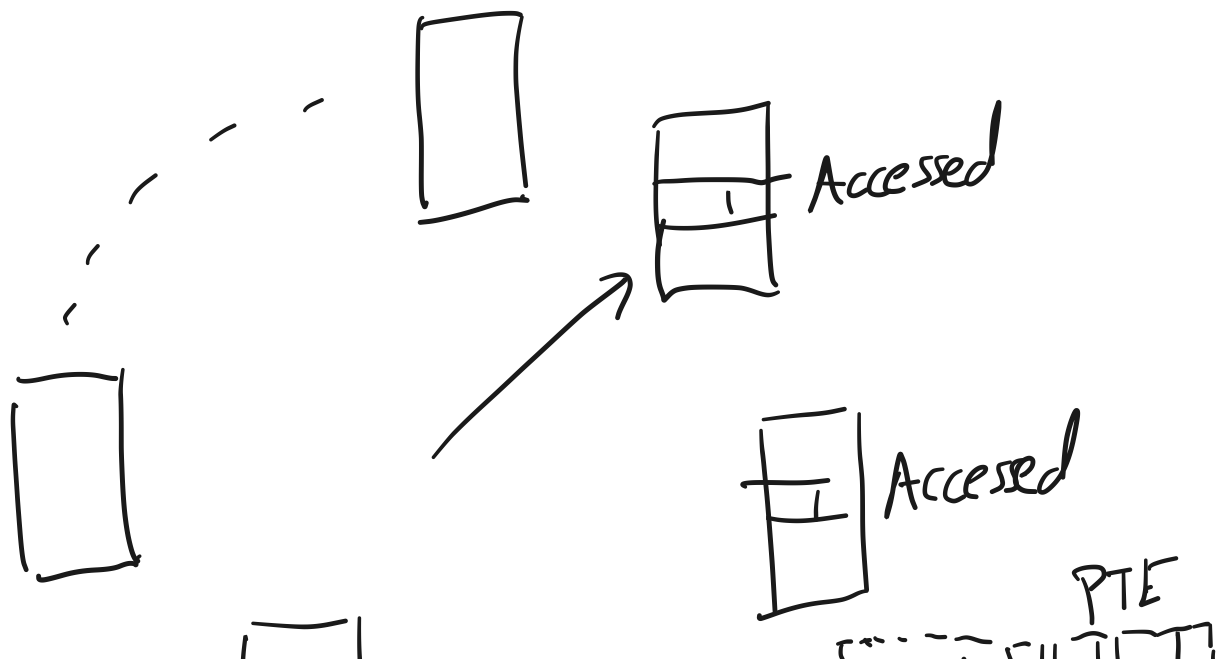
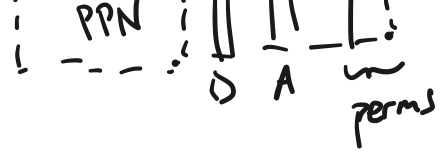
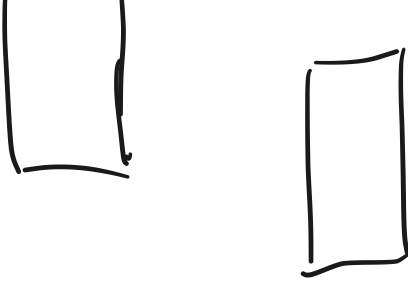


- ☑ 1. Last time
  - ☑ 2. Page replacement policies (continued)
  - ☑ 3. Thrashing
  - ☑ 4. mmap()
  - ☑ 5. Where does the OS live?
  - 6. Weensy OS
- 

## 2. Page replacement, contd.

- OPT minimizes misses/swaps/evictions
  - but can't be implemented in general
- LRU: approximates OPT (assuming what?)
- approximate LRU with CLOCK





H/w sets Accessed + Dirty bits

OS consumes these bits and clears them

- Generalization of CLOCK:  $N^{\text{th}}$  Chance (see notes)

---

### 3. Thrashing

Ex

- Prog. touches 50 pages, each equally likely
- 40 slots (pages) available in RAM

100ns/mem ref

10ms/page fault that reads from disk

$$\frac{4 \left( \frac{100\text{ns}}{\text{mem ref}} \right) + 1 \left( \frac{10\text{ms}}{\text{page fault}} \right)}{5} \approx \frac{2\text{ms}}{\text{ref}}$$

shed load:

- (a) working set
- (b) pg fault freq.

4. mmap()

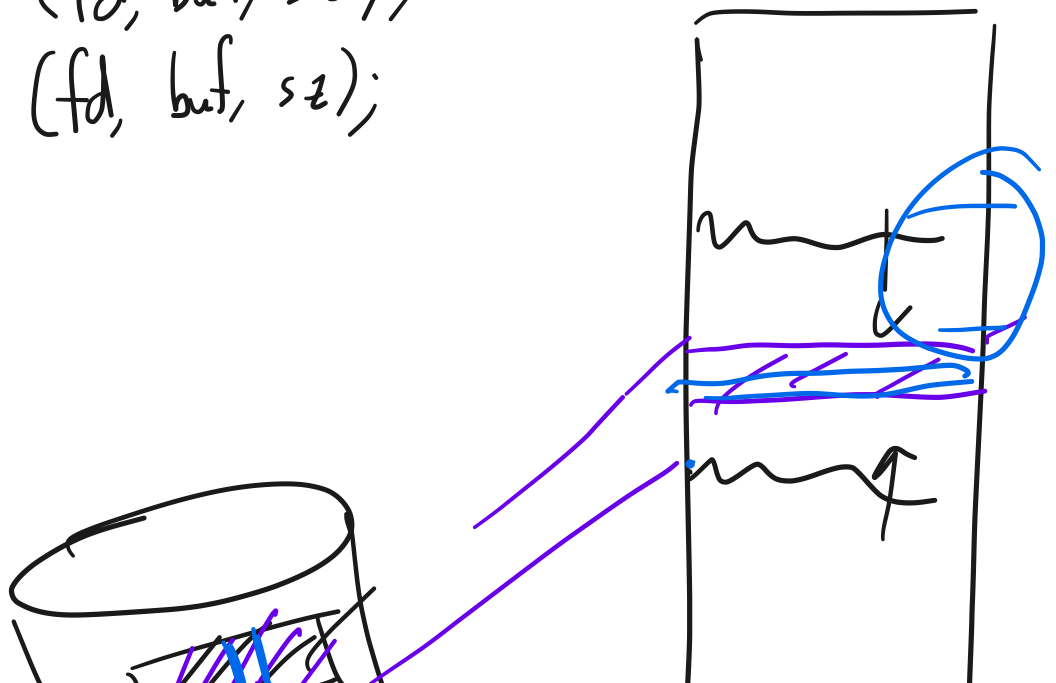
recall:

fd = open(pathname, mode);

rc = write(fd, buf, sz);

rc = read(fd, buf, sz);

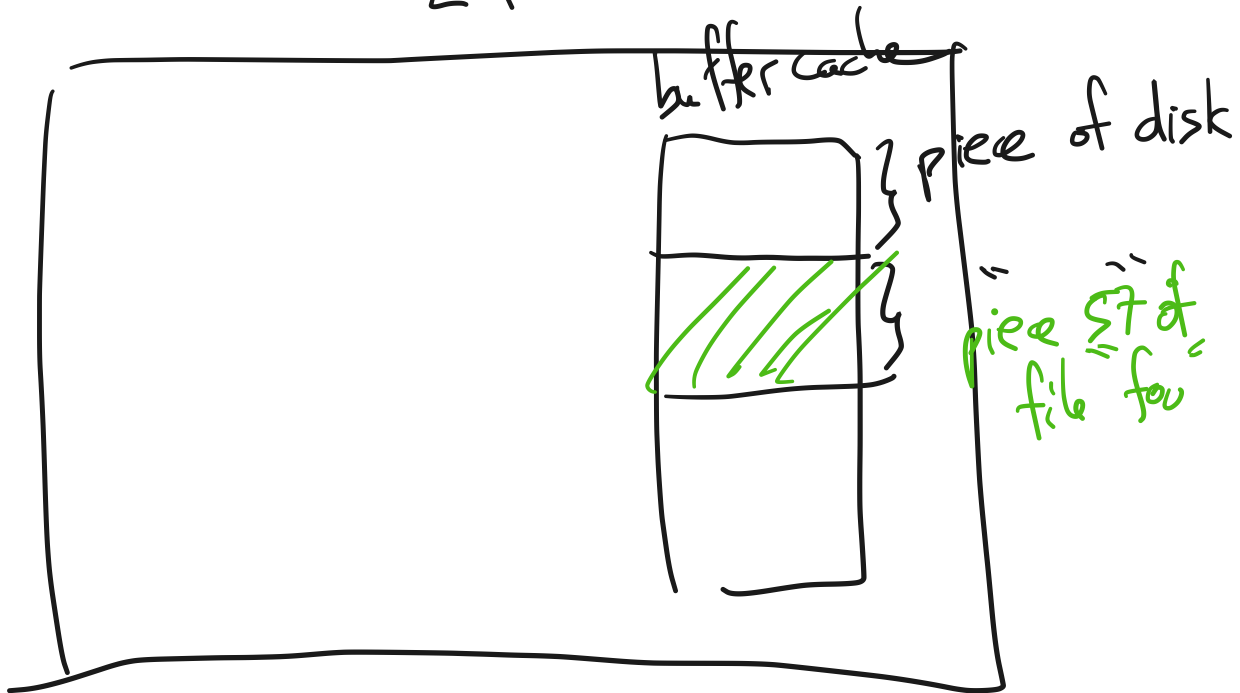
mmap()





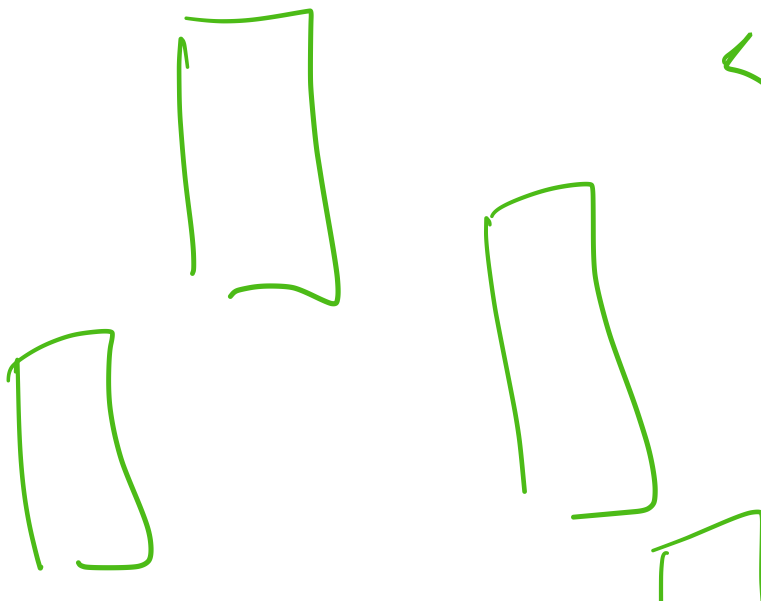
$\langle \langle \text{proc}, \text{VPN}^?, \text{entry}^? \rangle \rangle$

OS



page table for proc that mapped  
file "foo"

page table gets a new  
mapping  
 $\langle \langle \text{mapped ptr}, \text{VPN of the relevant buffer cache entry} \rangle \rangle$





---

median: 74

avg: 72

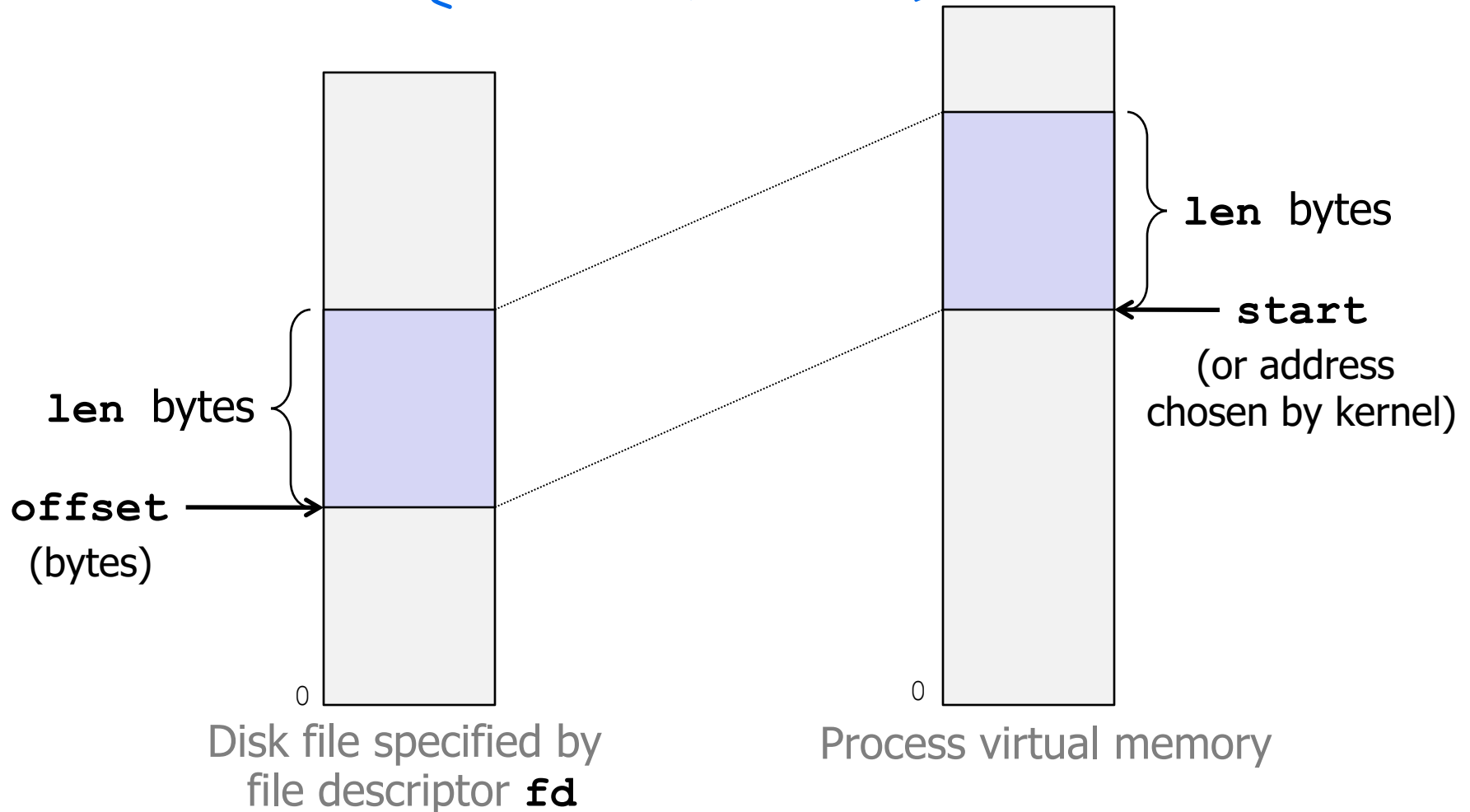
$\sigma$ : 18.2

high: 98



# User-Level Memory Mapping

```
void *mmap(void *start, int len,  
           int prot, int flags, int fd, int offset)
```



Nov 01, 21 1:13

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

```

\$ cat  
\$ cat bar

assert(fd != 0)

char buf[512];

int rc;  
int fd = open(argv[1], RD\_ONLY);  
while ((rc = read(fd, buf, sizeof(buf)))  
!= NULL) {

write(1, buf, rc);

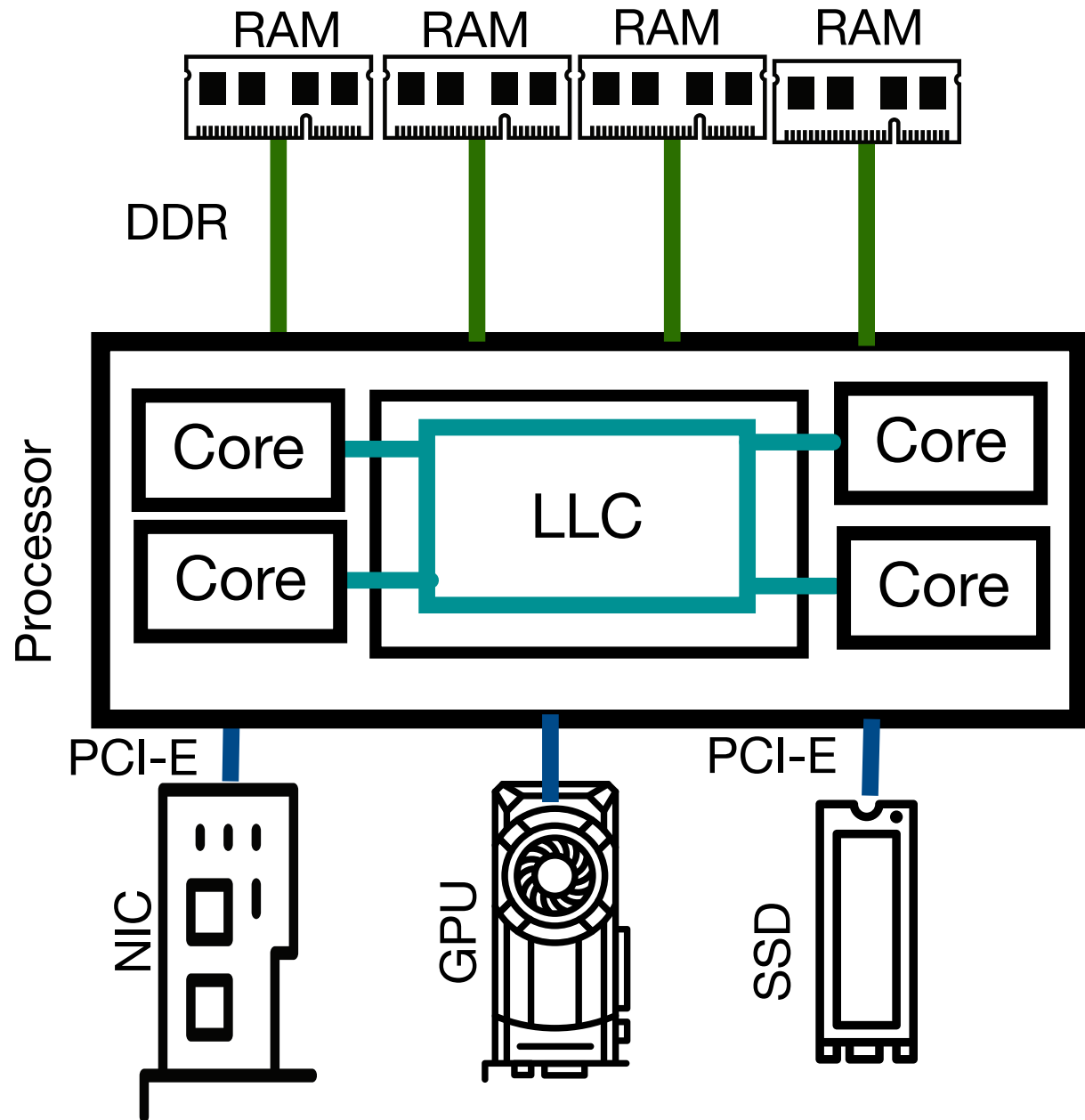
}

2 copies  
buffer cache → user-buf  
user-buf → buffer in front of terminal

1 memcopy:  
from buffer cache  
to buffer in front  
of terminal

writ  
sendfile();

# Machine





Nov 01, 21 1:12

handout09-2.txt

Page 1/5

```

1 CS 202, Fall 2021
2 Handout 9 (Class 15)
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    bootloader reads in the kernel from the disk.
12
13    See the functions boot_waitdisk() and boot_readsect(). Compare to Figures 36
14    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 510 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, PAGE_SIZE, PAGE_SIZE);
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    }

```

Nov 01, 21 1:12

handout09-2.txt

Page 2/5

```

73
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
93        // clear bss segment
94        for (; end_ptr < memsz; ++end_ptr) {
95            *(uint8_t*) end_ptr = 0;
96        }
97    }
98
99
100    // boot_waitdisk
101    // Wait for the disk to be ready.
102    static void boot_waitdisk(void) {
103        // Wait until the ATA status register says ready (0x40 is on)
104        // & not busy (0x80 is off)
105        while ((inb(0x1F7) & 0xC0) != 0x40) {
106            /* do nothing */
107        }
108    }
109
110
111    // boot_readsect(dst, src_sect)
112    // Read disk sector number 'src_sect' into address 'dst'.
113    static void boot_readsect(uintptr_t dst, uint32_t src_sect) {
114        // programmed I/O for "read sector"
115        boot_waitdisk();
116        outb(0x1F2, 1); // send 'count = 1' as an ATA argument
117        outb(0x1F3, src_sect); // send 'src_sect', the sector number
118        outb(0x1F4, src_sect >> 8);
119        outb(0x1F5, src_sect >> 16);
120        outb(0x1F6, (src_sect >> 24) | 0xE0);
121        outb(0x1F7, 0x20); // send the command: 0x20 = read sectors
122
123        // then move the data into memory
124        boot_waitdisk();
125        insl(0x1F0, (void*) dst, SECTORSIZE/4); // read 128 words from the disk
126    }
127
128

```

Nov 01, 21 1:12

handout09-2.txt

Page 3/5

```

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         ys
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
168     int ch = (unsigned char) keymap[data | escape];
169
170     if (ch >= 'a' && ch <= 'z') {
171         if (modifiers & MOD_CONTROL) {
172             ch -= 0x60;
173         } else if (!(modifiers & MOD_SHIFT) != !(modifiers & MOD_CAPSLOCK))
174         {
175             ch -= 0x20;
176         }
177     } else if (ch >= KEY_CAPSLOCK) {
178         modifiers ^= 1 << (ch - KEY_SHIFT);
179         ch = 0;
180     } else if (ch >= KEY_SHIFT) {
181         modifiers |= 1 << (ch - KEY_SHIFT);
182         ch = 0;
183     } else if (ch >= CKEY(0) && ch <= CKEY(21)) {
184         ch = complex_keymap[ch - CKEY(0)].map[modifiers & 3];
185     } else if (ch < 0x80 && (modifiers & MOD_CONTROL)) {
186         ch = 0;
187     }
188
189     return ch;
190 }

```

Nov 01, 21 1:12

handout09-2.txt

Page 4/5

```

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 in the upper left of
195 the 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
206     outb(0x3D4, 14); // Command 14 = upper byte of position
207     outb(0x3D5, 0 / 256); // row 0
208     outb(0x3D4, 15); // Command 15 = lower byte of position
209     outb(0x3D5, 0 % 256); // column 0
210
211 }
212
213
214

```

Nov 01, 21 1:12

handout09-2.txt

Page 5/5

## 215 3. Memory-mapped I/O

216

217 a. Here is a 32-bit PC's physical memory map:

218

219 +-----+ &lt;- 0xFFFFFFFF (4GB)

220 |

221 | 32-bit  
222 | memory mapped  
223 | devices

224 |

225 | /\/\/\/\/\/\/\/\/\

226 |

227 | /\/\/\/\/\/\/\/\/\

228 |

229 | Unused

230 |

231 | &lt;- depends on amount of RAM

232 |

233 | Extended Memory

234 |

235 | &lt;- 0x00100000 (1MB)

236 |

237 | BIOS ROM

238 |

239 | &lt;- 0x000F0000 (960KB)

240 |

241 | 16-bit devices,  
242 | expansion ROMs

243 |

244 | &lt;- 0x000C0000 (768KB)

245 |

246 | VGA Display

247 |

248 | &lt;- 0x000A0000 (640KB)

249 |

250 | Low Memory

251 |

252 | &lt;- 0x00000000

253 |

254 [Credit to Frans Kaashoek, Robert Morris, and Nickolai Zeldovich for  
255 this picture]

256

257

258 b. Loads and stores to the device memory "go to hardware".

259

260 An example is in the console printing code from WeensyOS. Here is an  
261 excerpt from link/shared.ld:

262

263 /\* Compare the address below to the map above. \*/

264 PROVIDE(console = 0xB8000);

265

266

267

268 /\*

269 \* prints a character to the console at the specified

270 \* cursor position in the specified color.

271 \* Question: what is going on in the check

272 \* if (c == '\n')

273 \* ?

274 \* Hint: '\n' is "C" for "newline" (the user pressed enter).

275 \*/

276 static void console\_putc(printer\* p, unsigned char c, int color) {

277 console\_printer\* cp = (console\_printer\*) p;

278 if (cp-&gt;cursor &gt;= console + CONSOLE\_ROWS \* CONSOLE\_COLUMNS) {

279 cp-&gt;cursor = console;

280 }

281 if (c == '\n') {

282 int pos = (cp-&gt;cursor - console) % 80;

283 for (; pos != 80; pos++) {

284 \*cp-&gt;cursor++ = ' ' | color;

285 }

286 }

287 \*cp-&gt;cursor++ = c | color;

288 }

289 }

290 }

291 }

292 }

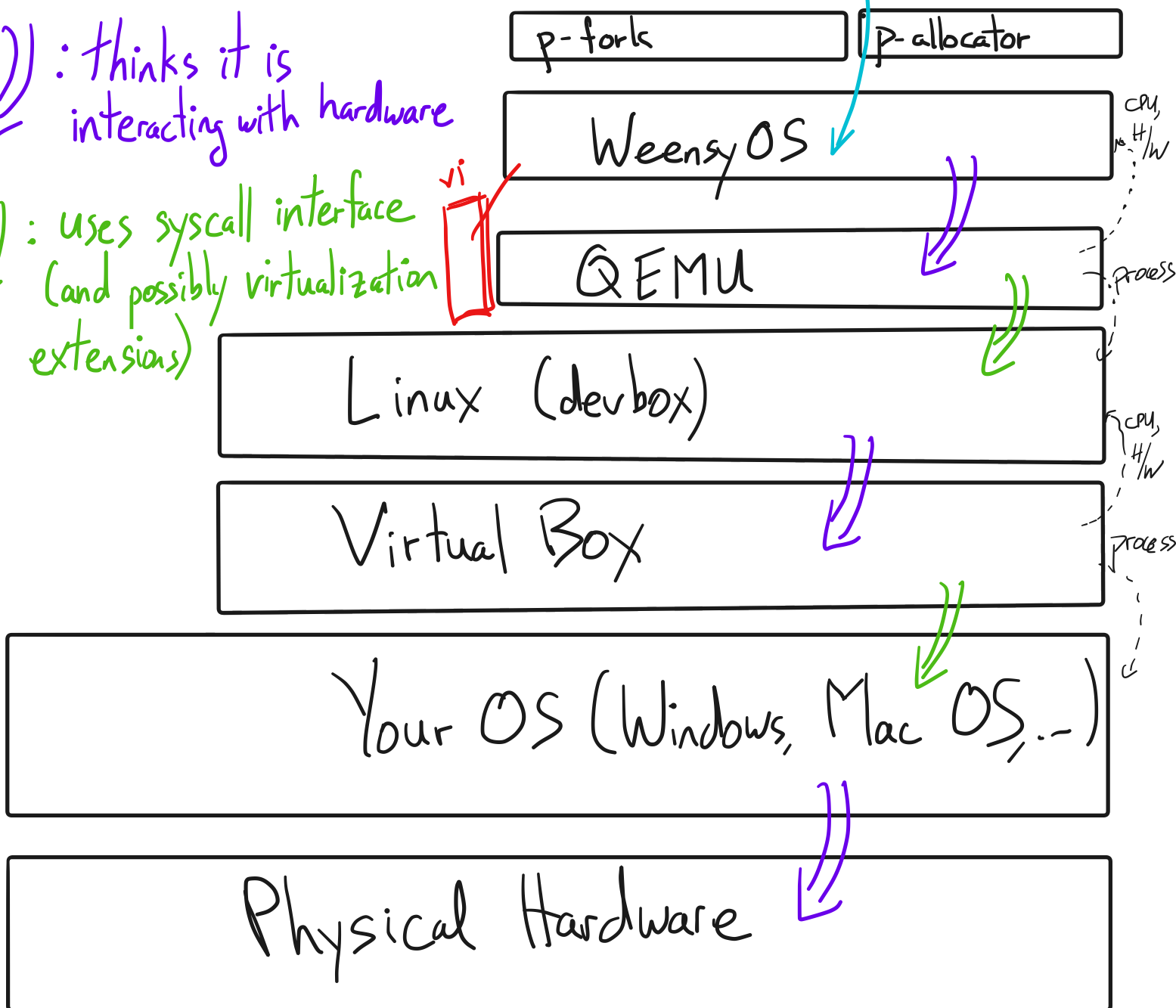
293 }

294 }

lab 4 work here

)) : thinks it is interacting with hardware

)) : uses syscall interface (and possibly virtualization extensions)



Hints:

- processes: files matching p-\*.c
- kernel code: files matching k-\*.c , k-\*.s
- system calls and returns `()` /\* cousin of `mmap()` \*/

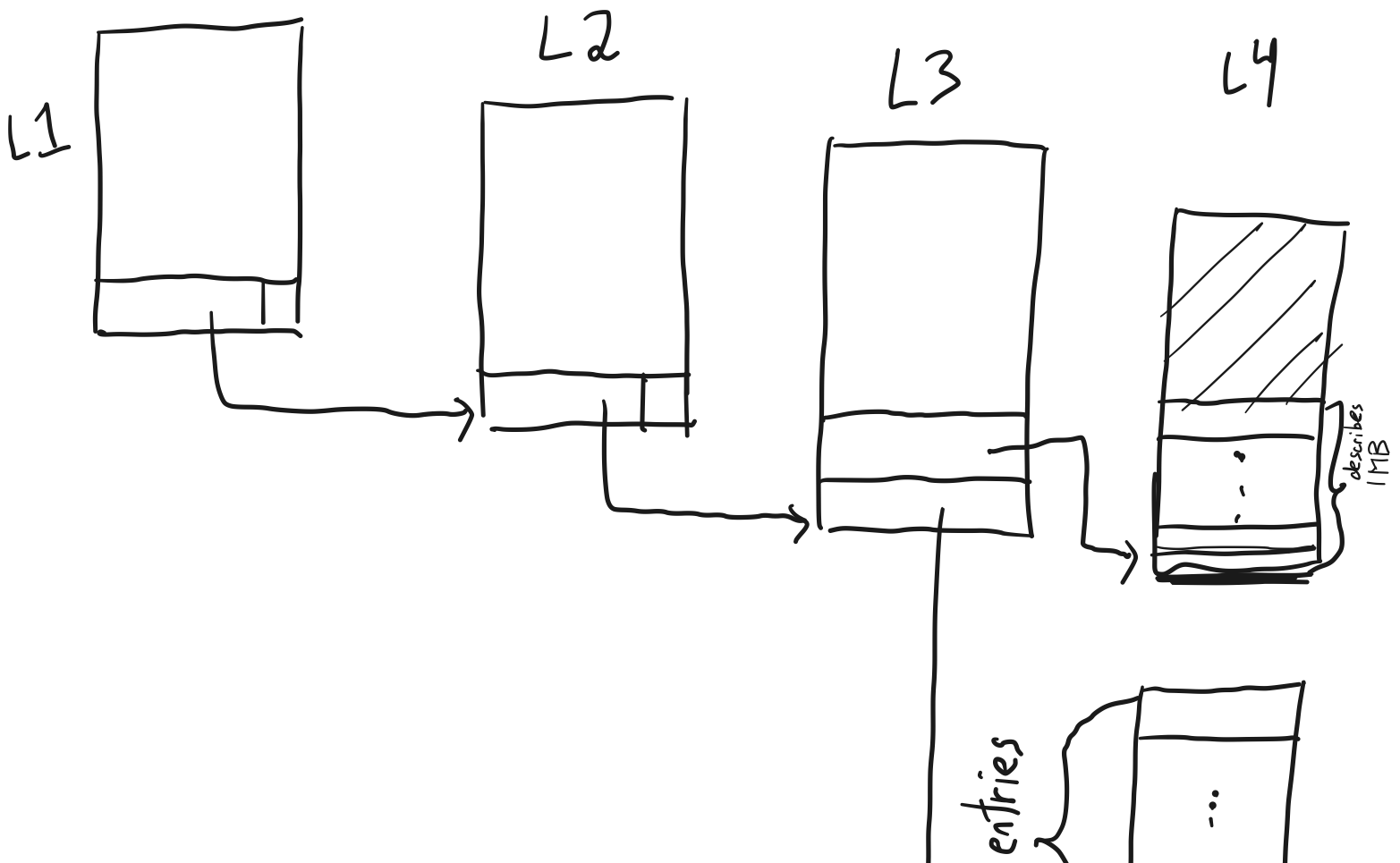
sys\_page\_alloc, lookin process.h

kernel returns: exception\_return();  
%rax contains return value of system call  
(errno, 0)

%rdi:  
arg to  
syscall

- you'll use virtual-memory-map() = (NULL vs. non-NULL)
- pay attention to the "allocator" argument
- make sure your allocator initializes the page table  
memset(addr, 0, len);

- a process's virtual address space: 3 MB. What's the page table structure?

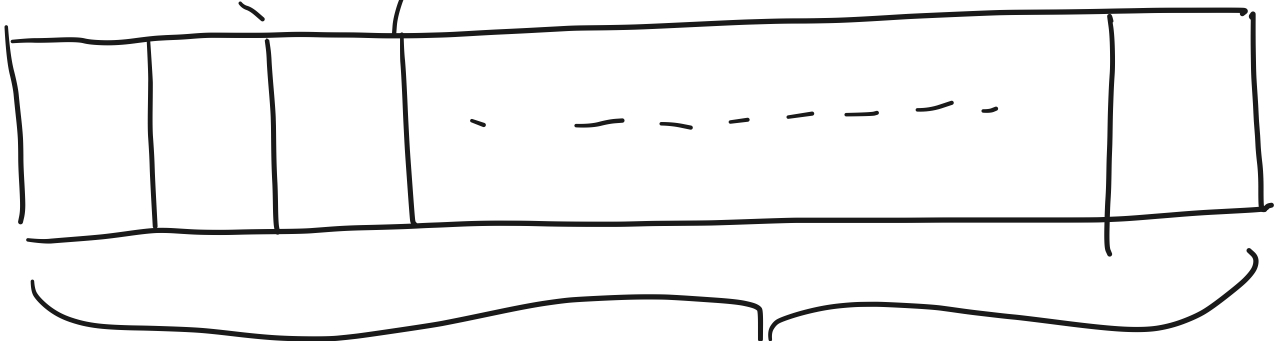




each entry  
describes a mapping  
for one page  
(4KB)

PCB  $\equiv$  struct proc (in kernel.h)

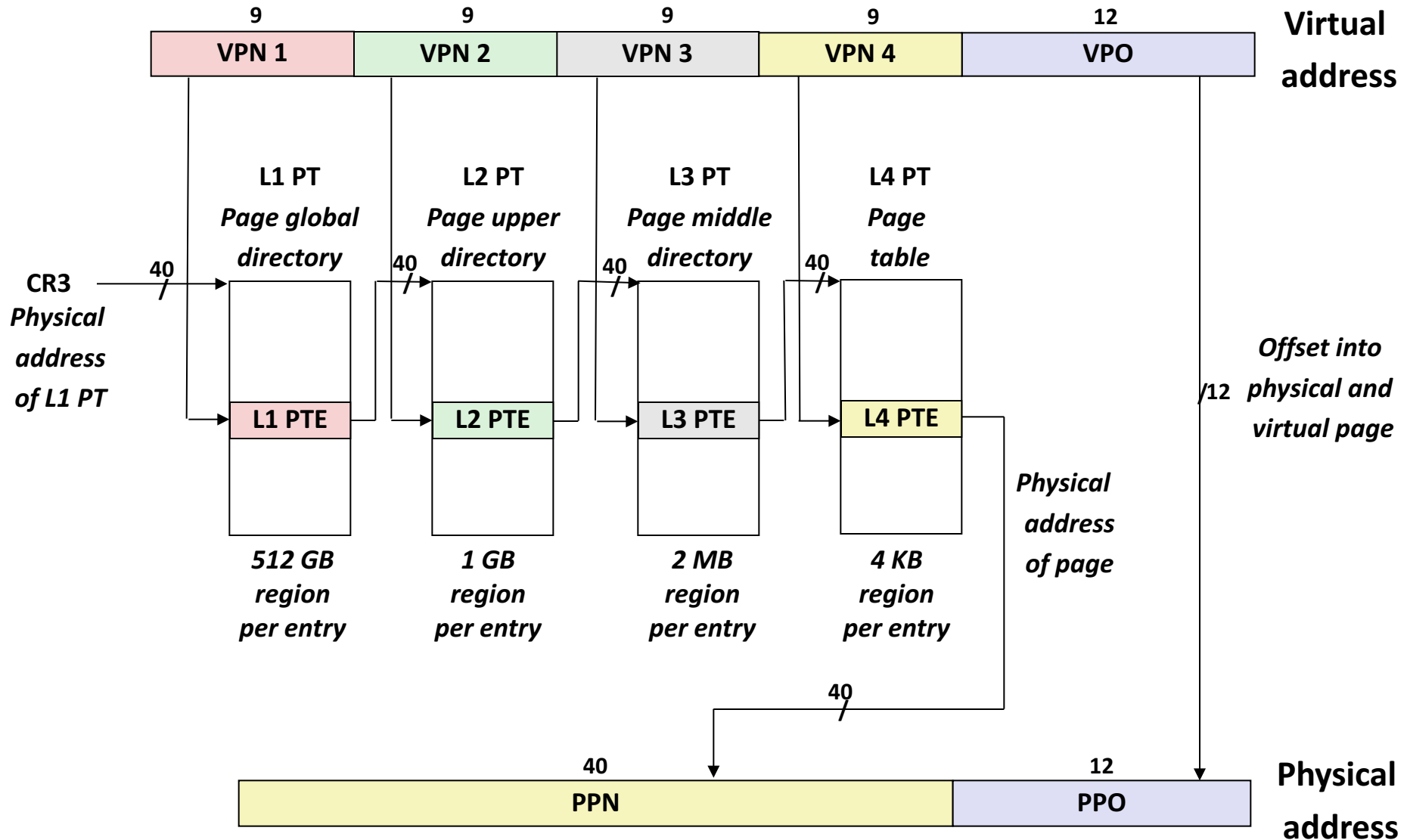
```
struct physical_pageinfo {  
    int8_t owner;  
    int8_t refcount;  
}
```



one per physical page

this is not a page table; it is bookkeeping.

# Core i7 Page Table Translation





# Review of Symbols

## ■ Basic Parameters

- $N = 2^n$  : Number of addresses in virtual address space
- $M = 2^m$  : Number of addresses in physical address space
- $P = 2^p$  : Page size (bytes)

## ■ Components of the virtual address (VA)

- TLBI: TLB index
- TLBT: TLB tag
- VPO: Virtual page offset
- VPN: Virtual page number

## ■ Components of the physical address (PA)

- PPO: Physical page offset (same as VPO)
- PPN: Physical page number
- CO: Byte offset within cache line
- CI: Cache index
- CT: Cache tag

# Core i7 Level 1-3 Page Table Entries

63	62	52	51	12	11	9	8	7	6	5	4	3	2	1	0
XD	Unused	Page table physical base address				Unused	G	PS		A	CD	WT	U/S	R/W	P=1
Available for OS															P=0

**Each entry references a 4K child page table. Significant fields:**

**P:** Child page table present in physical memory (1) or not (0).

**R/W:** Read-only or read-write access access permission for all reachable pages.

**U/S:** user or supervisor (kernel) mode access permission for all reachable pages.

**WT:** Write-through or write-back cache policy for the child page table.

**A:** Reference bit (set by MMU on reads and writes, cleared by software).

**PS:** Page size: if bit set, we have 2 MB or 1 GB pages (bit can be set in Level 2 and 3 PTEs only).

**Page table physical base address:** 40 most significant bits of physical page table address (forces page tables to be 4KB aligned)

**XD:** Disable or enable instruction fetches from all pages reachable from this PTE.

# Core i7 Level 4 Page Table Entries

63	62	52	51	12	11	9	8	7	6	5	4	3	2	1	0
XD	Unused	Page physical base address				Unused	G		D	A	CD	WT	U/S	R/W	P=1
Available for OS (for example, if page location on disk)															P=0

**Each entry references a 4K child page. Significant fields:**

**P:** Child page is present in memory (1) or not (0)

**R/W:** Read-only or read-write access permission for this page

**U/S:** User or supervisor mode access

**WT:** Write-through or write-back cache policy for this page

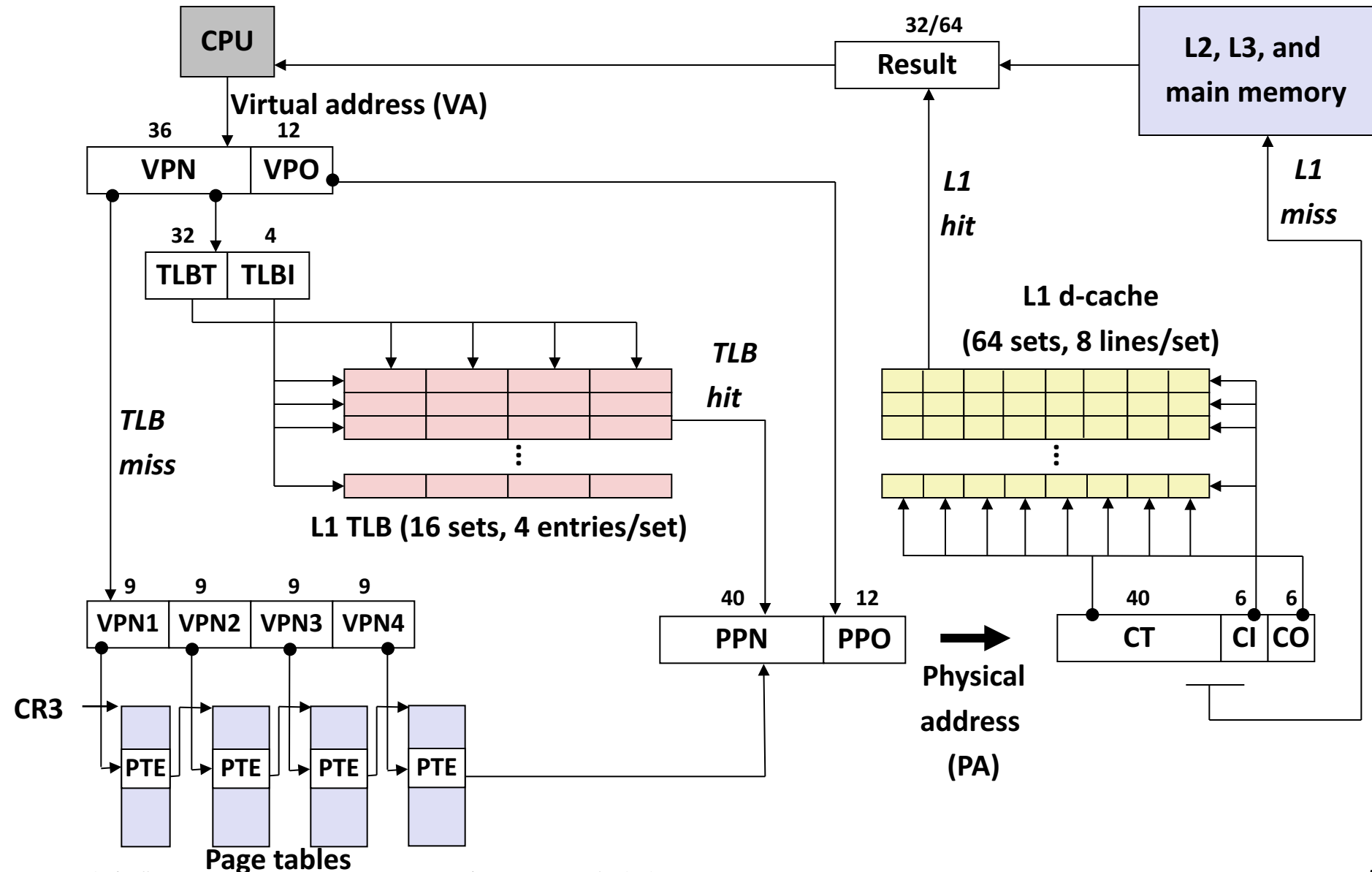
**A:** Reference bit (set by MMU on reads and writes, cleared by software)

**D:** Dirty bit (set by MMU on writes, cleared by software)

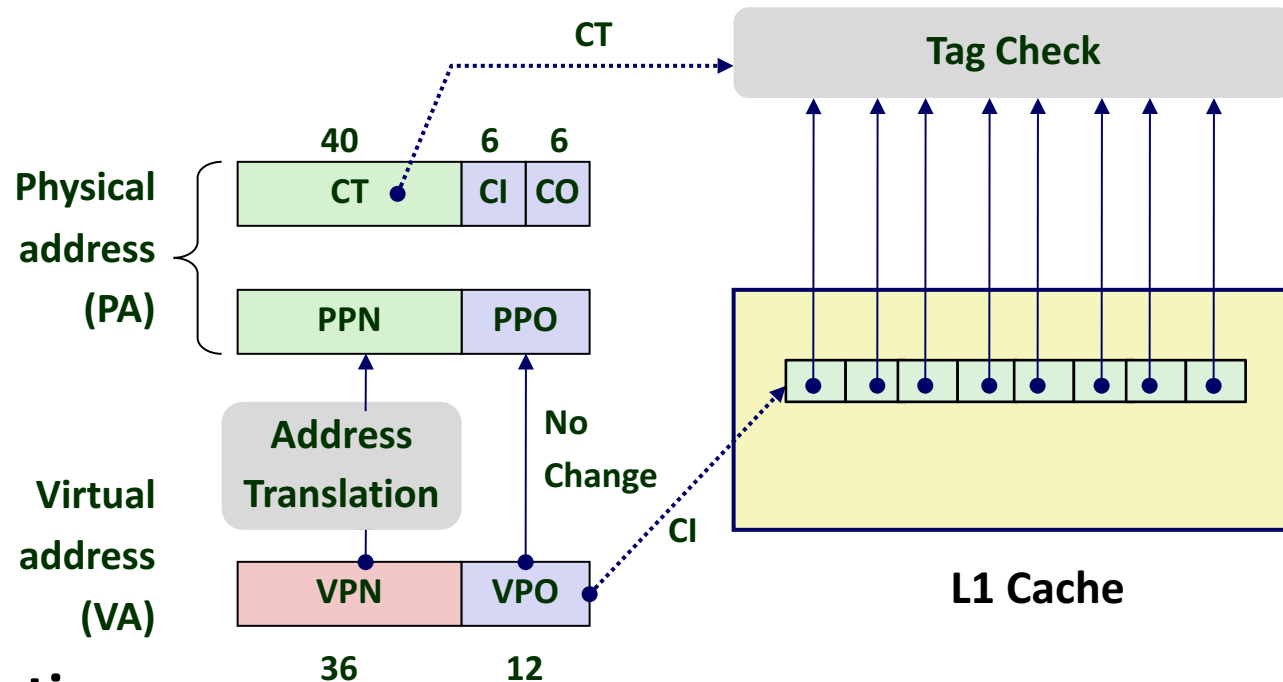
**Page physical base address:** 40 most significant bits of physical page address  
(forces pages to be 4KB aligned)

**XD:** Disable or enable instruction fetches from this page.

# End-to-end Core i7 Address Translation



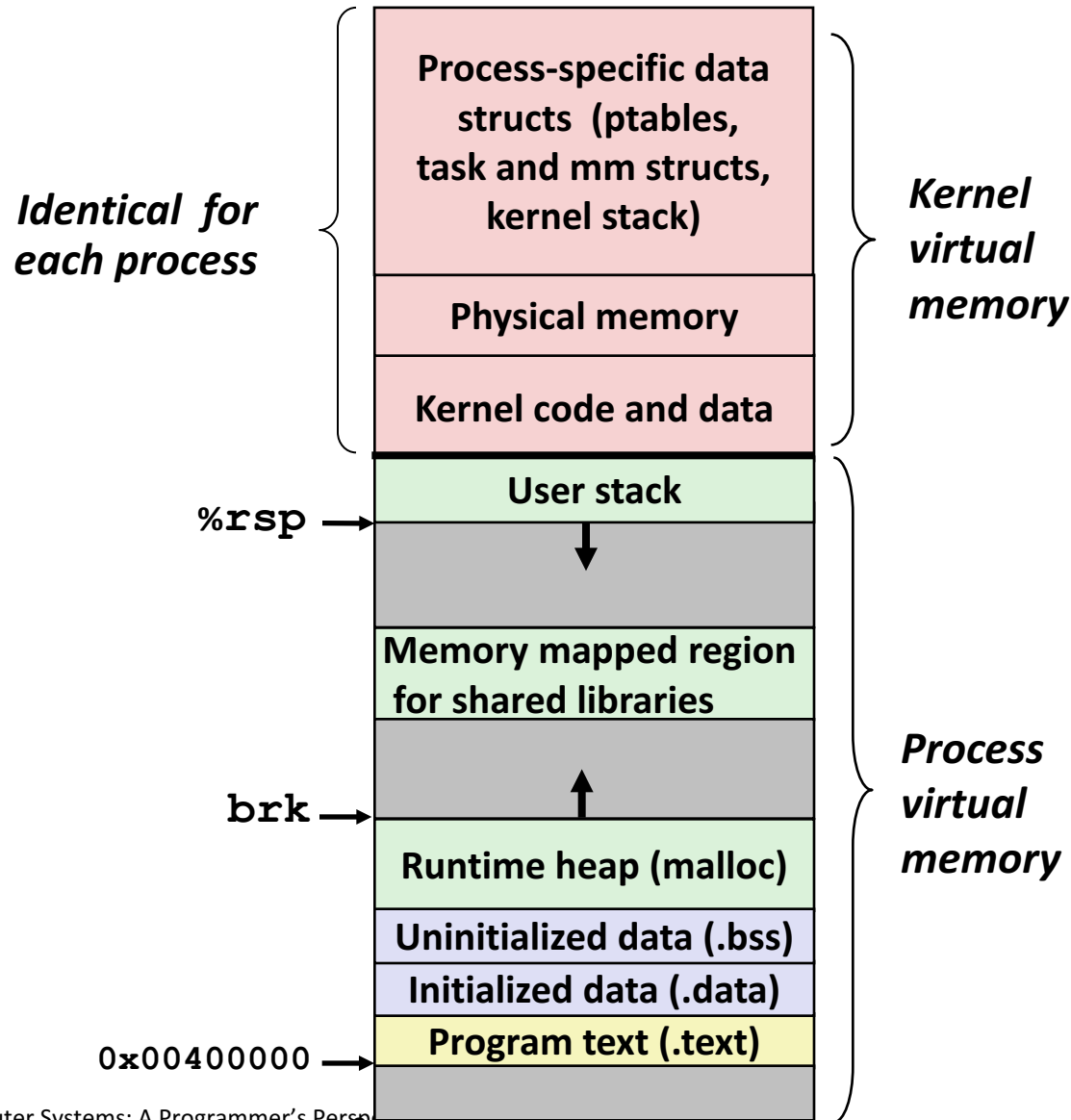
# Cute Trick for Speeding Up L1 Access



## ■ Observation

- Bits that determine CI identical in virtual and physical address
- Can index into cache while address translation taking place
- Cache carefully sized to make this possible: 64 sets, 64-byte cache blocks
- Means 6 bits for cache index, 6 for *cache* offset
- That's 12 bits; matches *VPO*, *PPO* → One reason pages are  $2^{12}$  bits = 4 KB

# Virtual Address Space of a Linux Process



melt down  
Spectre

31		15		5	4	3	2	1	0
Reserved		SGX	Reserved	PK	I/D	RSVD	U/S	W/R	P
P	0	The fault was caused by a non-present page.							
	1	The fault was caused by a page-level protection violation.							
W/R	0	The access causing the fault was a read.							
	1	The access causing the fault was a write.							
U/S	0	A supervisor-mode access caused the fault.							
	1	A user-mode access caused the fault.							
RSVD	0	The fault was not caused by reserved bit violation.							
	1	The fault was caused by a reserved bit set to 1 in some paging-structure entry.							
I/D	0	The fault was not caused by an instruction fetch.							
	1	The fault was caused by an instruction fetch.							
PK	0	The fault was not caused by protection keys.							
	1	There was a protection-key violation.							
SGX	0	The fault is not related to SGX.							
	1	The fault resulted from violation of SGX-specific access-control requirements.							

Figure 4-12. Page-Fault Error Code