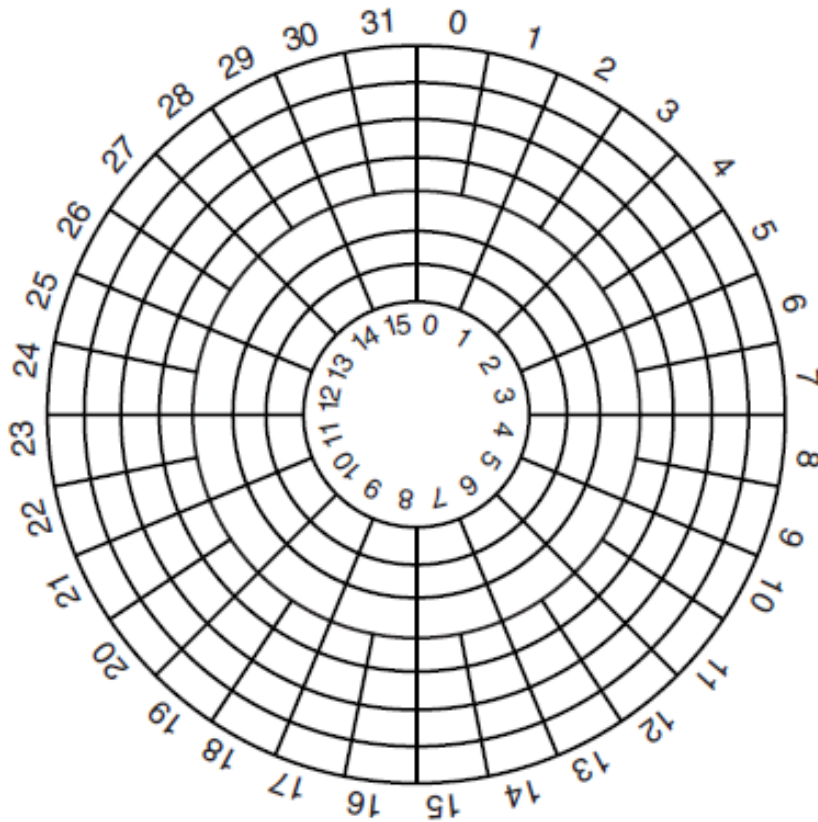
# Transfer Time

- » Consider a disk with 1MB per track, a rotation time of 8.33 ms and a seek time of 5 ms.
- » The time in ms to read a block of k bytes is the sum of the seek, rotational delay and transfer times:

$$5 + 4.165 + ( k / 1000000 ) * 8.33$$

# Disk Formatting

- » Cylinder skew - position of sector 0 is offset from previous track

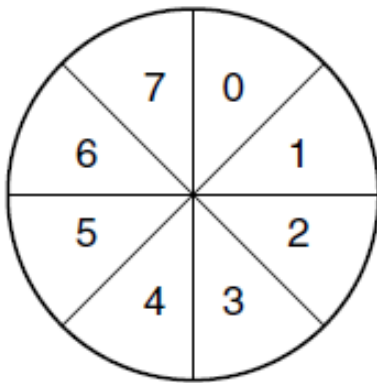


# Skew

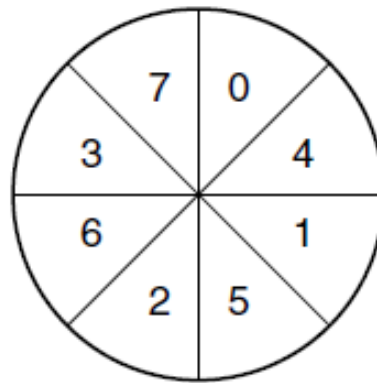
- » How much cylinder skew is needed for a 7200-RPM disk with a track-to-track seek time of 3 msec? The disk has 200 sectors of 512 bytes each on each track.

# Disk Formatting

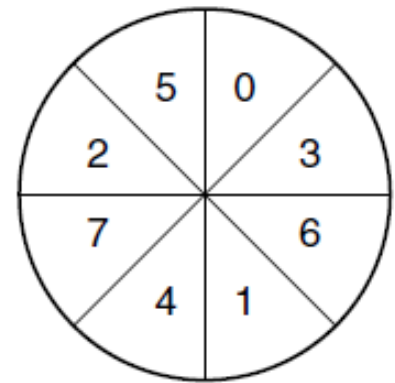
- » Transferring from controller to memory also causes delays that must be accounted for:
  - Single interleaving
  - Double interleaving



(a)



(b)



(c)

# Disk Scheduling

- ❑ The operating system is responsible for using hardware efficiently — for the disk drives, this means having a fast access time and disk bandwidth
- ❑ Access time has two major components
  - **Seek time** is the time for the disk are to move the heads to the cylinder containing the desired sector
  - **Rotational latency** is the additional time waiting for the disk to rotate the desired sector to the disk head
- ❑ Minimize seek time
- ❑ Seek time  $\approx$  seek distance
- ❑ Disk **bandwidth** is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer



# Disk Scheduling (Cont.)

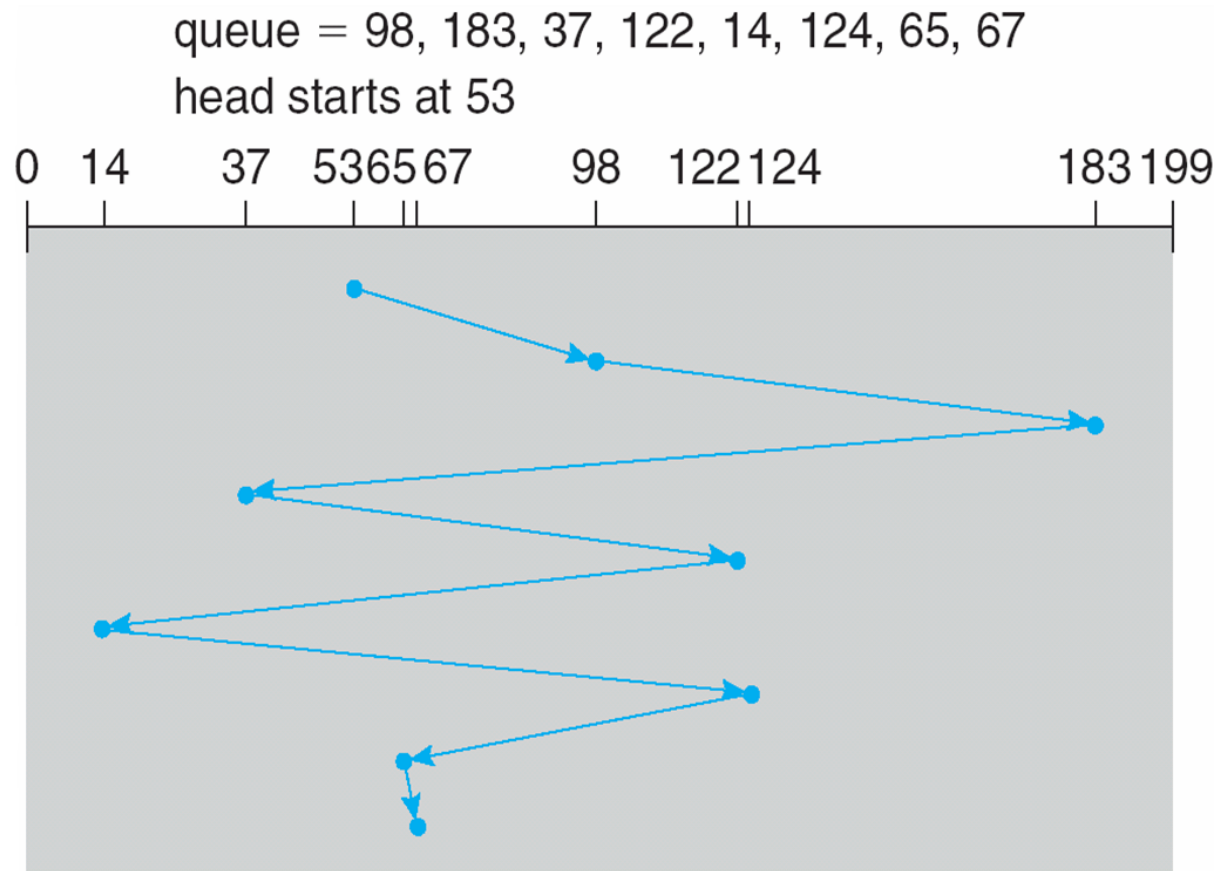
- ☐ Several algorithms exist to schedule the servicing of disk I/O requests
- ☐ We illustrate them with a request queue (0-199)

98, 183, 37, 122, 14, 124, 65, 67

Head pointer 53

# FCFS

Illustration shows total head movement of 640 cylinders



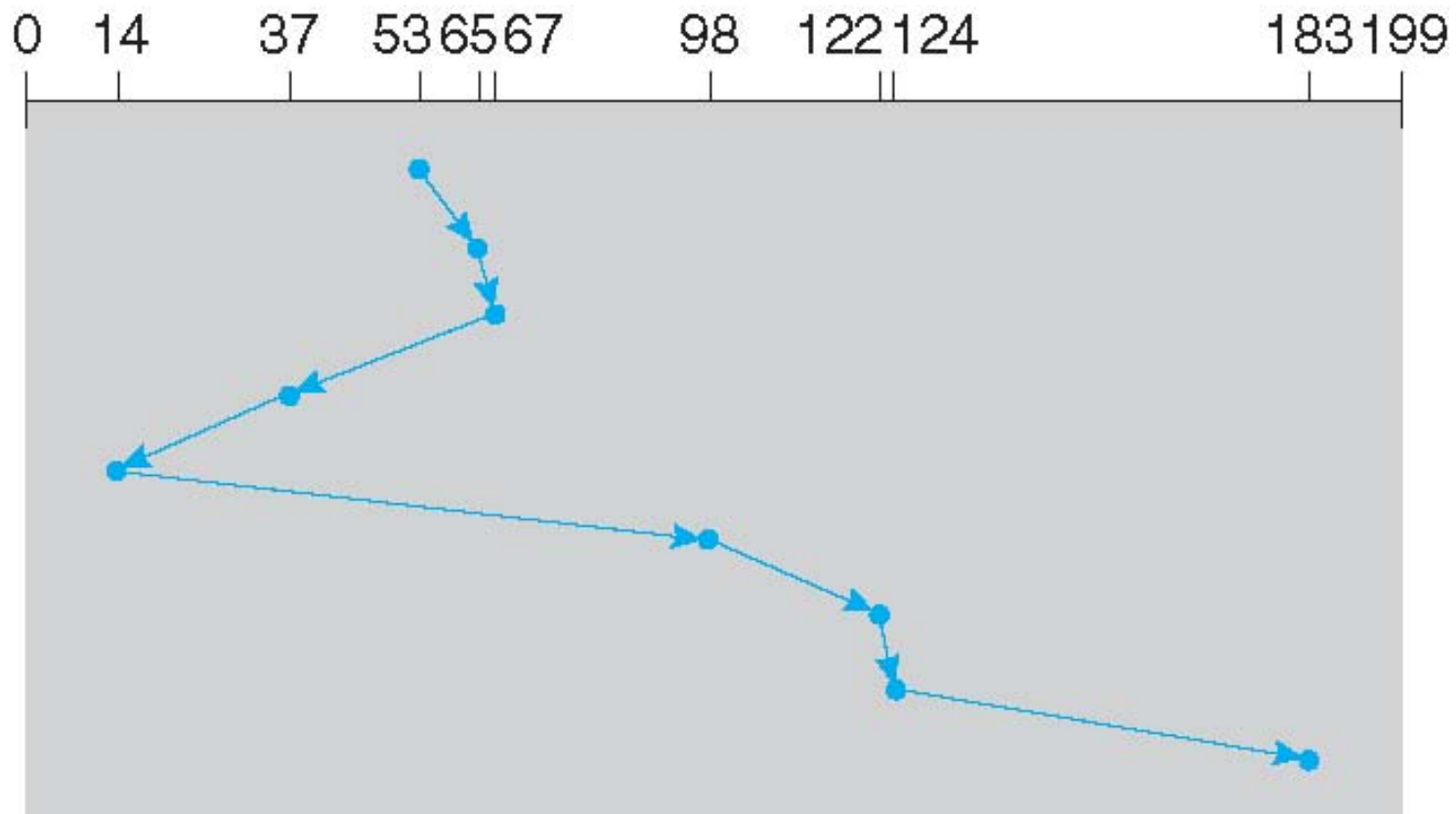
# SSF

- ❑ Selects the request with the minimum seek time from the current head position
- ❑ SSF scheduling is a form of SJF scheduling; may cause starvation of some requests
- ❑ Illustration shows total head movement of 236 cylinders

# SSF (Cont.)

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53



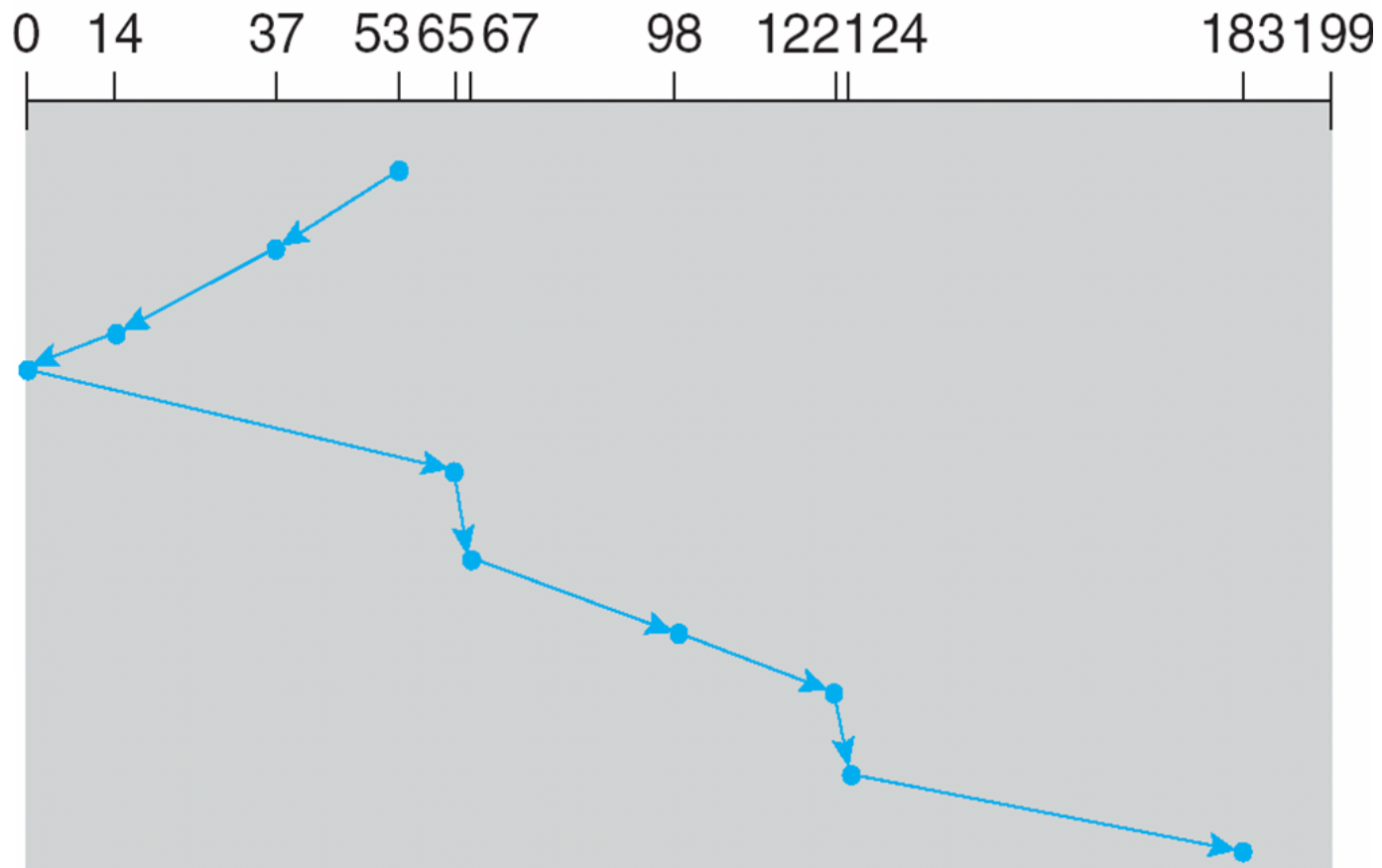
# SCAN/Elevator

- ❑ The disk arm starts at one end of the disk, and moves toward the other end, servicing requests until it gets to the other end of the disk, where the head movement is reversed and servicing continues.
- ❑ **SCAN algorithm** Sometimes called the **elevator algorithm**
- ❑ Illustration shows total head movement of 236 cylinders

# SCAN (Cont.)

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53



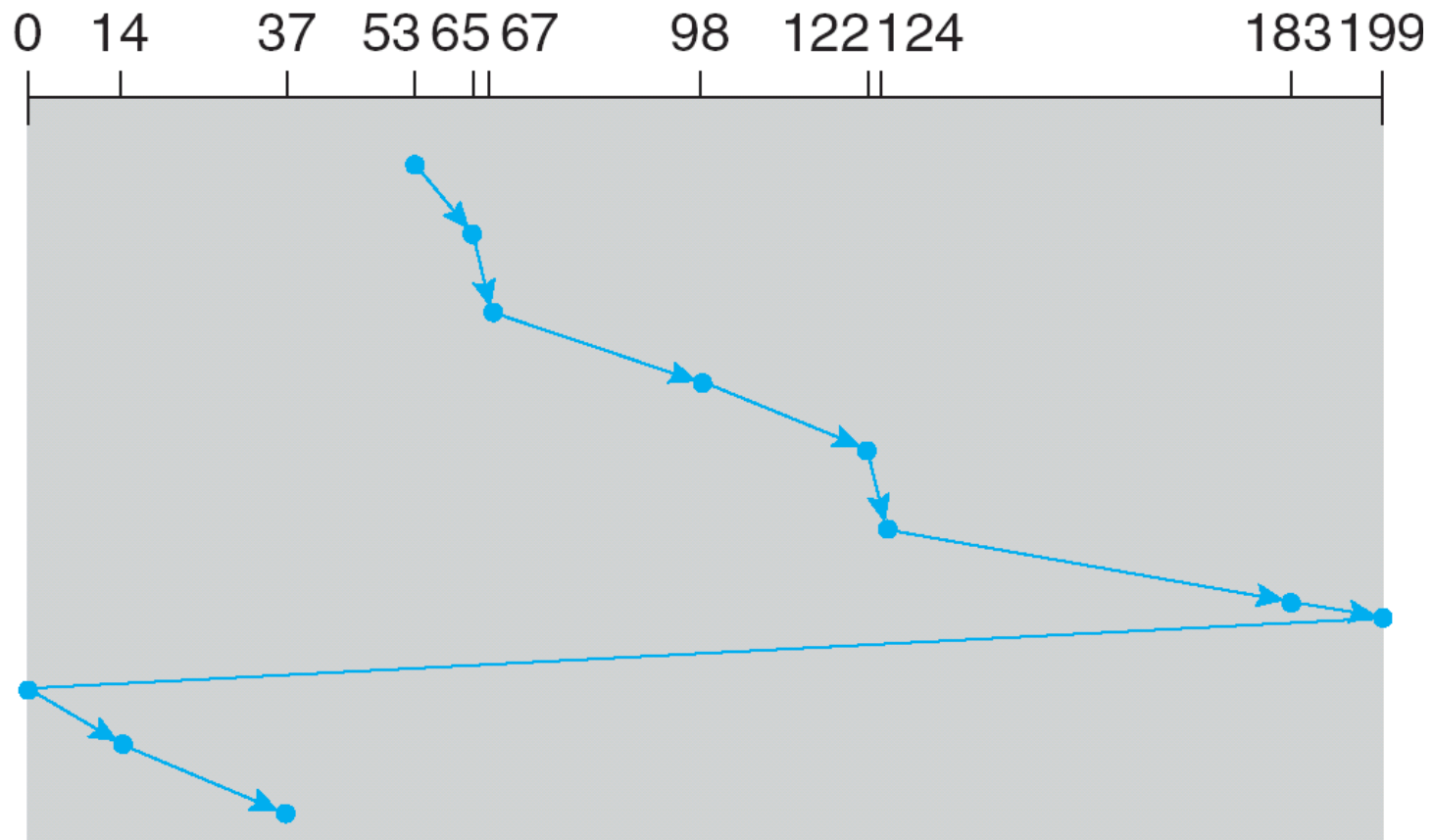
# C-SCAN

- ❑ Provides a more uniform wait time than SCAN
- ❑ The head moves from one end of the disk to the other, servicing requests as it goes
  - When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip
- ❑ Treats the cylinders as a circular list that wraps around from the last cylinder to the first one

# C-SCAN (Cont.)

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53





# C-LOOK

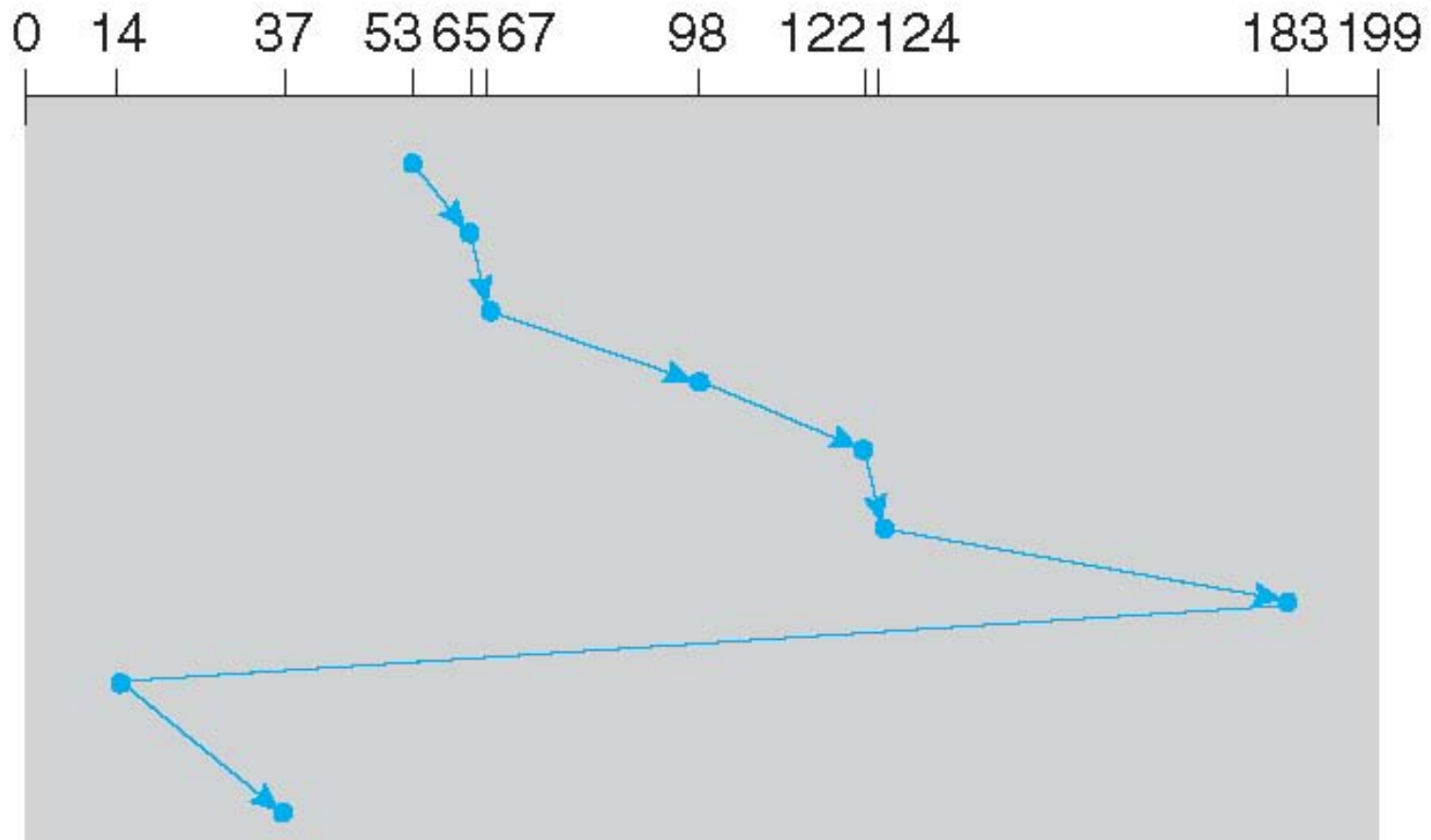
☐ Version of C-SCAN

☐ Arm only goes as far as the last request in each direction, then reverses direction immediately, without first going all the way to the end of the disk

# C-LOOK (Cont.)

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53



# Selecting a Disk-Scheduling Algorithm

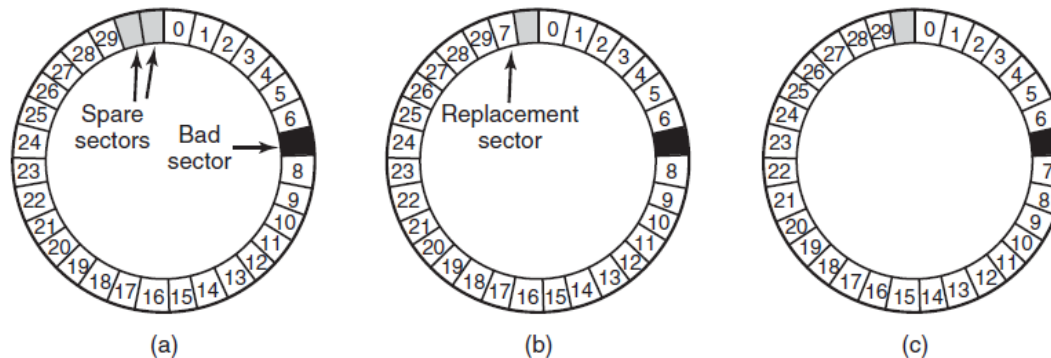
- ❑ SSTF is common and has a natural appeal
- ❑ SCAN and C-SCAN perform better for systems that place a heavy load on the disk
- ❑ Performance depends on the number and types of requests
- ❑ Requests for disk service can be influenced by the file-allocation method
- ❑ The disk-scheduling algorithm should be written as a separate module of the operating system, allowing it to be replaced with a different algorithm if necessary
- ❑ Either SSTF or LOOK is a reasonable choice for the default algorithm

# Disk Error Handling

- » Increasing linear bit densities constant drive by manufacturers
  - Defects introduced
- » Bad sector - sector that does not correctly read back the value written to them.
  - A few bits can be corrected by ECC
  - Larger defects can not be masked.
  - Deal with them in the controller or in the OS

# Disk Error Handling

- » Before a drive is shipped it is tested and bad sectors are written onto the disk
  - Spares substituted
    - Substitute or shift



(a) A disk track with a bad sector. (b) Substituting a spare for the bad sector. (c) Shifting all the sectors to bypass the bad one.

# Disk Error Handling

- » What about errors during runtime?
  - The first line of defense if ECC fails, try reading it again
    - Some read errors are transient: dust
  - Controller can switch to a spare before the sector dies completely

# Disk Error Handling

- » Option 2: Let the OS handle it
  - OS must acquire a list of bad sectors
    - Read them from the disk
    - Test entire disk
  - Can't let bad sectors reside in free block list or bitmap
    - Create a secret file
  - Backups can't backup bad block file
    - Sector by sector backups an issue

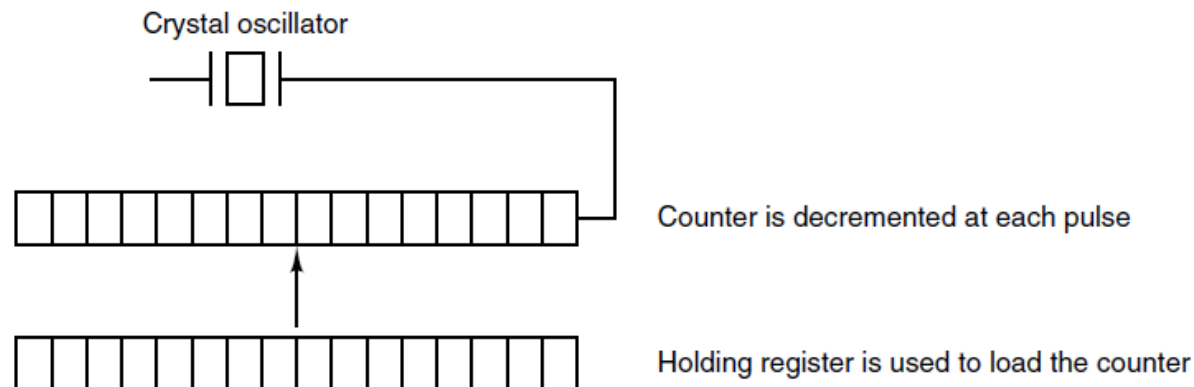
# Disk Error Handling

- » Seek errors also an issue.
  - Mechanical problem with the arm
  - Modern controllers handle arm issues automatically
  - 80's and 90's floppy drives set error bit and let the driver handle it
  - Driver issued a recalibrate command to move arm as far out as it would go and reset the current cylinder to 0



# Clocks

- » Essential to the operation of a system
- » Original clocks tied to 110 or 220 volt power and caused an interrupt every voltage cycle at 50 or 60 Hz
- » More recent built out of crystal oscillator, counter and holding register.



# Clocks

- » Two modes
  - One shot
  - square-wave mode
- » Interrupts called clock ticks
- » 500 Mhz crystal will pulse every 2 nsec
  - 32 bit register means interval from 2 sec to 8.6 sec
- » Battery backup
- » UNIX / Windows Epoch

# Clocks

- » Typical duties of a clock driver:
- » Maintaining the time of day.
- » Preventing processes from running longer than allowed.
- » Accounting for CPU usage.
- » Handling alarm system call from user processes.
- » Providing watchdog timers for parts of system itself.
- » Profiling, monitoring, statistics gathering.

# Clocks

- » 32 bit at 60Hz will overflow in a bit over 2 years
  - How can we store ticks since Jan 1, 1970?
  - 1. 64 bit counter
  - 2. Store time in seconds rather than ticks
    - Use secondary counter for ticks this second
    - 136 years
  - 3. Calculate ticks since boot
    - On boot store current time in memory
    - Use ticks as offset.

# Clocks

- » Processes can request clock notifications
- » If driver has enough clocks it can assign one to each process
  - Never enough clocks
  - Driver instead maintains a table with pending events
    - Store events as linked list indexed by time
- » Watchdog timers
  - Timer and procedure must be in the same address space

# Soft Timers

- » Interrupts have significant overhead.
  - Context switches
- » Polling has high latency
- » Soft timers avoid interrupts
  - Check real time mode every time before leaving kernel mode.
    - If timer expires then perform the event.
    - Already in kernel mode, no overhead

# Soft Timers

- » Soft timers stand or fall with the rate at which kernel entries are made for other reasons. These reasons include:
  - System calls.
  - TLB misses.
  - Page faults.
  - I/O interrupts.
  - The CPU going idle.
- » Average varies from 2 to 18 usec. 10 usec