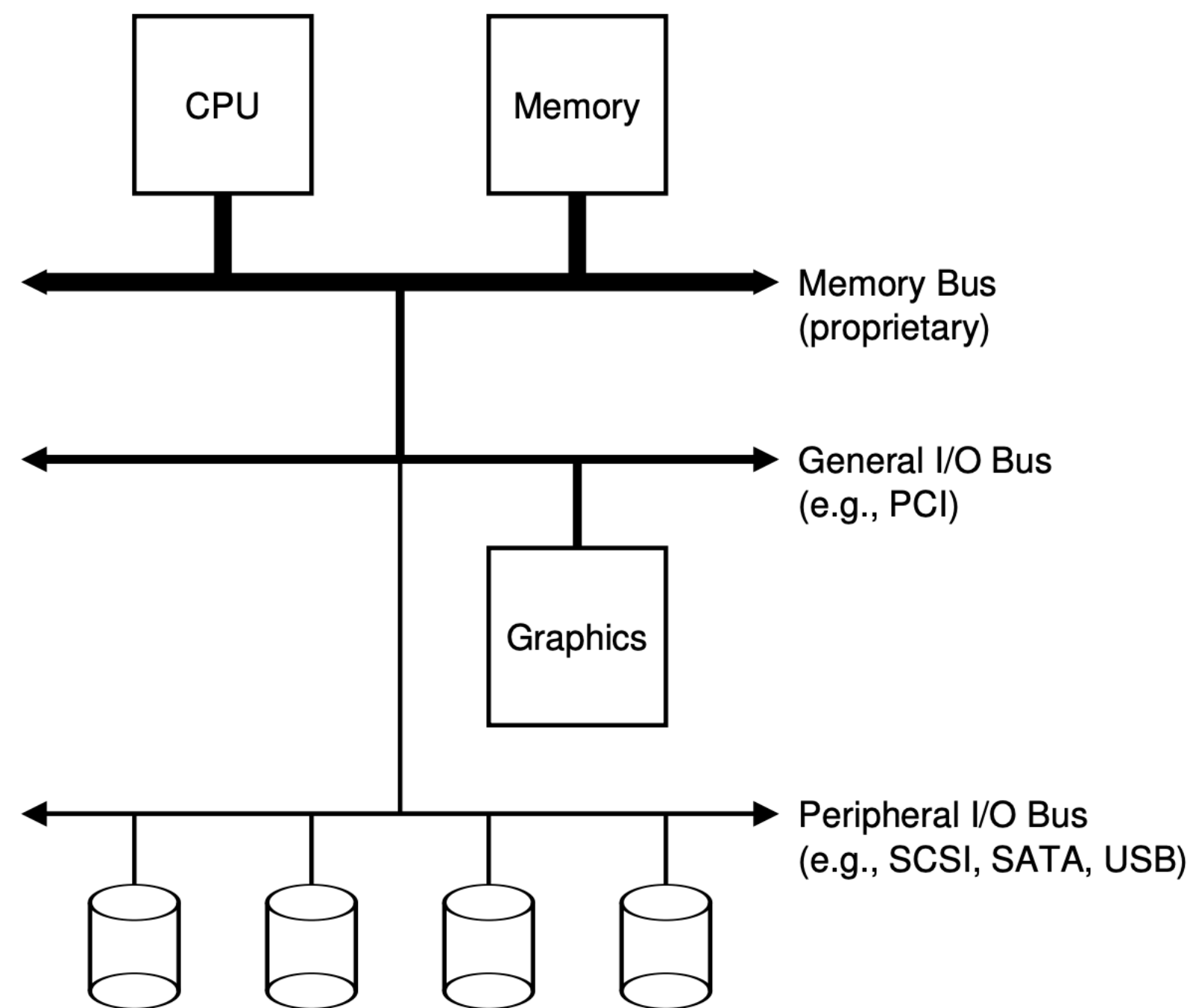


# CS202 (003): Operating Systems I/O

Instructor: Jocelyn Chen

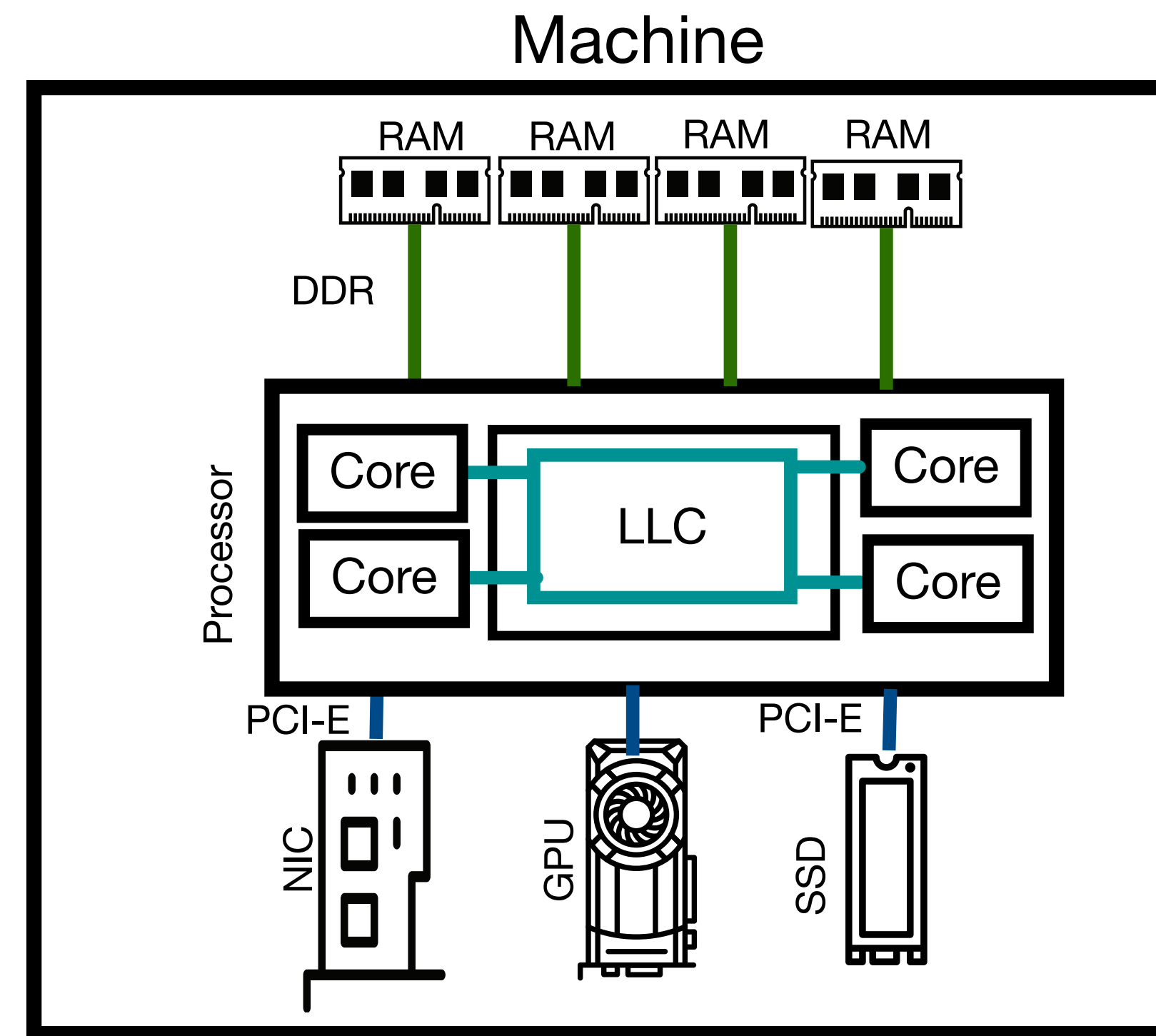
Last Time

# I/O Architecture at a high-level



Textbook (older)

Figure 36.1: Prototypical System Architecture



Newer

# CPU<sub>(kernel)</sub>/Device Interaction

## Mechanics of Communication

Explicit I/O instructions

`outb, inb, outw, inw`

WeesyOS boot.c (handout)

```

1 CS 202, Spring 2023
2 Handout 9 (Class 16)
3
4 1. Example use of I/O instructions: boot loader
5
6 Below is the WeensyOS boot loader
7
8 It may be helpful to understand the overall picture
9
10 This code demonstrates I/O, specifically with the disk: the
11 boot loader reads in the kernel from the disk.
12
13 See the functions boot_waitdisk() and boot_readsect(). Compare to Figures 36
14 .5 and 36.6 in OSTEP.
15
16 /* boot.c */
17 #include "x86-64.h"
18 #include "elf.h"
19
20 // boot.c
21 //
22 // WeensyOS boot loader. Loads the kernel at address 0x40000 from
23 // the first IDE hard disk.
24 //
25 // A BOOT LOADER is a tiny program that loads an operating system into
26 // memory. It has to be tiny because it can contain no more than 512 bytes
27 // of instructions: it is stored in the disk's first 512-byte sector.
28 //
29 // When the CPU boots it loads the BIOS into memory and executes it. The
30 // BIOS initializes devices and CPU state, reads the first 512-byte sector of
31 // the boot device (hard drive) into memory at address 0x7C00, and jumps to
32 // that address.
33 //
34 // The boot loader is contained in bootstart.S and boot.c. Control starts
35 // in bootstart.S, which initializes the CPU and sets up a stack, then
36 // transfers here. This code reads in the kernel image and calls the
37 // kernel.
38 //
39 // The main kernel is stored as an ELF executable image starting in the
40 // disk's sector 1.
41
42 #define SECTORSIZE 512
43 #define ELFHDR ((elf_header*) 0x10000) // scratch space
44
45 void boot(void) __attribute__((noreturn));
46 static void boot_readsect(uintptr_t dst, uint32_t src_sect);
47 static void boot_readseg(uintptr_t dst, uint32_t src_sect,
48 size_t filesz, size_t memsz);
49
50 // boot
51 // Load the kernel and jump to it.
52 void boot(void) {
53 // read 1st page off disk (should include programs as well as header)
54 // and check validity
55 boot_readseg((uintptr_t) ELFHDR, 1, PAGESIZE, PAGESIZE);
56 while (ELFHDR->e_magic != ELF_MAGIC) {
57 /* do nothing */
58 }
59
60 // load each program segment
61 elf_program* ph = (elf_program*) ((uint8_t*) ELFHDR + ELFHDR->e_phoff);
62 elf_program* eph = ph + ELFHDR->e_phnum;
63 for (; ph < eph; ++ph) {
64 boot_readseg(ph->p_va, ph->p_offset / SECTORSIZE + 1,
65 ph->p_filesz, ph->p_memsz);
66 }
67
68 // jump to the kernel
69 typedef void (*kernel_entry_t)(void) __attribute__((noreturn));
70 kernel_entry_t kernel_entry = (kernel_entry_t) ELFHDR->e_entry;
71 kernel_entry();
72 }

```

Integrated Drive Electronics.  
(connection btw a bus on  
motherboard & disk storage)

(executable & linkable format. Unix system)

```

73 program headers. (suitable chunks to load).
74
75 // boot_readseg(dst, src_sect, filesz, memsz)
76 // Load an ELF segment at virtual address 'dst' from the IDE disk's sector
77 // 'src_sect'. Copies 'filesz' bytes into memory at 'dst' from sectors
78 // 'src_sect' and up, then clears memory in the range
79 // '[dst+filesz, dst+memsz)'.
80 static void boot_readseg(uintptr_t ptr, uint32_t src_sect,
81 size_t filesz, size_t memsz) {
82 uintptr_t end_ptr = ptr + filesz;
83 memsz += ptr;
84
85 // round down to sector boundary
86 ptr &= ~(SECTORSIZE - 1);
87
88 // read sectors
89 for (; ptr < end_ptr; ptr += SECTORSIZE, ++src_sect) {
90 boot_readsect(ptr, src_sect);
91 }
92 // clear bss segment
93 for (; end_ptr < memsz; ++end_ptr) {
94 *(uint8_t*) end_ptr = 0;
95 }
96 }
97
98 // boot_waitdisk
99 // Wait for the disk to be ready.
100 static void boot_waitdisk(void) {
101 // Wait until the ATA status register says ready (0x40 is on)
102 // & not busy (0x80 is off) bit 7
103 while ((inb(0x1F7) & 0xC0) != 0x40) {
104 /* do nothing */
105 }
106 }
107
108 // boot_readsect(dst, src_sect)
109 // Read disk sector number 'src_sect' into address 'dst'.
110 static void boot_readsect(uintptr_t dst, uint32_t src_sect) {
111 // programmed I/O for "read sector"
112 boot_waitdisk();
113 outb(0x1F2, 1); // send 'count = 1' as an ATA argument
114 outb(0x1F3, src_sect); // send 'src_sect', the sector number
115 outb(0x1F4, src_sect >> 8);
116 outb(0x1F5, src_sect >> 16);
117 outb(0x1F6, (src_sect >> 24) | 0xE0);
118 outb(0x1F7, 0x20); // send the command: 0x20 = read sectors
119
120 // then move the data into memory
121 boot_waitdisk();
122 insl(0x1F0, (void*) dst, SECTORSIZE/4); // read 128 words from the disk
123
124 }
125
126
127
128

```

(block started by symbol)

group section by attribute to load

- takes no space in exe  
- need space in memory  
- uninit global/static variables  
- guaranteed to be zero before exe in (kth)

bit 6

read 1 sector at a time.

outb: writing to I/O port  
insl: read data.

Control Register:

Address 0x3F6 = 0x08 (0000 1RE0): R=reset,  
E=0 means "enable interrupt"

Command Block Registers: **each register stores 8-bit data**

Address 0x1F0 = Data Port

Address 0x1F1 = Error

Address 0x1F2 = Sector Count

Address 0x1F3 = LBA low byte

Address 0x1F4 = LBA mid byte

Address 0x1F5 = LBA hi byte

Address 0x1F6 = 1B1D TOP4LBA: B=LBA, D=drive

Address 0x1F7 = Command/status

Status Register (Address 0x1F7):

7	6	5	4	3	2	1	0
BUSY	READY	FAULT	SEEK	DRQ	CORR	IDDEX	ERROR

Error Register (Address 0x1F1): (check when ERROR==1)

7	6	5	4	3	2	1	0
BBK	UNC	MC	IDNF	MCR	ABRT	T0NF	AMNF

BBK = Bad Block

UNC = Uncorrectable data error

MC = Media Changed

IDNF = ID mark Not Found

MCR = Media Change Requested

ABRT = Command aborted

T0NF = Track 0 Not Found

AMNF = Address Mark Not Found

LBA: Linear block address (28-bit in the handout)

describes storage locations

each sector gets a unique number starting from 0

Status is read-only

Command is used to add instructions to the disk  
(e.g. 0xE0 is the write instruction)

Figure 36.5: **The IDE Interface** (Integrated Device Electronics): connection between a bus on the motherboard and disk storage

```

static int ide_wait_ready() {
    while (((int r = inb(0x1f7)) & IDE_BSY) || !(r & IDE_DRDY))
        ; // loop until drive isn't busy
    // return -1 on error, or 0 otherwise
}

static void ide_start_request(struct buf *b) {
    ide_wait_ready();
    outb(0x3f6, 0); // generate interrupt
    outb(0x1f2, 1); // how many sectors?
    outb(0x1f3, b->sector & 0xff); // LBA goes here ...
    outb(0x1f4, (b->sector >> 8) & 0xff); // ... and here
    outb(0x1f5, (b->sector >> 16) & 0xff); // ... and here!
    outb(0x1f6, 0xe0 | ((b->dev&1)<<4) | ((b->sector>>24)&0x0f));
    if(b->flags & B_DIRTY){
        outb(0x1f7, IDE_CMD_WRITE); // this is a WRITE
        outsl(0x1f0, b->data, 512/4); // transfer data too!
    } else {
        outb(0x1f7, IDE_CMD_READ); // this is a READ (no data)
    }
}

void ide_rw(struct buf *b) {
    acquire(&ide_lock);
    for (struct buf **pp = &ide_queue; *pp; pp=&(*pp)->qnext)
        ; // walk queue
    *pp = b; // add request to end
    if (ide_queue == b) // if q is empty
        ide_start_request(b); // send req to disk
    while ((b->flags & (B_VALID|B_DIRTY)) != B_VALID)
        sleep(b, &ide_lock); // wait for completion
    release(&ide_lock);
}

void ide_intr() {
    struct buf *b;
    acquire(&ide_lock);
    if (!(b->flags & B_DIRTY) && ide_wait_ready() >= 0)
        insl(0x1f0, b->data, 512/4); // if READ: get data
    b->flags |= B_VALID;
    b->flags &= ~B_DIRTY;
    wakeup(b); // wake waiting process
    if ((ide_queue = b->qnext) != 0) // start next request
        ide_start_request(ide_queue); // (if one exists)
    release(&ide_lock);
}

```

**handle read/write request**

**read: data is read and valid bit is set**  
**write: data is written and dirty bit is cleaned**

**Figure 36.6: The xv6 IDE Disk Driver (Simplified)**

# CPU<sub>(kernel)</sub>/Device Interaction

## Mechanics of Communication

Explicit I/O instructions

`outb, inb, outw, inw`

WeesyOS boot.c (handout)

`keyboard_readc()`

reading keyboard input (handout)



# CPU<sub>(kernel)</sub>/Device Interaction

## Mechanics of Communication

Explicit I/O instructions

`outb, inb, outw, inw`

WeesyOS boot.c (handout)

`keyboard_readc()`

reading keyboard input (handout)

`console_show_cursor()`

setting blinking cursor (handout)

```

129 2. Two more examples of I/O instructions
130
131 (a) Reading keyboard input
132
133 The code below is an excerpt from WeensyOS's k-hardware.c
134
135 This reads a character typed at the keyboard (which shows up on the
136 "keyboard data port" (KEYBOARD_DATAREG)).
137
138 /* Excerpt from WeensyOS x86-64.h */
139 // Keyboard programmed I/O
140 #define KEYBOARD_STATUSREG 0x64
141 #define KEYBOARD_STATUS_READY 0x01
142 #define KEYBOARD_DATAREG 0x60
143
144 int keyboard_readc(void) {
145     static uint8_t modifiers;
146     static uint8_t last_escape;
147
148     if ((inb(KEYBOARD_STATUSREG) & KEYBOARD_STATUS_READY) == 0) {
149         return -1;
150     }
151
152     uint8_t data = inb(KEYBOARD_DATAREG);
153     uint8_t escape = last_escape;
154     last_escape = 0;
155
156     if (data == 0xE0) { // mode shift
157         last_escape = 0x80;
158         return 0;
159     } else if (data & 0x80) { // key release: matters only for modifier ke
160
161         int ch = keymap[(data & 0x7F) | escape];
162         if (ch >= KEY_SHIFT && ch < KEY_CAPSLOCK) {
163             modifiers &= ~(1 << (ch - KEY_SHIFT));
164         }
165         return 0;
166
167         int ch = (unsigned char) keymap[data | escape];
168
169         if (ch >= 'a' && ch <= 'z') {
170             if (modifiers & MOD_CONTROL) {
171                 ch -= 0x60;
172             } else if (!(modifiers & MOD_SHIFT) != !(modifiers & MOD_CAPSLOCK))
173             {
174                 ch -= 0x20;
175             }
176             } else if (ch >= KEY_CAPSLOCK) {
177                 modifiers ^= 1 << (ch - KEY_SHIFT);
178                 ch = 0;
179             } else if (ch >= KEY_SHIFT) {
180                 modifiers |= 1 << (ch - KEY_SHIFT);
181                 ch = 0;
182             } else if (ch >= CKEY(0) && ch <= CKEY(21)) {
183                 ch = complex_keymap[ch - CKEY(0)].map[modifiers & 3];
184             } else if (ch < 0x80 && (modifiers & MOD_CONTROL)) {
185                 ch = 0;
186             }
187         }
188     }
189     return ch;

```

*hardware ports / status flag for interact w/ keyboard.*

*is keyboard ready to be read?*

*read data.*

*reset.*

*signal for going into special key mode. (E0 + some other byte).*

*map data to actual char using keymap.*

*control key logic.*

*special keys.*

*function keys.*

*Shift caps lock*

```

190
191 (b) Setting the cursor position
192
193 The code below is also excerpted from WeensyOS's k-hardware.c. It
194 uses I/O instructions to set a blinking cursor somewhere on a 25 x 80
195 screen.
196
197 // console_show_cursor(cpos)
198 // Move the console cursor to position 'cpo', which should be between 0
199 // and 80 * 25.
200
201 void console_show_cursor(int cpos) {
202     if (cpo < 0 || cpo > CONSOLE_ROWS * CONSOLE_COLUMNS) {
203         cpo = 0;
204     }
205     outb(0x3D4, 14); // Command 14 = upper byte of position
206     outb(0x3D5, cpo / 256);
207     outb(0x3D4, 15); // Command 15 = lower byte of position
208     outb(0x3D5, cpo % 256);
209 }
210
211
212
213
214

```

*VGA controller command port.*

*VGA controller data port.*

*split 16 bit # to .2 8 bit #.*

*Special key mode.*

# CPU<sub>(kernel)</sub>/Device Interaction

## Mechanics of Communication

Memory-mapped I/O

Most of the physical address space contains regular RAM  
However, the lower memory addresses (650K-1MB) are special - they don't refer to actual RAM

WeesyOS console printing (handout)

Mar 26, 24 22:43

copyout.c

Page 1/1

```
1 #include <fcntl.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <sys/mman.h>
5 #include <sys/stat.h>
6 #include <sys/types.h>
7 #include <unistd.h>
8
9 void mmapcopy(int fd, int size);
10
11 int main(int argc, char **argv) {
12     struct stat stat;
13     int fd;
14
15     /* Check for required cmd line arg */
16     if (argc != 2) {
17         printf("usage: %s <filename>\n", argv[0]);
18         exit(0);
19     }
20
21     /* Copy input file to stdout */
22     if ((fd = open(argv[1], O_RDONLY, 0)) < 0)
23         perror("open");
24
25     fstat(fd, &stat);
26     mmapcopy(fd, stat.st_size);
27
28     close(fd);
29
30     return 0;
31 }
32
33 void mmapcopy(int fd, int size) {
34     /* Ptr to memory mapped area */
35     char *bufp;
36
37     bufp = mmap(NULL, size, PROT_READ, MAP_PRIVATE, fd, 0);
38
39     write(STDOUT_FILENO, bufp, size);
40
41     return;
42 }
43 }
```

# CPU<sub>(kernel)</sub>/Device Interaction

## Mechanics of Communication

Explicit I/O instructions

Memory-mapped I/O

Interrupts

# CPU<sub>(kernel)</sub>/Device Interaction

## Mechanics of Communication

Explicit I/O instructions

Memory-mapped I/O

Interrupts

Through memory

Both CPU and the device see the same memory, so they can use shared memory to communicate. (eg. DMA)

# Polling vs. Interrupt (vs. busy waiting)

Busy Waiting

```
while (!device_is_ready()) {  
    // CPU is stuck here,  
    repeatedly checking  
    // Consuming CPU cycles doing  
    nothing useful  
}
```

Simple but inefficient  
CPU stuck here running at full speed

# Polling vs. Interrupt (vs. busy waiting)

## Busy Waiting

```
while (!device_is_ready()) {  
    // CPU is stuck here,  
    repeatedly checking  
    // Consuming CPU cycles doing  
    nothing useful  
}
```

## Polling

```
void poll_device() {  
    while (1) {  
        if (device_is_ready())  
            { process_device_data(); }  
  
        // Sleep/delay for a set  
        interval  
        sleep_ms(100); // Check every  
        100ms  
    }  
}
```

System checks device status at  
regular intervals  
Lower CPU usage



# Polling vs. Interrupt (vs. busy waiting)

## Busy Waiting

```
while (!device_is_ready()) {  
    // CPU is stuck here,  
    repeatedly checking  
    // Consuming CPU cycles doing  
    nothing useful  
}
```

## Interrupts

```
// CPU sets up interrupt handler  
set_interrupt_handler  
    (device_interrupt_handler);  
  
// CPU goes on to do other useful work  
do_other_work();  
  
// When device needs attention,  
it triggers interrupt  
// and the handler runs automatically  
void device_interrupt_handler() {  
    // Handle the device's needs  
    process_device_data();  
}
```

## Polling

```
void poll_device() {  
    while (1) {  
        if (device_is_ready())  
            { process_device_data(); }  
  
        // Sleep/delay for a set  
        interval  
        sleep_ms(100); // Check every  
        100ms  
    }  
}
```

Most sophisticated approach  
If interrupt rate is high, we can get **livelock**

# Programmed I/O vs. DMA (Direct memory access)

Programmed I/O

CPU writes data directly to device, and reads data directly from device.

# Programmed I/O vs. DMA (Direct memory access)

## Programmed I/O

CPU writes data directly to device, and reads data directly from device.

## DMA

CPU places some buffers in main memory.

Tells device where the buffers are

Then "pokes" the device by writing to register

Then device uses DMA to read or write the

The CPU can poll to see if the DMA completed (or the device can interrupt the CPU when done).

CPU don't have to constantly deal with small amount of data transfer, the device can write the contents straight into memory

# Software architecture: device drivers

Device drivers act as a bridge between hardware devices and the operating system kernel

# Software architecture: device drivers

Device drivers act as a bridge between hardware devices and the operating system kernel

```
reset(): Initializes or resets the device  
ioctl(): Provides device-specific controls  
read()/write(): Standard data transfer operations  
handle_interrupt(): Manages hardware interrupts
```

# Software architecture: device drivers

Device drivers act as a bridge between hardware devices and the operating system kernel

`reset()`: Initializes or resets the device  
`ioctl()`: Provides device-specific controls  
`read()/write()`: Standard data transfer operations  
`handle_interrupt()`: Manages hardware interrupts

Example interface

Advantages: don't have to worry about the specific hardware implementation

Issues:

1. Device drivers is per-OS and per-device (hard part cannot be reused)
2. Bugs in device drivers often bring down the entire machine

# Synchronous vs. Asynchronous I/O

## Synchronous I/O

When a process makes a system call (like `read()` or `write()`), it blocks (suspends) until the operation completes

The process enters a sleep state and cannot execute other codes

Control returns to process after operation finishes

Code is often more readable, but it is slow

# Synchronous vs. Asynchronous I/O

## Synchronous I/O

When a process makes a system call (like `read()` or `write()`), it blocks (suspends) until the operation completes

The process enters a sleep state and cannot execute other codes

Control returns to process after operation finishes

## Async I/O

System calls return immediately, even if the operation isn't completed

Instead of blocking, the call returns a status indicating what would have happened

Check through polling or interrupt

Need to use platform-specific extensions to POSIX to do async I/O for files