# ECE391 Computer System Engineering Lecture 15

Dr. Zbigniew Kalbarczyk
University of Illinois at Urbana- Champaign

Spring 2021

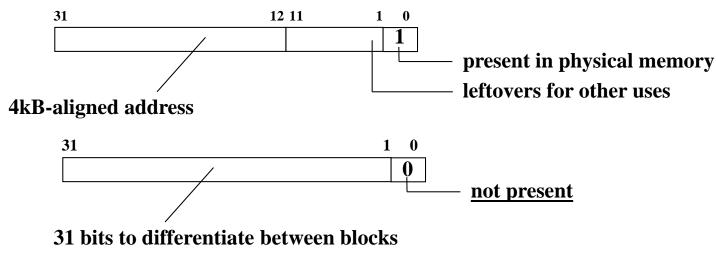
#### Lecture Topics

- x86 support for VM
  - protection model
  - segmentation
  - paging

Paging is the second level of indirection in x86

- Each page of a virtual address space is in one of three states
  - doesn't exist
  - exists and is in physical memory
  - exists, but is now on the disk rather than in memory

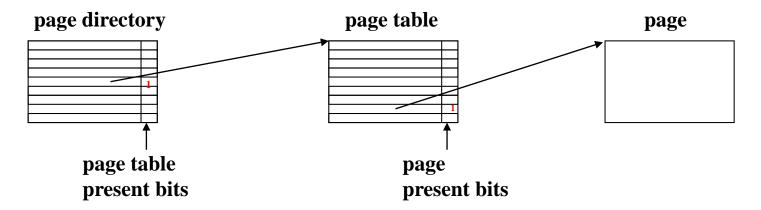
We can encode these possibilities as follows using 4B



- on disk & blocks that don't exist
- These 4B are called a page table entry (PTE); a group of them is a page table
- Question: if we use 4B for every 4kB, how big is the page table for a single virtual address space?

$$(4 / 4096) \times 2^{32} = 4MB$$
 too big...

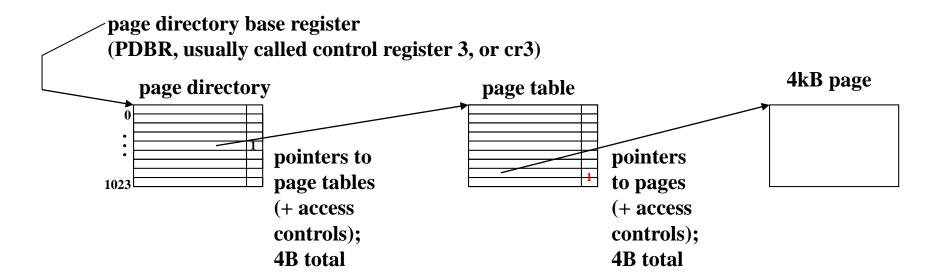
- Solution?
  - page the page table
  - i.e., use a hierarchical structure



- The page table is just another 4kB page
  - it holds 4096 / 4 = 1024 PTEs

- What about the page directory?
  - given
    - 2<sup>32</sup> bytes total (32-bit address space)
    - 2<sup>10</sup> PTEs per table
    - 2<sup>12</sup> bytes per page
  - the page directory needs  $2^{32}/(2^{10} \times 2^{12}) = 2^{10}$  entries total
  - which also fits in one page
     (and <u>could</u> be paged out to disk, although Linux does not)

31	22 21	12	11
directory #		page#	offset



- To translate a virtual address into a physical address
  - start with the PDBR (cr3)
  - look up page directory entry (PDE) using the 10 MSb of virtual address
  - extract page table address from PDE
  - look up page table entry (PTE) using next 10 bits of virtual address
  - extract page address from PTE
  - use last 12 bits of virtual address as offset into page

Way too slow to do on every memory access!

- Hence the translation lookaside buffers (TLBs)
  - keep translations of first 20 bits around and reuse them
  - only walk tables when necessary (in x86, OS manages tables, but hardware walks them)
  - TLBs flushed when cr3 is reloaded

- Remember the 11 free bits in the PTEs?
- What should we use them to do?
  - protect
  - optimize to improve performance

- Protect
  - User/Supervisor (U/S) page or page table
    - User means accessible to anyone (any privilege level)
    - Supervisor requires PL < 3 (i.e., MAX (CPL,RPL) < 3)</li>
  - Read-Only or Read/Write

#### Optimize

- TLBs must be fast, so you can't use many (~32 or 64)
- nice if
  - some translations are the same for all programs
  - bigger translations could be used when possible
     (e.g., use one translation for 4MB rather than 1024 translations)

#### x86 supports both

#### G flag—global

- TLB not flushed when changing to new program or address space (i.e., when cr3 changes)
- used for kernel pages (in Linux)

#### 4MB pages

- skip the second level of translation
- indicated by PS (page size) bit in PDE
- PS=1 means that the PDE points directly to a 4MB page
- remaining 22 bits of virtual address used as offset
- x86 provides separate TLBs for 4kB & 4MB translations

#### Important Bit Settings to Enable/Use Paging

- Paging is optional and must be enabled for use
  - PG (page enable): bit 31 in CR0 (control register 0)

- Allow mixed page sizes (i.e., pages of size 4kB and 4MB)
  - PSE (page size extension) bit 4 in CR4 (control register 4)

- Allow 4MB pages
  - PS (page size), bit 7 in PDE (page directory entry) must be set to 1

#### MP3: System Initialization (example)

```
gdt:
# x86_desc.S - Set up x86 segment descriptors, descriptor tables
                                                                        _gdt:
#define ASM 1
                                                                        # First GDT entry cannot be used
#include "x86 desc.h"
                                                                        .quad 0
                                                                        # NULL entry
                                                                        .quad 0
tss size:
                                                                        . . . . . . . . . . . . . . . .
.long tss_bottom - tss - 1
                                                                        . gdt_bottom:
ldt size:
                                                                        ldt:
.long ldt_bottom - ldt - 1
                                                                        .rept 4
                                                                        .quad 0
ldt_desc:
                                                                        .endr
.word KERNEL LDT
                                                                        ldt bottom:
.long ldt
                                                                        idt_desc_ptr:
gdt_desc:
                                                                        .word idt_bottom - idt - 1
// <gdt size>
                                                                        .long idt
// <gdt address>
                                                                        .align 16
                                                                        idt:
tss:
                                                                        idt:
tss:
.rept 104
                                                                        .rept NUM_VEC
                                                                        .quad 0
.byte 0
.endr
                                                                        .endr
tss_bottom:
                                                                        idt_bottom:
```

#### MP3: System Initialization (example)

```
# boot.S - start point for the kernel after GRUB gives us control
             1#include "multiboot.h"
#define ASM
#include "x86_desc.h"
# Load the GDT
Igdt (????)
```