## 15-213
### *"The course that gives CMU its Zip!"*

## P6 / Linux Memory System
## October 23, 2003

**Topics**
- P6 address translation
- Linux memory management
- Linux page fault handling
- Memory mapping

`class18.ppt`

---

## Intel P6
### (Bob Collwel's Chip, CMU Alumni)

**Internal Designation for Successor to Pentium**
- Which had internal designation P5

**Fundamentally Different from Pentium**
- Out-of-order, superscalar operation
- Designed to handle server applications
  - Requires high performance memory system

**Resulting Processors**
- PentiumPro (1996)
- Pentium II (1997)
  - Incorporated MMX instructions
    - » special instructions for parallel processing
  - L2 cache on same chip
- Pentium III (1999)
  - Incorporated Streaming SIMD Extensions
    - » More instructions for parallel processing

---

## P6 Memory System

```
      DRAM
       ▲
       │     external
       │     system bus
       ▼     (e.g. PCI)
   ┌─────────────────────────────┐
   │        L2                   │
   │       cache                 │
   │         ▲                   │
   │         │ cache bus         │
   │         ▼                   │
   │  bus interface unit    inst │
   │                        TLB  │
   │                             │
   │                        data │
   │  instruction    L1     TLB  │
   │  fetch unit   i-cache       │
   │                         L1  │
   │                       d-cache│
   └─────────────────────────────┘
     processor package
```

**32 bit address space**

**4 KB page size**

**L1, L2, and TLBs**
- 4-way set associative

**inst TLB**
- 32 entries
- 8 sets

**data TLB**
- 64 entries
- 16 sets

**L1 i-cache and d-cache**
- 16 KB
- 32 B line size
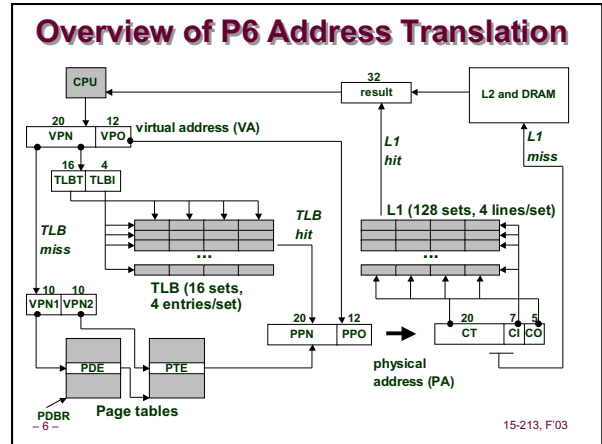- 128 sets

**L2 cache**
- unified
- 128 KB -- 2 MB
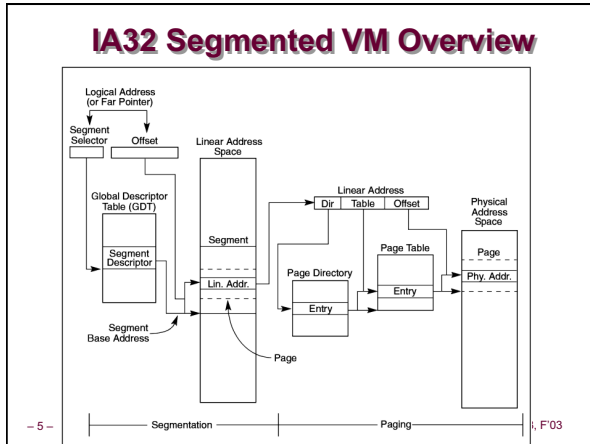
---

## Review of Abbreviations

**Symbols:**
- **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

---

## IA32 Segmented VM Overview

Logical Address
(or Far Pointer)

Segment Selector | Offset

Global Descriptor Table (GDT)

Segment Descriptor

Segment Base Address

Linear Address Space

Segment

Lin. Addr.

Page

Linear Address

Dir | Table | Offset

Page Directory

Entry

Page Table

Entry

Physical Address Space

Page

Phy. Addr.

| Segmentation | Paging |

i, F'03

---

## Overview of P6 Address Translation

CPU

20 VPN | 12 VPO  virtual address (VA)

16 TLBT | 4 TLBI

*TLB miss*

TLB (16 sets, 4 entries/set)

*TLB hit*

10 VPN1 | 10 VPN2

PDE | PTE

PDBR  Page tables

20 PPN | 12 PPO

physical address (PA)

32 result

*L1 hit*

*L1 miss*

L2 and DRAM

L1 (128 sets, 4 lines/set)

20 CT | 7 CI | 5 CO

---

## P6 2-level Page Table Structure

### Page directory
- 1024 4-byte page directory entries (PDEs) that point to page tables
- one page directory per process.
- page directory must be in memory when its process is running
- always pointed to by PDBR

### Page tables:
- 1024 4-byte page table entries (PTEs) that point to pages.
- page tables can be paged in and out.

page directory
1024 PDEs

Up to 1024 page tables

1024 PTEs
...
1024 PTEs
...
1024 PTEs

---

## P6 Page Directory Entry (PDE)

| 31 ... 12 | 11 ... 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| Page table physical base addr | Avail | G | PS | | A | CD | WT | U/S | R/W | P=1 |

**Page table physical base address**: 20 most significant bits of physical page table address (forces page tables to be 4KB aligned)

**Avail**: These bits available for system programmers

**G**: global page (don't evict from TLB on task switch)

**PS**: page size 4K (0) or 4M (1)

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

**CD**: cache disabled (1) or enabled (0)

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

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

**R/W**: read-only or read-write access

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

| 31 ... 1 | 0 |
|---|---|
| Available for OS (page table location in secondary storage) | P=0 |

## P6 Page Table Entry (PTE)

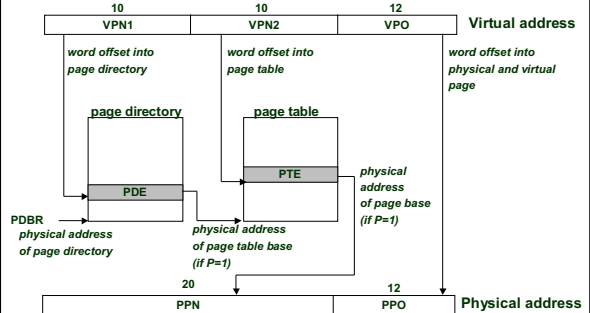| 31 | 12 | 11 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Page physical base address | | Avail | | G | 0 | D | A | CD | WT | U/S | R/W | P=1 |

**Page base address**: 20 most significant bits of physical page address (forces pages to be 4 KB aligned)
**Avail**: available for system programmers
**G**: global page (don't evict from TLB on task switch)
**D**: dirty (set by MMU on writes)
**A**: accessed (set by MMU on reads and writes)
**CD**: cache disabled or enabled
**WT**: write-through or write-back cache policy for this page
**U/S**: user/supervisor
**R/W**: read/write
**P**: page is present in physical memory (1) or not (0)

| 31 | 1 | 0 |
|---|---|---|
| Available for OS (page location in secondary storage) | | P=0 |

– 9 –                                                                 15-213, F'03

## How P6 Page Tables Map Virtual Addresses to Physical Ones



– 10 –                                                                 15-213, F'03

## 4Mbyte PDE's

**Page-Directory Entry (4-MByte Page)**

| 31 | 22 | 21 | 13 | 12 | 11 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Page Base Address | | Reserved | | PAT | Avail. | | G | PS | D | A | PCD | PWT | U/S | R/W | P |

Page Table Attribute Index
Available for system programmer's use
Global page
Page size (1 indicates 4 MBytes)
Dirty
Accessed
Cache disabled
Write-through
User/Supervisor
Read/Write
Present

– 11 –                                                                 15-213, F'03

## Support for 4Mbyte Pages

Linear Address



1024 PDE = 1024 Pages

*32 bits aligned onto a 4-KByte boundary.

– 12 –                                                                 15-213, F'03

Page 3

## Representation of VM Address Space

PT 3

Page Directory

PT 2

PT 0

| | | |
|---|---|---|
| P=1, M=1 | ● | |
| P=0, M=0 | ● | |
| P=1, M=1 | ● | |
| P=0, M=1 | ● | |

Page Directory
| | |
|---|---|
| P=1, M=1 | ● |
| P=1, M=1 | ● |
| P=0, M=0 | ● |
| P=0, M=1 | ● |

| | |
|---|---|
| P=1, M=1 | ● |
| P=0, M=0 | ● |
| P=1, M=1 | ● |
| P=0, M=1 | ● |

| | |
|---|---|
| P=0, M=1 | ● |
| P=0, M=1 | ● |
| P=0, M=0 | ● |
| P=0, M=0 | ● |

Page 15
Page 14
Page 13
Page 12
Page 11
Page 10
Page 9
Page 8
Page 7
Page 6
Page 5
Page 4
Page 3
Page 2
Page 1
Page 0

→ Mem Addr
--- Disk Addr
☐ In Mem
▨ On Disk
■ Unmapped

### Simplified Example
  - 16 page virtual address space

### Flags
  - P: Is entry in physical memory?
  - M: Has this part of VA space been mapped?

– 13 –    15-213, F'03

---

## P6 TLB Translation

CPU

result ← L2 and DRAM

20 VPN  12 VPO   virtual address (VA)

16   4
TLBT TLBI

L1 hit

L1 miss

TLB

TLB miss

TLB hit

L1 (128 sets, 4 lines/set)

10   10
VPN1 VPN2

TLB (16 sets, 4 entries/set)

20 PPN  12 PPO   →   20 CT  7 CI  5 CO

PDE   PTE

physical address (PA)

PDBR   Page tables

– 14 –    15-213, F'03

---

## P6 TLB

### TLB entry (not all documented, so this is speculative):

| 32 | 16 | 1 | 1 |
|---|---|---|---|
| PDE/PTE | Tag | PD | V |

  - V: indicates a valid (1) or invalid (0) TLB entry
  - PD: is this entry a PDE (1) or a PTE (0)?
  - tag: disambiguates entries cached in the same set
  - PDE/PTE: page directory or page table entry

- ● Structure of the data TLB:
  - 16 sets, 4 entries/set

| entry | entry | entry | entry | set 0 |
|---|---|---|---|---|
| entry | entry | entry | entry | set 1 |
| entry | entry | entry | entry | set 2 |
| ... | | | | |
| entry | entry | entry | entry | set 15 |

– 15 –    15-213, F'03

---

## Translating with the P6 TLB

CPU

20 VPN  12 VPO   virtual address

16   4
TLBT TLBI  ①         ②

TLB miss          PDE    PTE   TLB hit  ③
...
20 PPN  12 PPO

page table translation

physical address  ④

1. Partition VPN into TLBT and TLBI.
2. Is the PTE for VPN cached in set TLBI?
   - 3. Yes: then build physical address.
4. No: then read PTE (and PDE if not cached) from memory and build physical address.

– 16 –    15-213, F'03

Page 4

## P6 Page Table Translation

## Translating with the P6 Page Tables (case 1/1)



**Case 1/1: page table and page present.**

**MMU Action:**
- MMU builds physical address and fetches data word.

● **OS action**
- none

## Translating with the P6 Page Tables (case 1/0)



**Case 1/0: page table present but page missing.**

**MMU Action:**
- page fault exception
- handler receives the following args:
  ● VA that caused fault
  ● fault caused by non-present page or page-level protection violation
  ● read/write
  ● user/supervisor

## Translating with the P6 Page Tables (case 1/0, cont)



**OS Action:**
- Check for a legal virtual address.
- Read PTE through PDE.
- Find free physical page (swapping out current page if necessary)
- Read virtual page from disk and copy to virtual page
- Restart faulting instruction by returning from exception handler.

## Translating with the P6 Page Tables (case 0/1)

```
        20    12
       VPN  VPO

      VPN1 VPN2

Mem        PDE |p=0          data

    PDBR
        Page              Data
        directory         page
- - - - - - - - - - - - - - - - - - - - - -
                  PTE |p=1
Disk
                  Page
                  table
```

**Case 0/1: page table missing but page present.**

**Introduces consistency issue.**
- potentially every page out requires update of disk page table.

**Linux disallows this**
- if a page table is swapped out, then swap out its data pages too.

---

## Translating with the P6 Page Tables (case 0/0)

```
        20    12
       VPN  VPO

      VPN1 VPN2

Mem        PDE |p=0

    PDBR
        Page
        directory
- - - - - - - - - - - - - - - - - - - - - -
            PTE |p=0        data
Disk
            Page            Data
            table           page
```

**Case 0/0: page table and page missing.**

**MMU Action:**
- page fault exception

---

## Translating with the P6 Page Tables (case 0/0, cont)

```
        20    12
       VPN  VPO

      VPN1 VPN2

            PDE |p=1   PTE |p=0
Mem

    PDBR
        Page       Page
        directory  table
- - - - - - - - - - - - - - - - - - - - - -
                        data
Disk
                        Data
                        page
```

**OS action:**
- swap in page table.
- restart faulting instruction by returning from handler.

**Like case 0/1 from here on.**

---

## P6 L1 Cache Access

```
   CPU                         result      L2 and DRAM

   20    12                      L1          L1
  VPN  VPO  virtual address (VA) hit         miss

  16    4
 TLBT TLBI                               L1 (128 sets, 4 lines/set)

TLB                      TLB
miss                     hit

  10    10
 VPN1 VPN2    TLB (16 sets,
             4 entries/set)      PPN   PPO        CT      CI CO

   PDE    PTE                         physical
                                      address (PA)
  PDBR   Page tables
```

## L1 Cache Access



**32**
data ← L2 andDRAM

*L1 hit* *L1 miss*

L1 (128 sets, 4 lines/set)

• • •

| 20 | 7 | 5 |
|----|----|----|
| CT | CI | CO |

physical
address (PA)

Partition physical address into CO, CI, and CT.

Use CT to determine if line containing word at address PA is cached in set CI.

If no: check L2.
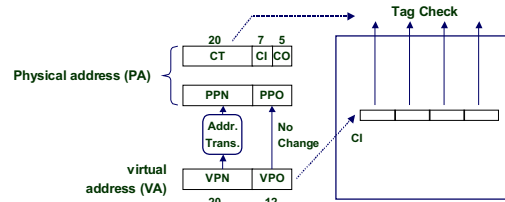
If yes: extract word at byte offset CO and return to processor.

– 25 –                                           15-213, F'03

---

## Speeding Up L1 Access

Tag Check

|  | 20 | 7 | 5 |
|--|----|----|----|
|  | CT | CI | CO |

**Physical address (PA)**

| PPN | PPO |
|-----|-----|

Addr. Trans.      No Change

**virtual address (VA)**

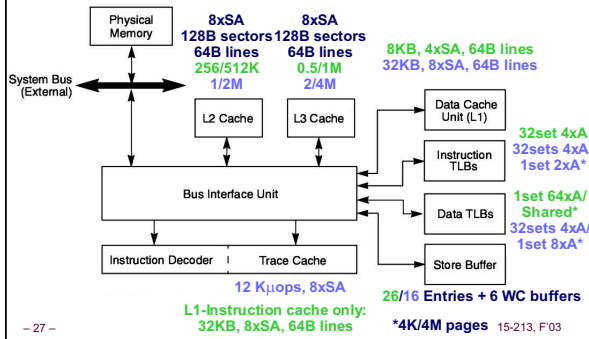| VPN | VPO |
|-----|-----|
| 20 | 12 |

CI

### Observation
- Bits that determine CI identical in virtual and physical address
- Can index into cache while address translation taking place
- Then check with CT from physical address
- "Virtually indexed, physically tagged"
- Cache carefully sized to make this possible

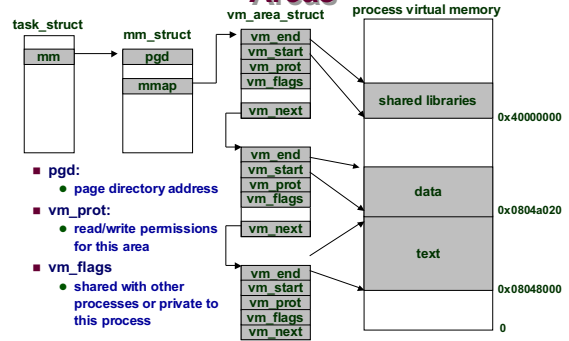– 26 –                                           15-213, F'03

---

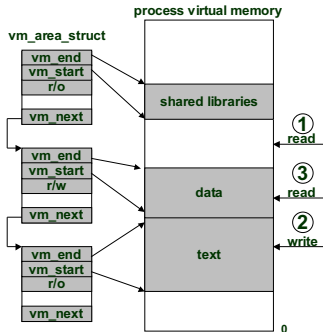## Pentium 4 Xeon Changes

**Pentium 4 Xeon** / **Pentium 4 Xeon MP**

Physical Memory

System Bus (External)

L2 Cache      L3 Cache

Bus Interface Unit

Instruction Decoder      Trace Cache

Data Cache Unit (L1)

Instruction TLBs

Data TLBs

Store Buffer

**8xSA
128B sectors
64B lines
256/512K
1/2M**

**8xSA
128B sectors
64B lines
0.5/1M
2/4M**

**8KB, 4xSA, 64B lines
32KB, 8xSA, 64B lines**

**32set 4xA
32sets 4xA/
1set 64xA\***

**1set 64xA/
Shared\*
32sets 4xA/
1set 8xA\***

**12 Kμops, 8xSA**

**26/16 Entries + 6 WC buffers**

**L1-Instruction cache only:
32KB, 8xSA, 64B lines**

**\*4K/4M pages**

– 27 –                                           15-213, F'03

---

## Linux Organizes VM as Collection of "Areas"

task_struct → mm_struct → vm_area_struct → process virtual memory

| mm |

| pgd |
| mmap |

vm_end
vm_start
vm_prot
vm_flags
vm_next

vm_end
vm_start
vm_prot
vm_flags
vm_next

vm_end
vm_start
vm_prot
vm_flags
vm_next

shared libraries      0x40000000

data      0x0804a020

text      0x08048000

0

- **pgd:**
  - page directory address
- **vm_prot:**
  - read/write permissions for this area
- **vm_flags**
  - shared with other processes or private to this process

– 28 –                                           15-213, F'03

## Linux Page Fault Handling

**process virtual memory**

**vm_area_struct**

| vm_end |
| vm_start |
| r/o |

| vm_next |

shared libraries

| vm_end |
| vm_start |
| r/w |

| vm_next |

data

| vm_end |
| vm_start |
| r/o |

| vm_next |

text

① read

③ read

② write

0

**Is the VA legal?**
- i.e. is it in an area defined by a vm_area_struct?
- if not then signal segmentation violation (e.g. (1))

**Is the operation legal?**
- i.e., can the process read/write this area?
- if not then signal protection violation (e.g., (2))

**If OK, handle fault**
- e.g., (3)

– 29 –                    15-213, F'03

---

## Memory Mapping

**Creation of new VM *area* done via "memory mapping"**
- create new vm_area_struct and page tables for area
- area can be backed by (i.e., get its initial values from) :
  - regular file on disk (e.g., an executable object file)
    - » initial page bytes come from a section of a file
  - nothing (e.g., bss)
    - » initial page bytes are zeros
- dirty pages are swapped back and forth between a special swap file.

**Key point:** no virtual pages are copied into physical memory until they are referenced!
- known as "demand paging"
- crucial for time and space efficiency

– 30 –                    15-213, F'03

---

## User-Level Memory Mapping

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

- map `len` bytes starting at offset `offset` of the file specified by file description `fd`, preferably at address `start` (usually 0 for don't care).
  - `prot`: MAP_READ, MAP_WRITE
  - `flags`: MAP_PRIVATE, MAP_SHARED
- return a pointer to the mapped area.
- Example: fast file copy
  - useful for applications like Web servers that need to quickly copy files.
  - `mmap` allows file transfers without copying into user space.

– 31 –                    15-213, F'03

---

## `mmap()` Example: Fast File Copy

```
#include <unistd.h>
#include <sys/mman.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <fcntl.h>

/*
 * mmap.c - a program that uses mmap
 * to copy itself to stdout
 */
```

```
int main() {
    struct stat stat;
    int i, fd, size;
    char *bufp;

    /* open the file & get its size*/
    fd = open("./mmap.c", O_RDONLY);
    fstat(fd, &stat);
    size = stat.st_size;
    /* map the file to a new VM area */
    bufp = mmap(0, size, PROT_READ,
      MAP_PRIVATE, fd, 0);

    /* write the VM area to stdout */
    write(1, bufp, size);
}
```
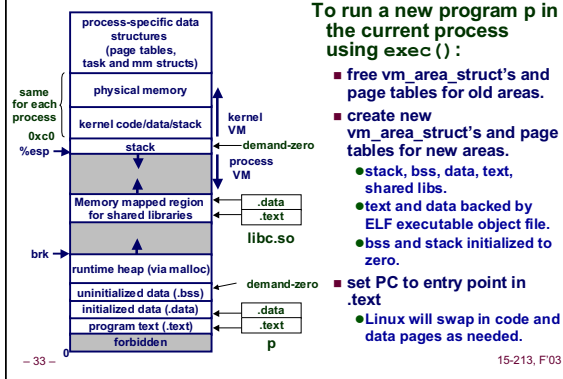
– 32 –                    15-213, F'03

## Exec() Revisited

| | |
|---|---|
| process-specific data structures (page tables, task and mm structs) | |
| physical memory | |
| kernel code/data/stack | kernel VM |
| stack | |
| Memory mapped region for shared libraries | .data / .text |
| runtime heap (via malloc) | libc.so |
| uninitialized data (.bss) | |
| initialized data (.data) | .data / .text |
| program text (.text) | |
| forbidden | p |

same for each process
0xc0
%esp
brk
demand-zero
process VM
demand-zero

**To run a new program p in the current process using exec():**

- free vm_area_struct's and page tables for old areas.
- create new vm_area_struct's and page tables for new areas.
  - stack, bss, data, text, shared libs.
  - text and data backed by ELF executable object file.
  - bss and stack initialized to zero.
- set PC to entry point in .text
  - Linux will swap in code and data pages as needed.

– 33 –    0    15-213, F'03

---

## Fork() Revisited

**To create a new process using `fork()`:**

- make copies of the old process's mm_struct, vm_area_struct's, and page tables.
  - at this point the two processes are sharing all of their pages.
  - How to get separate spaces without copying all the virtual pages from one space to another?
    » "copy on write" technique.
- copy-on-write
  - make pages of writeable areas read-only
  - flag vm_area_struct's for these areas as private "copy-on-write".
  - writes by either process to these pages will cause page faults.
    » fault handler recognizes copy-on-write, makes a copy of the page, and restores write permissions.
- Net result:
  - copies are deferred until absolutely necessary (i.e., when one of the processes tries to modify a shared page).

– 34 –    15-213, F'03

---

## Memory System Summary

**Cache Memory**

- Purely a speed-up technique
- Behavior invisible to application programmer and OS
- Implemented totally in hardware

**Virtual Memory**

- Supports many OS-related functions
  - Process creation
    » Initial
    » Forking children
  - Task switching
  - Protection
- Combination of hardware & software implementation
  - Software management of tables, allocations
  - Hardware access of tables
  - Hardware caching of table entries (TLB)

– 35 –    15-213, F'03