# Last Class

① Page Faults
  ↳ What They Are
    → How They Occur
      → Information They Provide
        ↳ Faulting Address (%CR2)
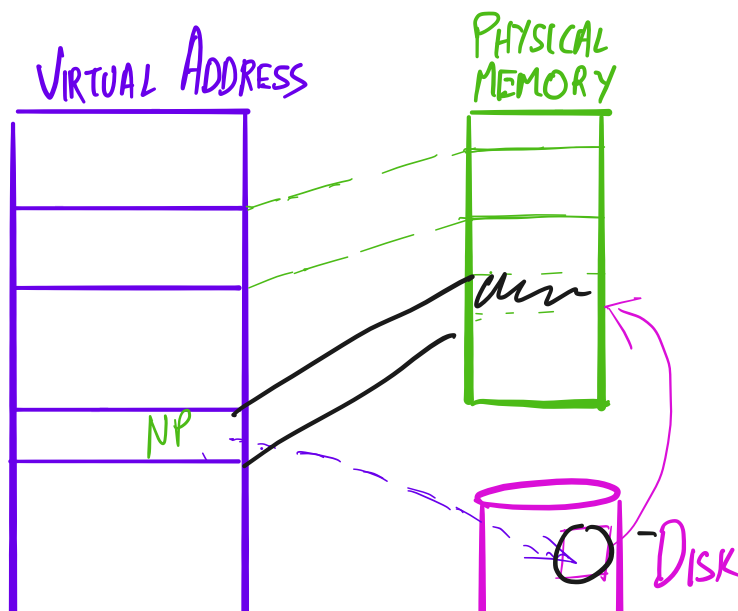          → Why The Fault Occurred

② Uses
    ① Copy-On-Write

# Today

Uses: <u>Demand Paging</u>

## Goal

make run

# Why?

① Loading Executables

⇑

② MMAP(2)

③ Using More Memory Than Physically Available

# How?
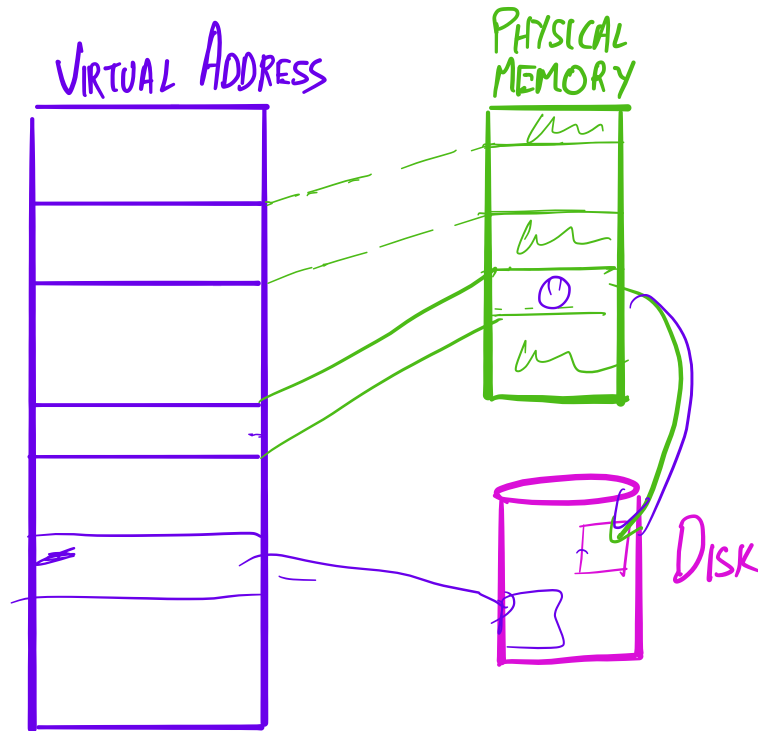
① PF Occurs
  ↳ Check If Page Should Be Loaded From Disk

② Allocate Physical Page
  Check If FREE PHYSICAL PAGES Are Available
  └→ Evict To Free Physical Pages
  No

YES → ALLOCATE PAGE

③ LOAD & MAP        /dev/zero

EVICTION

VIRTUAL ADDRESS          PHYSICAL MEMORY

DISK

TODAY'S QUESTION: WHAT PAGE TO EVICT?

PAGE REPLACEMENT POLICY

ASIDE: DOES THIS EVEN MATTER ANYMORE?

- CDNs

- DATABASES

# Policies

- FIFO : First In First Out

- Min (Optimal): Clairvoyant (Knows The Future)
  Evict Page That Won't Be Accessed
  For The Longest Time

## Example

3 Physical Pages

Access: ABC ABD ADB CB      (Page)

FIFO

| | A | B | C | A | B | D | A | D | B | C | B |
|---|---|---|---|---|---|---|---|---|---|---|---|
| P0 | A | | | | | h | | **D** | | h | **C** |
| P1 | | B | | | | h | | **A** | | | |
| P2 | | | C | | | | | | **B** | h | |

# Of Page-Ins / Misses      7 } ⁷/₁₁

# Of Hits      4

## MIN

|  | # OF HITS | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  | A | B | C | A | B | D | A | D | B | C | B |
| P0 | A |  |  | h |  |  | h |  |  | C |  |
| P1 |  | B |  |  | h |  |  | h |  | h | h |
| P2 |  |  | C |  |  | D |  | h |  |  |  |

# OF PAGE-INS / MISSES  5 ⎫
# OF HITS  6 ⎬ — 5/11

# WHY IT MATTERS

Avg MEMORY ACCESS TIME = $\boxed{P}$ · page fault time  ← Page-in/miss

+ $\boxed{(1-p)}$ · memory access latency
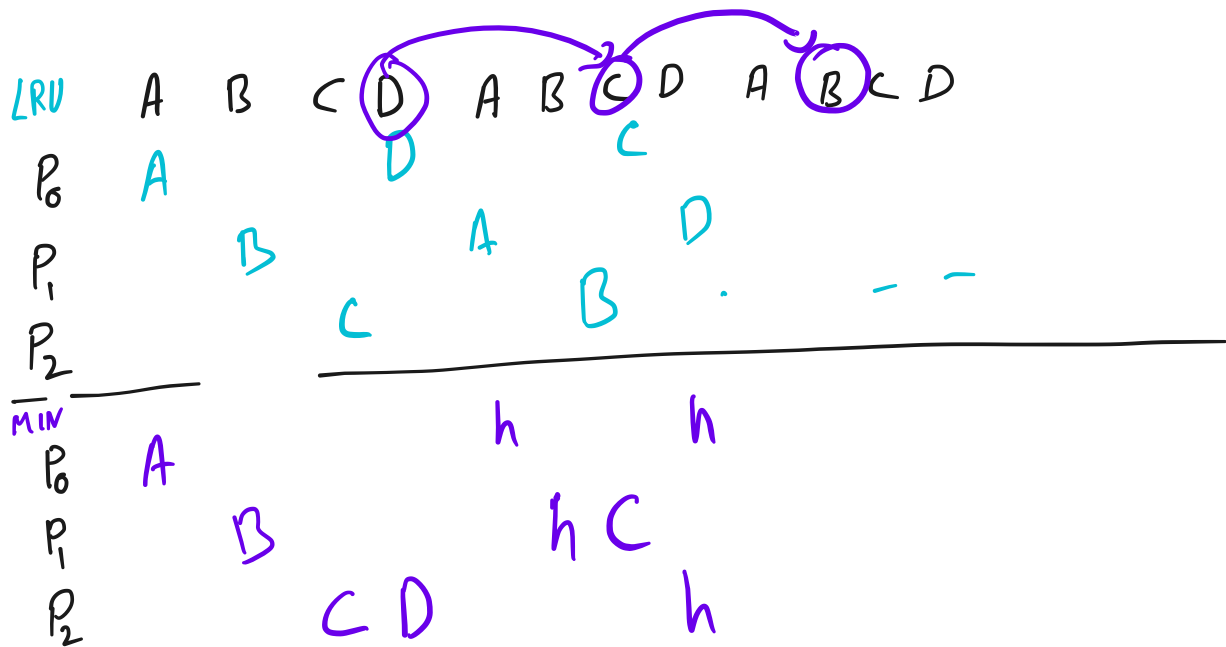
ASSUMPTION   Memory Access Lat. $\ll$ Page Fault Time
$(\sim 100\text{ns})$   $(\sim 10\text{ms} = 10^7\text{ns})$

# LRU: Least Recently Used

|  | A | B | C | A | B | D | A | D | B | C | B |
|---|---|---|---|---|---|---|---|---|---|---|---|
| P0 | A |  |  | h |  |  | h |  |  | C |  |
| P1 |  | B |  |  | h |  |  | h |  | h | h |
| P2 |  |  | C |  |  | D |  | h |  |  |  |

# OF PAGE-INS / MISSES  5
# OF HITS

# OF HITS                          6

| LRU | A | B | C | D | A | B | C | D | A | B | C | D |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $P_0$ | A | | | D | | | C | | | B | | |
| $P_1$ | | B | | | A | | | D | | | | |
| $P_2$ | | | C | | | B | | | | | | |

MIN

| | A | B | C | D | A | B | C | D | A | B | C | D |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $P_0$ | A | | | | | | h | | | h | | |
| $P_1$ | | B | | | | | h | C | | | | |
| $P_2$ | | | C | D | | | | h | | | | |

# TRADE-OFFS

- FIFO    **+ Simple**

    — High miss rates

Min
- ~~OPT~~    + Optimal

    — $\mathbb{P}$ know the future

- LRU    + ~~Emp~~ Empirically close to Min

    —

# APPROXIMATING LRU — CLOCK

## REFERENCE BIT (A)

# MISCELLANEOUS

- BELADI'S ANOMALY
  - ↳ FIFO

A B C D A B E A B C D E

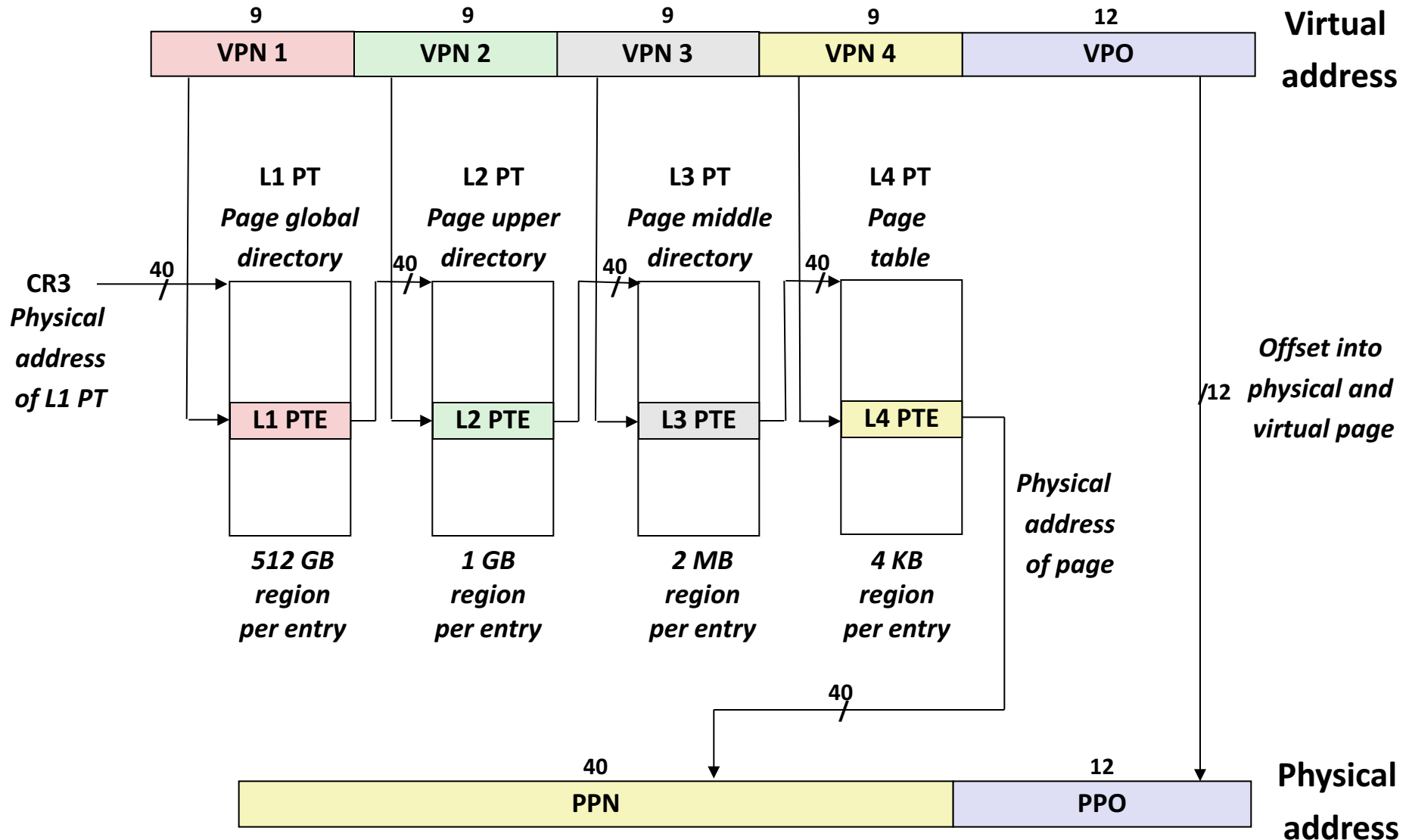$P_0$

$P_1$

$P_2$

---

$P_0$

$P_1$

$P_2$

$P_3$

3 Hit                                         4 Hit

# Thrashing

## Page Faults for Accounting

# Core i7 Page Table Translation

| 9 | 9 | 9 | 9 | 12 | **Virtual** |
|---|---|---|---|---|---|
| VPN 1 | VPN 2 | VPN 3 | VPN 4 | VPO | **address** |

**L1 PT**
*Page global directory*

**L2 PT**
*Page upper directory*

**L3 PT**
*Page middle directory*

**L4 PT**
*Page table*

CR3
*Physical address of L1 PT*

40

40

40

40

*Offset into physical and virtual page*

/12

| L1 PTE | | L2 PTE | | L3 PTE | | L4 PTE |

*Physical address of page*

*512 GB region per entry*

*1 GB region per entry*

*2 MB region per entry*

*4 KB region per entry*

40

| 40 | 12 | **Physical** |
|---|---|---|
| PPN | PPO | **address** |

# 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 |

| | | | | | | | | | | | | | | | | | | P=0 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| | Available for OS (for example, if page location on disk) | | | | | | | | | | | | | | | | | |

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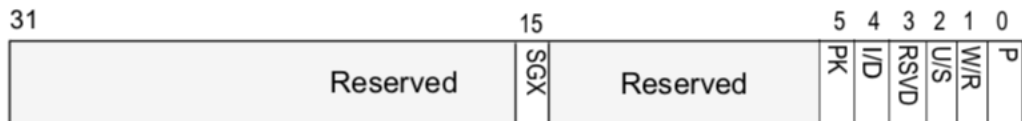
**U/S:** User or supervisor mode access

**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.

# 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.

```
31                                      15          5 4 3 2 1 0
                                            S          P I R U W P
          Reserved                      G  Reserved  K / S / /
                                            X          K D V S R
```

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**

# 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



Identical for each process

Process-specific data structs (ptables, task and mm structs, kernel stack)

Physical memory

Kernel code and data

Kernel virtual memory

%rsp →

User stack

Memory mapped region for shared libraries

brk →

Runtime heap (malloc)

Uninitialized data (.bss)

Initialized data (.data)

0x00400000 →

Program text (.text)

Process virtual memory

0

Code | Blame   810 lines (675 loc) · 27.9 KB    Raw   ↑ Top

```
120         set_gate(&interrupt_descriptors[INT_TIMER], X86GATE_INTERRUPT, 0,
121               (uint64_t) timer_int_handler);
122
123     // GPF and page fault
124         set_gate(&interrupt_descriptors[INT_GPF], X86GATE_INTERRUPT, 0,
125               (uint64_t) gpf_int_handler);
126         set_gate(&interrupt_descriptors[INT_PAGEFAULT], X86GATE_INTERRUPT, 0,
127               (uint64_t) pagefault_int_handler);
128
129     // System calls get special handling.
130     // Note that the last argument is '3'.  This means that unprivileged
131     // (level-3) applications may generate these interrupts.
132     for (unsigned i = INT_SYS; i < INT_SYS + 16; ++i) {
133         set_gate(&interrupt_descriptors[i], X86GATE_INTERRUPT, 3,
134               (uint64_t) sys_int_handlers[i - INT_SYS]);
135     }
136
137 x
138 i
139 i
140
141 /
142 a
143
144
145
146
147
148
149
150 /
151 u
152 c
153 l
154 }
155
156 // in
157 // in
```

Code | Blame   199 lines (168 loc) · 4.12 KB

```
133
134 generic_exception_handler:
135         pushq %gs
136         pushq %fs
137         pushq %r14
138         pushq %r13
139         pushq %r12
140         pushq %r11
141         pushq %r10
142         pushq %r9
143         pushq %r8
144         pushq %rdi
145         pushq %rsi
146         pushq %rbp
147         pushq %rbx
148         pushq %rdx
149         pushq %rcx
150         pushq %rax
151         pushq %rax
152         movq %rsp, %rdi
153         call exception
154     # `exception` should never return.
155
156
```

Code | Blame   199 lines (168 loc) · 4.12 KB    ↑ Top

```
26  2:      jmp kernel
27
28
29     # Interrupt handlers
30     .align 2
31
32         .globl gpf_int_handler
33 gpf_int_handler:
34         pushq $13              // trap number
35         jmp generic_exception_handler
36
37         .globl pagefault_int_handler
38 pagefault_int_handler:
39         pushq $14
40         jmp generic_exception_handler
```

func set_sys_segment
func set_gate

Code | Blame   575 lines (468 loc) · 18.2 KB    Raw ↑ Top

```
60     typedef enum pageowner {
167
168 void exception(x86_64_registers* reg) {
169     // Copy the saved registers into the `current` process descriptor
170     // and always use the kernel's page table.
171     current->p_registers = *reg;
172     set_pagetable(kernel_pagetable);
173
174     // It can be useful to log events using `log_printf`.
175     // Events logged this way are stored in the host's `log.txt` file.
176     /*log_printf("proc %d: exception %d\n", current->p_pid, reg->reg_intno);*/
177
178     // Show the current cursor location and memory state
179     // (unless this is a kernel fault).
180     console_show_cursor(cursorpos);
181     if (reg->reg_intno != INT_PAGEFAULT || (reg->reg_err & PFERR_USER)) {
182         check_virtual_memory();
183         memshow_physical();
184         memshow_virtual_animate();
185     }
186 #if TICK_LIMIT
187     if (ticks == TICK_LIMIT) {
188         poweroff();
189     }
```