

# Systems 3

## Input/Output

Marcel Waldvogel

Department of Computer and Information Science  
University of Konstanz

Winter 2019/2020

# Chapter Goals

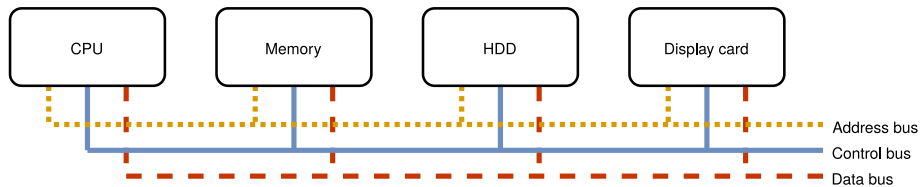
- What are some typical I/O requirements (rates, latencies)?
- How is I/O achieved? How has it improved over time?
- What is the motivation for device drivers?
- What is their interface/interaction with the OS?
- What is their operation?
- Explain the delays related to hard disk I/O.

# I/O Devices

Device	Data rate
Keyboard	10 B/s
Mouse	100 B/s
52x CD-ROM	8 MB/s
USB 2.0	60 MB/s
Gigabit Ethernet	125 MB/s
SATA hard disk	100...200 MB/s
SATA 6G bus	480 MB/s
PCI bus	528 MB/s
USB 3.0	625 MB/s
PCIe 4.0 x4	8 GB/s
PCIe 6.0 x16	128 GB/s

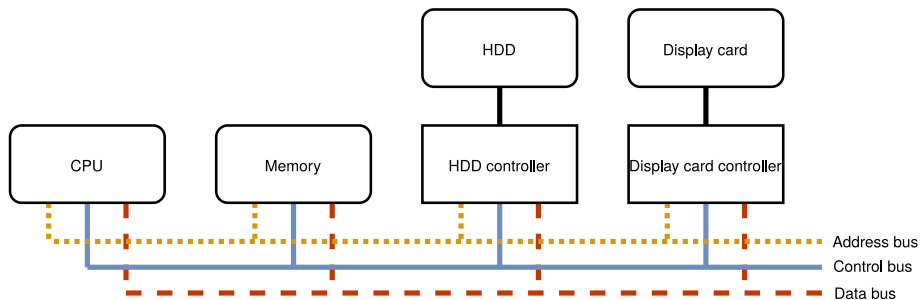
**Table:** Some typical device, network, and bus data rates.

# Simple Bus



**Figure:** Simple system bus.

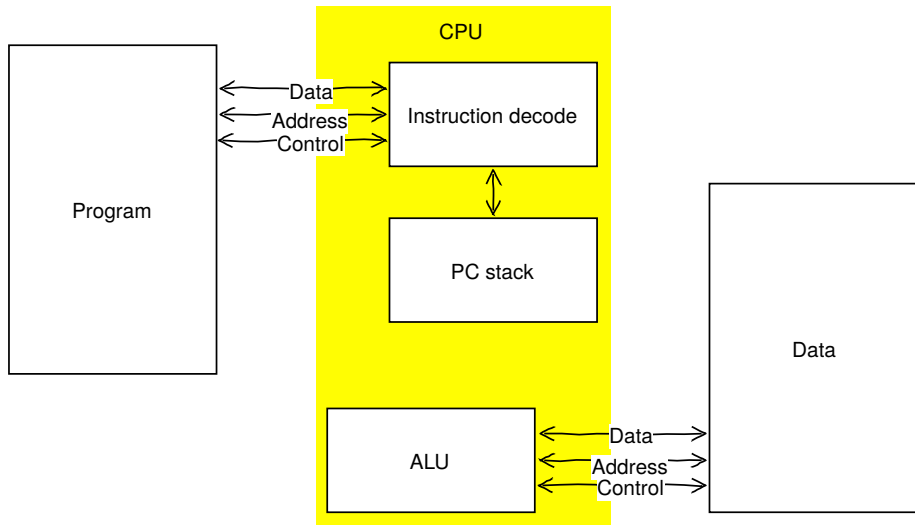
# Device Controllers



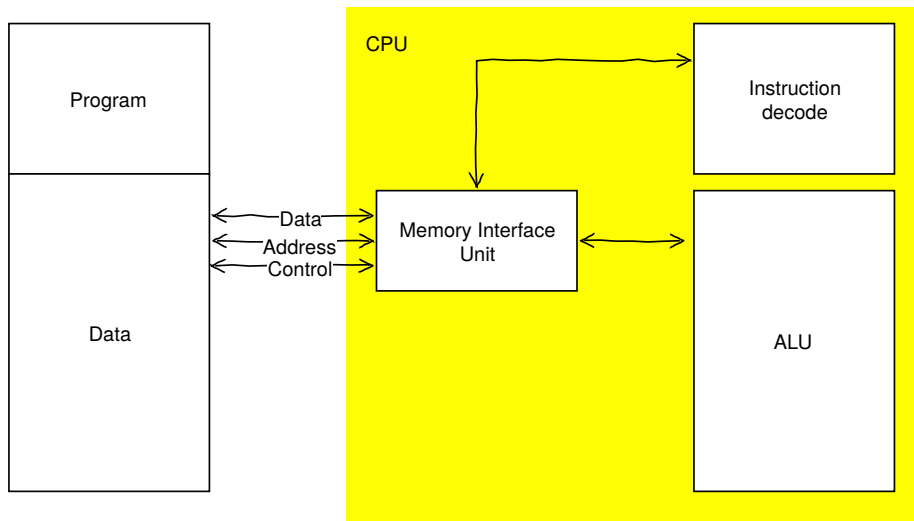
**Figure:** Simple system bus with device controllers.

Today's computer buses are more hierarchical. See e.g. [AMD Zen architecture](#).

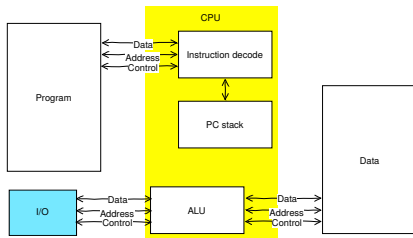
# Recapitulation: Harvard architecture



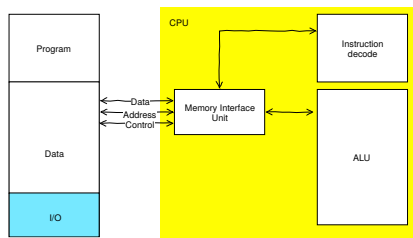
# Recapitulation: Von-Neumann architecture



# How to access I/O?



Harvard-style I/O: Separate I/O address space (and bus)



Von-Neumann-style I/O: Single I/O address space (and bus)

(Many hybrid forms exist as well)



# MIPS architecture

Register	Assembly name	Comment
r0	\$zero	Always zero
r1	\$at	<b>A</b> sembler <b>T</b> emp: Reserved for assembler
r2-r3	<b>\$v0-\$v1</b>	<b>V</b> alue
r4-r7	<b>\$a0-\$a3</b>	Function call <b>A</b> rguments
r8-r15	<b>\$t0-\$t7</b>	<b>T</b> emporary values (not saved)
r16-r23	<b>\$s0-\$s7</b>	<b>S</b> aved values
r24-r25	<b>\$t8-\$t9</b>	<b>T</b> emporary values (not saved)
r26-r27	\$k0-\$k1	Reserved for OS <b>K</b> ernel
r28	\$gp	<b>G</b> lobal <b>P</b> ointer
r29	<b>\$sp</b>	<b>S</b> tack <b>P</b> ointer
r30	<b>\$fp</b>	<b>F</b> rame <b>P</b> ointer
r31	<b>\$ra</b>	<b>R</b> eturn <b>A</b> ddress

# MIPS function call register lifetime

Register	Function entry	Function exit	Saved by
\$v0-\$v1	Undefined	Return value or undefined	caller
\$a0-\$a3	Arguments or undefined	clobbered	caller
\$t0-\$t9	Undefined	clobbered	caller
\$s0-\$s7	Undefined	unmodified	callee

```

1 float log_base(float x, float b)
2 {
3     /* Save $s0-$s7 if used in this function */
4     float lx, lb, lr;
5     /* Save $a1 and any of $t0-$t9 which should survive */
6     lx = log(x);
7     /* Restore $a1, $t0-$t9 */
8     /* Save $t0-$t9, if needed */
9     lb = log(b);
10    /* Restore $t0-$t9 */
11    lr = lx / lb;
12    /* Restore $s0-$s7 if they were used in this function */
13    return lr; /* Fill $v0 */
14 }

```

# Someone's gotta save<sup>1</sup> (massively simplified MIPS-32)

```

1 float log_base(float x, float b)
2 { float lx, lb, lr;   lx = log(x);   lb = log(b);   lr = lx / lb;   return lr; }

1 log_base:
2     add    $sp,$sp,-12        /* Reserve space on stack */
3
4     st     $ra,8($sp)         /* Save $ra, $a1 */
5     st     $a1,4($sp)
6     jal    log                /* Call log() with $a0 to $v0; saves $pc in $ra */
7
8     st     $v0,0($sp)         /* Save $v0 */
9     ld     $a0,4($sp)         /* Restore original $a1 as $a0 */
10    jal    log                /* Call log with $a0 to $v0; saves $pc in $ra */
11
12    ld     $t0,0($sp)          /* Restore saved $v0 as $t0 */
13    div.s   $v0,$t0,$v0        /* $v0 = $t0 / $v0 */
14
15    ld     $ra,8($sp)          /* Restore $ra */
16    add     $sp,$sp,12         /* Free space on stack */
17    j       $ra               /* Restores $pc from $ra */

```

<sup>1</sup>With more local variables, we would need to save/restore \$s0-\$s7 at entry/exit

# I/O for single-process machines

## Programmed I/O (polling or delay loops)

```
1  /* Disk I/O; diskdevicespace points to I/O address of disk controller */
2  diskdevicespace->sector = s;
3  diskdevicespace->track  = t;
4  diskdevicespace->head   = h;
5  diskdevicespace->command = DISK_READ;
6  while (!(diskdevicespace->status & DISK_OP_COMPLETE)) {
7      while (!(diskdevicespace->status & DISK_BYTE_READY)) {
8          /* Busy waiting */
9      }
10     *dst++ = diskdevicespace->data;
11 }
12 switch (diskdevicespace->status & DISK_ERROR_MASK) {
13 case OK:
14     ...
15 }
```

# I/O for multi-process machines: Slave device

## Interrupt+PIO

```
1  /* Disk I/O; diskdevicespace points to I/O address of disk controller */
2  diskdevicespace->sector = s;
3  diskdevicespace->track  = t;
4  diskdevicespace->head   = h;
5  diskdevicespace->command = DISK_READ;
6  switch_to_other_process();
7
8  void interrupt_handler(void)
9  {
10     ...
11     if (diskdevicespace->status & DISK_INTERRUPTED) {
12         while (!(diskdevicespace->status & DISK_OP_COMPLETE)) {
13             *dst++ = diskdevicespace->data;
14         }
15         switch (diskdevicespace->status & DISK_ERROR_MASK) {
16             case OK:
17                 ...
18             }
19     }
20     ...
21 }
```

# Intermission: What is an interrupt?

## Interrupt handler invocation

‘Involuntary’ subprogramm call

- **Not** triggered by an opcode being executed
- Triggered by external hardware (interrupt pin),  
“a **call** opcode inserted between unsuspecting opcodes”
- Executes with privileges

## Interrupt handler structure

- 1 Save registers
- 2 Check/handle device activity
- 3 Restore registers
- 4 Return to (unsuspecting) calling program

# I/O for multi-process machines: Slave device

## Interrupt+DMA

```
1  /* Disk I/O; diskdevicespace points to I/O address of disk controller */
2  dmadevicespace->base_address = dst;
3  diskdevicespace->sector   = s;
4  diskdevicespace->track    = t;
5  diskdevicespace->head     = h;
6  diskdevicespace->command = DISK_READ;
7  switch_to_other_process();
8
9  void interrupt_handler(void)
10 {
11     ...
12     if (diskdevicespace->status & DISK_INTERRUPTED) {
13         switch (diskdevicespace->status & DISK_ERROR_MASK) {
14             case OK:
15                 ...
16             }
17         }
18     ...
19 }
```

# I/O for multi-process machines: Master device

## Interrupt+Bus Master

```
1  /* Disk I/O; diskdevicespace points to I/O address of disk controller */
2  diskdevicespace->transfer_address = dst;
3  diskdevicespace->sector   = s;
4  diskdevicespace->track    = t;
5  diskdevicespace->head     = h;
6  diskdevicespace->command = DISK_READ;
7  switch_to_other_process();
8
9  void interrupt_handler(void)
10 {
11     ...
12     if (diskdevicespace->status & DISK_INTERRUPTED) {
13         switch (diskdevicespace->status & DISK_ERROR_MASK) {
14             case OK:
15                 ...
16             }
17     }
18     ...
19 }
```



# I/O evolution

## Comparison

**Polling+PIO** CPU busy-waits for device to be ready and transmits all bytes itself

**Interrupt+PIO** CPU works on other stuff<sup>2</sup> and is interrupted by the device, when it is ready. Bytes are transferred by the CPU

**Interrupt+DMA** CPU works on other stuff. When the device is ready, the DMA controller does the byte transfer instead of the CPU. The CPU is only interrupted at the end.

**Interrupt+Bus master** The device itself can access the memory as needed<sup>3</sup>.

---

<sup>2</sup>or goes to sleep, if no other activity is pending

<sup>3</sup>could also read the next job from a job queue

# I/O software in the OS: Goals

- 1 device independence
- 2 uniform naming
- 3 error handling
- 4 synchronous vs. asynchronous
- 5 buffering

# I/O Software layers

I/O Software is often organized in four layers:

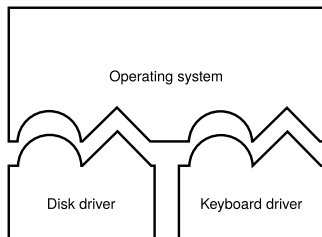
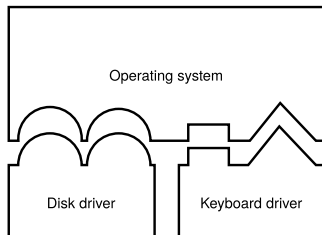
Hardware

- 1 Interrupt handlers
- 2 Device drivers
- 3 Device-independent OS software (next slide)
- 4 User-level I/O software

User

# Device-Independent I/O Software

- 1 Uniform interfacing
- 2 Buffering
- 3 Error reporting
- 4 Dedicated devices
- 5 Block size



# User-level I/O software

- library procedures (e.g. I/O calls, formatting)
- spooling

# Storage devices

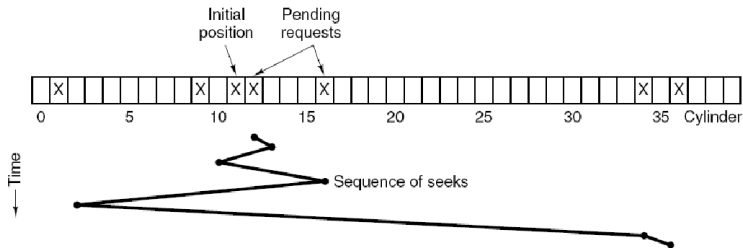
- Tape
- DVD/CD/.../WORM
- HDD
- SSD/Flash
- RAM disk

# Disk Software

## Read/Write timing factors

- 1 Seek time: Arm onto right track  
1 ms (track-to-track)... 10 ms (average random seek)
- 2 Rotational delay: Sector start under head  
 $\frac{1}{2}$  rotation:  $1/(\text{rotation speed [RPM]}/60)$
- 3 Data transfer time: Sector(s) passing by  
bit density [b/cm] \* rotation speed [RPM]/60 \* circumference [cm]

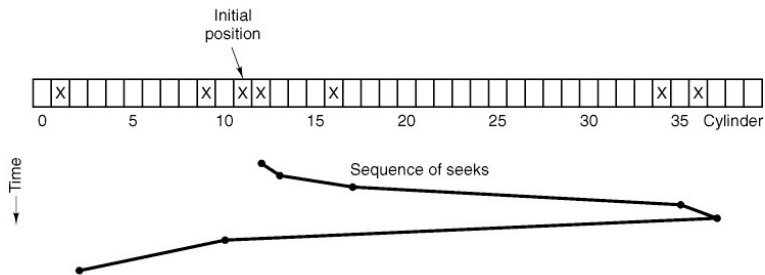
# Disk Arm Movement (1)



**Figure:** Shortest Seek First (SSF) disk scheduling algorithm. (Tanenbaum fig. 3.21)



# Disk Arm Movement (2)



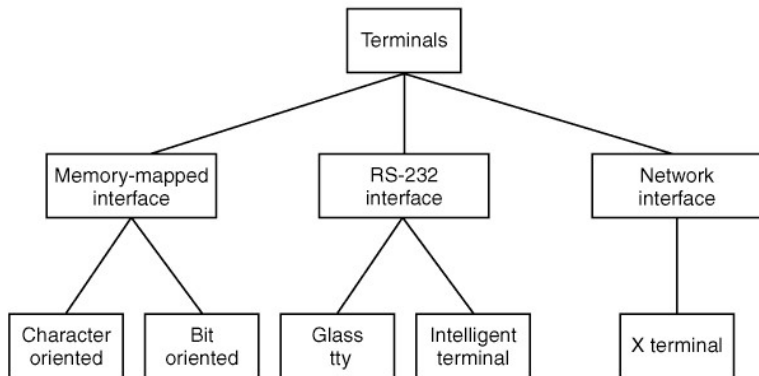
**Figure:** The elevator algorithm for scheduling disk requests. (Tanenbaum fig. 3.22)

# Common Hard Drive Errors

- 1** Programming error  
e.g. request for nonexistent sector
- 2** Transient checksum error  
e.g. caused by vibration during read
- 3** Permanent checksum error  
e.g. disk block physically damaged
- 4** Seek error  
e.g. arm was sent to cylinder 6 but it went to 7
- 5** Controller error  
e.g. controller refuses to accept commands

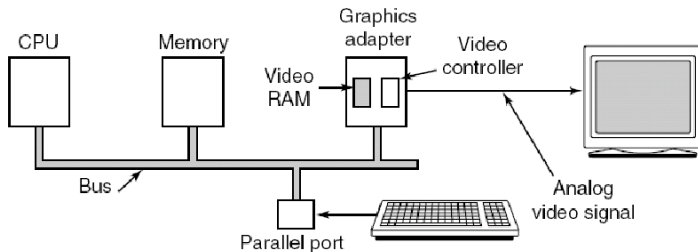
# Terminals

# Terminal Hardware



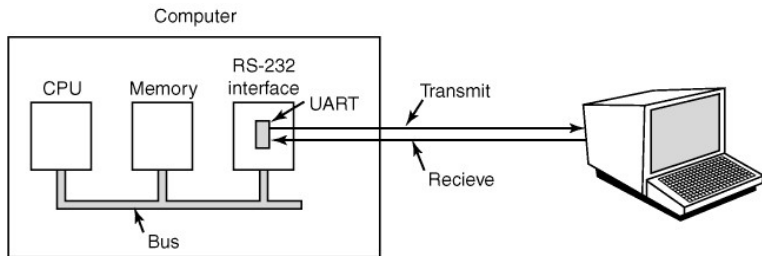
**Figure:** Different types of terminals. (Tanenbaum fig. 3.24)

# Memory-mapped interface



**Figure:** Memory-mapped terminals write directly into video RAM. (Tanenbaum fig. 3.25)

# RS-232 interface



**Figure:** An RS-232 terminal communicates with a computer over a communication line, one bit at a time. The computer and the terminal are completely independent. (Tanenbaum fig. 3.27)

# Input software

Character	POSIX name	Comment
CTRL-D	EOF	End of file
	EOL	End of line
CTRL-C	INTR	Interrupt process (SIGINT)
CTRL-U	KILL	Erase entire line beeing typed

**Table:** Characters that are handled specially in canonical (cooked) mode.

# Control/Escape sequences

Code	...0	...1	...2	...3	...4	...5	...6	...7
0...	<b>NUL</b>	SOH	STX	ETX	EOT	ENQ	ACK	<b>BEL</b>
1...	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB
Code	...8	...9	...A	...B	...C	...D	...E	...F
0...	<b>BS</b>	<b>HT</b>	<b>LF</b>	VT	FF	<b>CR</b>	SO	SI
1...	CAN	EM	SUB	<b>ESC</b>	FS	GS	RS	US
7...	x	y	z	{		}	~	<b>DEL</b>

Sequence	Comment
ESC [31;42m	red letters on green background
ESC [0m	reset all attributes
ESC [1E	move cursor to beginning of next line
ESC [5T	scroll page down by 5 lines

**Table:** ANSI escape sequences<sup>5</sup> are used to control the terminal and are not interpreted as text.

<sup>5</sup>Introduced by 'ESC [' aka 'CSI' (control sequence introducer); ended by letter



# Input codes

**ASCII codes** Only 7 bits, legacy<sup>6</sup>  
*USASCII code chart*

The diagram shows a 3D perspective of a bit plane with bits b7, b6, b5, b4, b3, b2, b1, and b0. Arrows indicate the flow of data from the bit plane into a 16x8 grid. The grid is labeled with 'Column' and 'Row' headers. The first column contains bit patterns (00, 01, 10, 11) and the second column contains character codes (0-15). The grid contains the following characters:

Column	Row	0	1	2	3	4	5	6	7
00	0	NUL	DLE	SP	0	@	P	\	p
01	1	SOH	DC1	!	1	A	Q	a	q
10	2	STX	DC2	"	2	B	R	b	r
11	3	ETX	DC3	#	3	C	S	c	s
00	4	EOT	DC4	\$	4	D	T	d	t
01	5	ENQ	NAK	%	5	E	U	e	u
10	6	ACK	SYN	&	6	F	V	f	v
11	7	BEL	ETB	'	7	G	W	g	w
00	8	BS	CAN	(	8	H	X	h	x
01	9	HT	EM	)	9	I	Y	i	y
10	10	LF	SUB	*	:	J	Z	j	z
11	11	VT	ESC	+	;	K	[	k	{
00	12	FF	FS	,	<	L	\	l	
01	13	CR	GS	=	=	M	]	m	}
10	14	SO	RS	.	>	N	^	n	~
11	15	SI	US	/	?	O	_	o	DEL

**Unicode** 20.1 bits(!), current (files, screen); often as encoded as UTF-8 (backward compatible to ASCII)

**Scan codes** Position on keyboard; unfortunately current (USB, Bluetooth input devices)

<sup>6</sup>Dozens of mostly incompatible 8 bit extensions exist