---
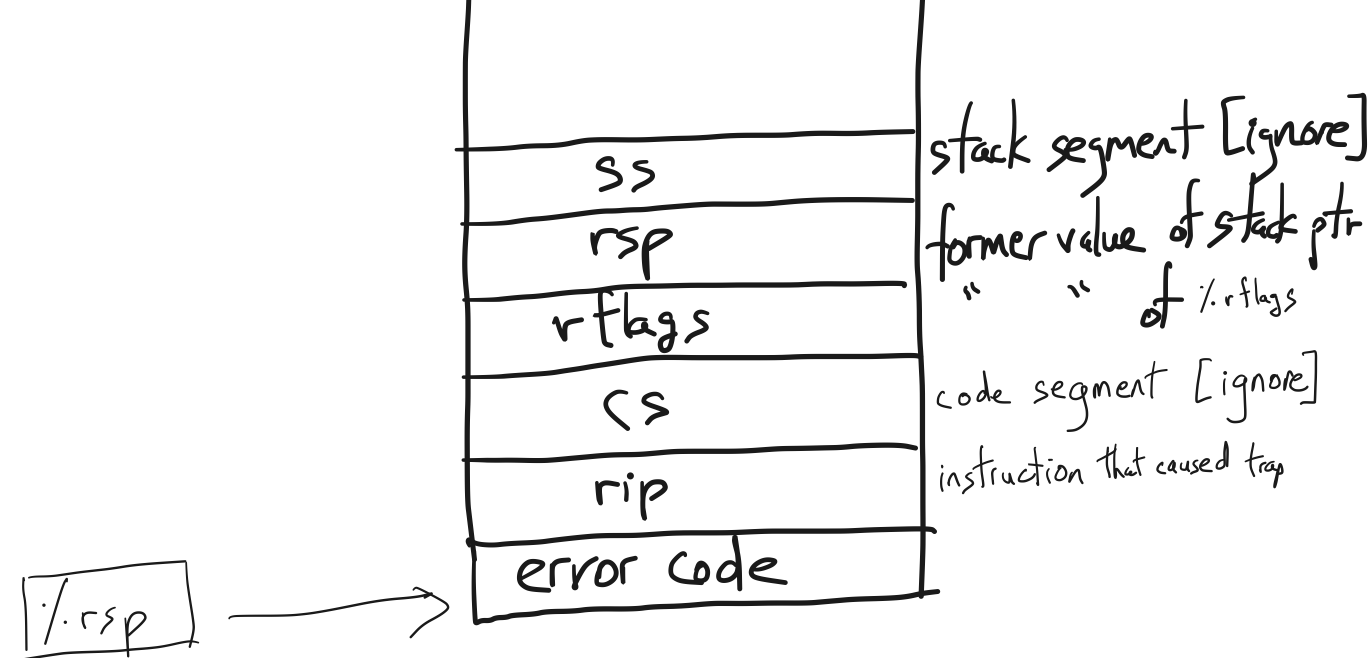
## 2. Page faults: intro + mechanics

Concept: illegal virtual memory reference:

<u>hardware</u> thinks it's illegal (though it might be valid for the process)

OS has to get involved

Mechanics:

- processor constructs <u>trap frame</u> and transfers execution to an interrupt or trap handler

| | |
|---|---|
| ss | stack segment [ignore] |
| rsp | former value of stack ptr |
| rflags | " " of %.rflags |
| cs | code segment [ignore] |
| rip | instruction that caused trap |
| error code | |

%.rsp $\longrightarrow$

%.rip $\longrightarrow$ code to handle the trap

error code : see last pg of handout

$$[ \_ \_ \_ \_ \_ \overset{I}{b} \_ \overset{U}{S} \overset{W}{R} \overset{P}{} ]$$
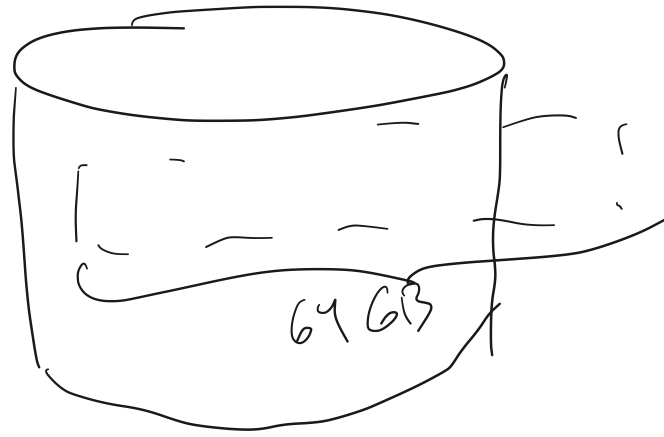
%.cr2 holds faulting virtual address

intent : on pg fault, kernel sets up the process's
page entries properly, or kills process.

3. Uses of page faults

- Classic example : overcommitting physical memory

prog: 64 GB          H/w: 16 GB

16GB          RAM

64 GB

- Copy on write          8 GB          cop1, mark
                                        all mem read-only
  fork()

- Accounting

- Store memory across the network (DSM)

machine

ABC

- Paging in day to day use

  - demand paging

  - growing the stack

  - BSS page allocation

BSS

data

text/code

- Shared text (code)

- Shared libraries

- Shared memory

Randy Pausch
time management

4. Page faults: costs

look at AMAT (avg. memory access time)

$$\text{AMAT} = (1-p) * (\text{mem access time}) + p * (\text{page fault time})$$

$t_M$ (under mem access time)     $t_D$ (under page fault time)

$p$ is probability (or frequency) of a page fault.

mem access time ~ 100ns $t_M$

disk access time ~ 10ms $= 10^7$ns $t_D$

QUESTION: what is p s.t. paging hurts performance by less than 10%?

w/out paging: $t_M$

w/ paging: $1.1 * t_M = (1-p) * t_M + p * t_D$

$.1 * t_M = (p)(-t_M) + p * t_D$

$$p = \frac{.1 * t_M}{t_D - t_M} = \frac{10 \text{ ns}}{10^7 \text{ ns}} = \frac{1}{10^6} = 10^{-6}$$

$$p = 10^{-6} \sim \frac{1}{1,000,000}$$

5. Page replacement policies

by VANTZ

RAM

PPN=46

$P_x$, VPN 12     $P_y$, VPN 23     $P_x$ VPN 14

- FIFO: eject oldest

- MIN (OPT): eject entry that won't be referenced
  for the <u>longest</u> time

  input:
    reference string
    cache size
  output:
    number of evictions

FIFO

| phys_slot | A | B | C | A | B | D | A | D | B | C | B |
|---|---|---|---|---|---|---|---|---|---|---|---|
| S1 | A |   |   | h |   | D | h |   |   | c |   |
| S2 |   | B |   |   | h |   | A |   |   |   |   |
| S3 |   |   | C |   |   |   |   |   | B |   | h |

7 evictions, 4 hits

## OPTIMAL

| phys_slot | A | B | C | A | B | D | A | D | B | C | B |
|---|---|---|---|---|---|---|---|---|---|---|---|
| S1 | A |   |   |   | h |   | h |   |   | c |   |
| S2 |   | B |   |   |   | h |   |   | h |   | h |
| S3 |   |   | C |   |   | D |   | h |   |   |   |

5 evictions, 6 hits

## LRU

| phys_slot | A | B | C | A | B | D | A | D | B | C | B |
|---|---|---|---|---|---|---|---|---|---|---|---|
| S1 | A |   |   | h |   | h |   |   |   | C |   |
| S2 |   | B |   |   | h |   |   | h |   |   | h |
| S3 |   |   | C |   |   |   |   | D |   |   |   |

5 evictions, 6 hits

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

phys-slot

S1   A        D        A        D        C
S2      B        A        D        B        C
S3         C        B        A        D

11 evictions, 0 hits

- - - - - - - - - - - - - -

FIFO
3entries     A  B  C  D  A  B  E  A  B  C  D  E
phys-slot
  S1
  S2
  S3

4entries     A  B  C  D  A  B  E  A  B  C  D  E
phys-slot
  S1
  S2
  S3
  S4

worse

Belady's anomaly

OPT ~ LRU

We are motivated to implement LRU

CLOCK

# 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

L1 PTE

L2 PTE

L3 PTE

L4 PTE

*512 GB region per entry*

*1 GB region per entry*

*2 MB region per entry*

*4 KB region per entry*

*Physical address of page*

*Offset into physical and virtual page*

/12

40

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

# 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        Page physical base address        12 | 11        Unused        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



CPU

32/64
Result

L2, L3, and
main memory

Virtual address (VA)

36
VPN

12
VPO

L1
hit

L1
miss

32
TLBT

4
TLBI

TLB
hit

L1 d-cache
(64 sets, 8 lines/set)

TLB
miss

L1 TLB (16 sets, 4 entries/set)

9
VPN1

9
VPN2

9
VPN3

9
VPN4

40
PPN

12
PPO

40
CT

6
CI

6
CO

CR3

PTE    PTE    PTE    PTE

Physical
address
(PA)

Page tables

# 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

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

**Physical memory**

**Kernel code and data**

*Identical for each process*

*Kernel virtual memory*

**User stack**

`%rsp` →

**Memory mapped region for shared libraries**

`brk` →

**Runtime heap (malloc)**

**Uninitialized data (.bss)**

**Initialized data (.data)**

`0x00400000` →

**Program text (.text)**

*Process virtual memory*

**0**

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

| 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