

"The course that gives CMU its Zip!"

Memory System Case Studies Mar. 20, 2008

Topics

- P6 address translation
- x86-64 extensions
- Linux memory management
- Linux page fault handling
- Memory mapping

class17.ppt

Intel P6

(Bob Colwell's Chip, CMU Alumni)

Internal designation for successor to Pentium

- Which had internal designation P5

Fundamentally different from Pentium

- Out-of-order, superscalar operation

Resulting processors

- Pentium Pro (1996)
- Pentium II (1997)
 - L2 cache on same chip
- Pentium III (1999)
 - The freshwater fish machines

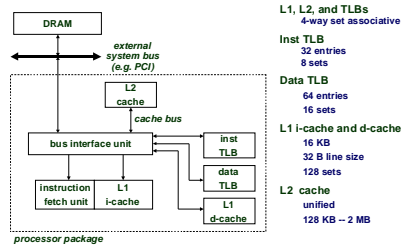
Saltwater fish machines: Pentium 4

- Different operation, but similar memory system
- Abandoned by Intel in 2005 for P6-based Core 2 Duo

2

15-213, S'08

P6 Memory System



- 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

3

15-213, S'08

Review of Abbreviations

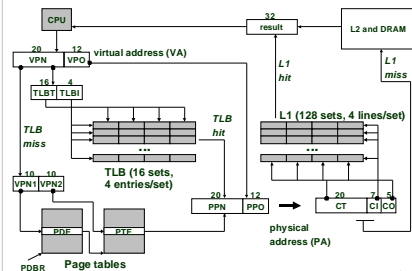
Symbols:

- **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

4

15-213, S'08

Overview of P6 Address Translation



5

15-213, S'08

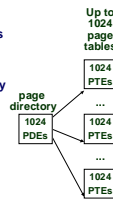
P6 2-level Page Table Structure

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.



6

15-213, S'08

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	US	RW	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
US: user or supervisor mode access
RW: 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

7

15-213, S'08

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	US	RW	P=1		

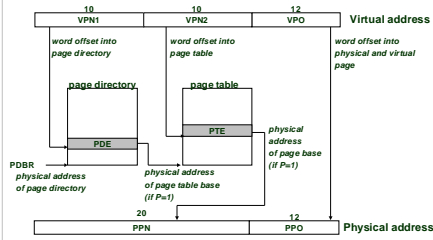
Page base address: 20 most significant bits of physical page address (forces pages to be 4 KB aligned)
Avail: available for system programmers
G: global page (don't evict from TLB on task switch)
D: dirty (set by MMU on writes)
A: accessed (set by MMU on reads and writes)
CD: cache disabled or enabled
WT: write-through or write-back cache policy for this page
US: user/supervisor
RW: 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

8

15-213, S'08

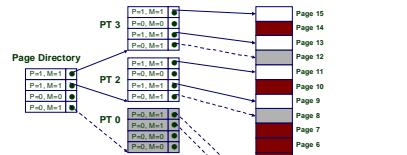
How P6 Page Tables Map Virtual Addresses to Physical Ones



9

15-213, S'08

Representation of VM Address Space



Simplified Example

- 16 page virtual address space

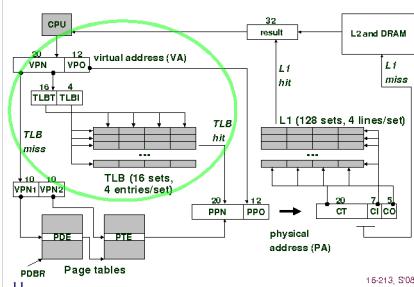
Flags

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

10

15-213, S'08

P6 TLB Translation



11

15-213, S'08

P6 TLB

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

20	16	1	1	1	1
PPN	Tag	V	G	S	W

- V: indicates a valid (1) or invalid (0) TLB entry
- Tag: disambiguates entries cached in the same set
- PPN: translation of the address indicated by index & tag
- G: page is "global" according to PDE, PTE
- S: page is "supervisor-only" according to PDE, PTE
- W: page is writable according to PDE, PTE
- D: PTE has already been marked "dirty" (once is enough)

Structure of the data TLB:

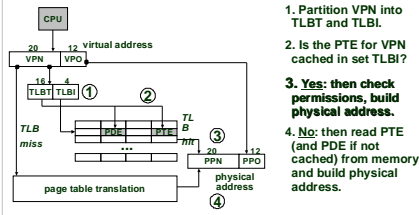
- 16 sets, 4 entries/set

entry	entry	entry	entry	set 0
entry	entry	entry	entry	set 1
...
entry	entry	entry	entry	set 15

12

15-213, S'08

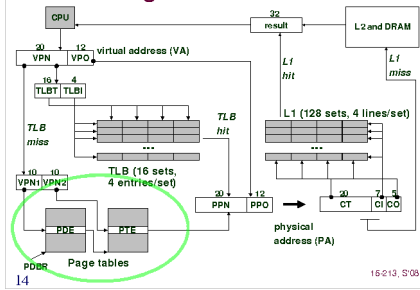
Translating with the P6 TLB



13

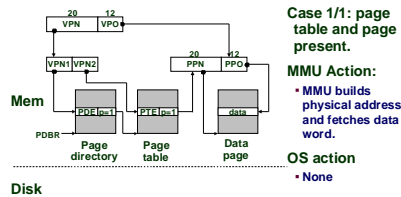
15-213, S'08

P6 Page Table Translation



16-213, S'08

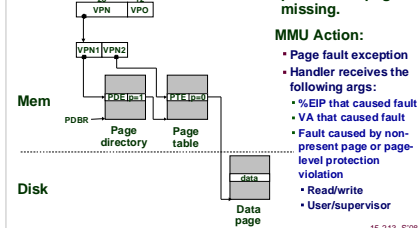
Translating with the P6 Page Tables (case 1/1)



15

15-213, S'08

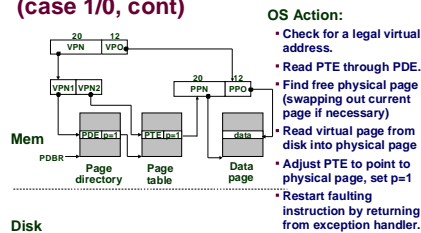
Translating with the P6 Page Tables (case 1/0)



16

15-213, S'08

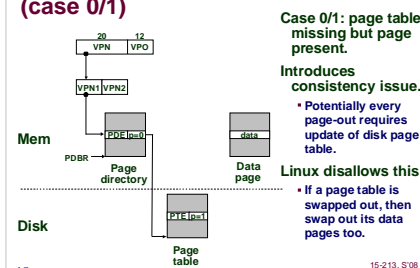
Translating with the P6 Page Tables (case 1/0, cont)



17

15-213, S'08

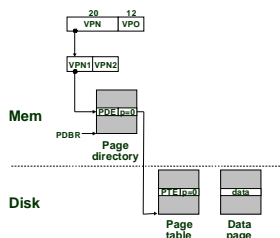
Translating with the P6 Page Tables (case 0/1)



18

15-213, S'08

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



Case 0/0: page table and page missing.

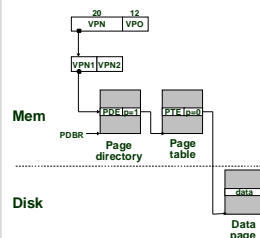
MMU Action:

- Page fault exception

19

15-213, S'08

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



OS action:

- Swap in page table.
- Restart faulting instruction by returning from handler.

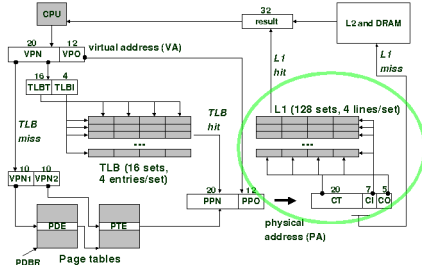
Like case 0/1 from here on.

- Two disk reads (8 sectors each) will probably require >15 ms

20

15-213, S'08

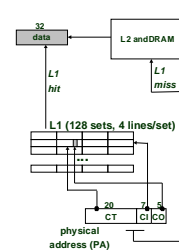
P6 L1 Cache Access



21

15-213, S'08

L1 Cache Access



Partition physical address into CO, CI, and CT.

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

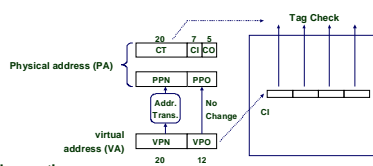
If no: check L2.

If yes: extract word at byte offset CO and return to processor.

22

15-213, S'08

Speeding Up L1 Access



Observation

- Bits that determine CI identical in virtual and physical address
- Can index into cache while address translation taking place
- Generally we hit in TLB, so PPN bits (CT bits) available next
- "Virtually indexed, physically tagged"
- Cache carefully sized to make this possible

23

15-213, S'08

x86-64 Paging

Origin

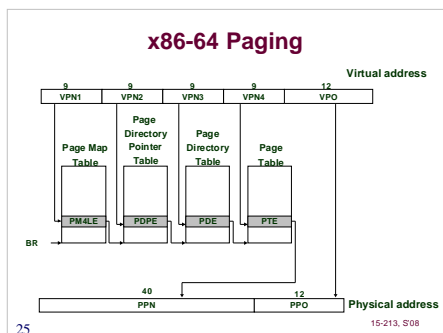
- AMD's way of extending x86 to 64-bit instruction set
- Intel has followed with "EM64T"

Requirements

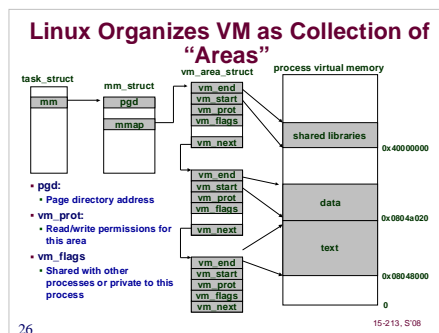
- 48-bit virtual address
- 256 terabytes (TB)
- Not yet ready for full 64 bits
 - Nobody can buy that much DRAM yet
 - Mapping tables would be huge
 - Multi-level array map may not be the right data structure
- 52-bit physical address
- Requires 64-bit table entries
- Keep traditional x86 4KB page size
 - (4096 bytes per PT) / (8 bytes per PTE) = only 512 entries per page

24

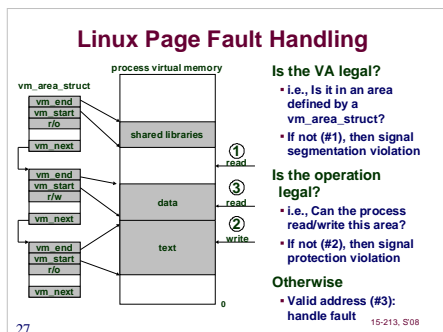
15-213, S'08



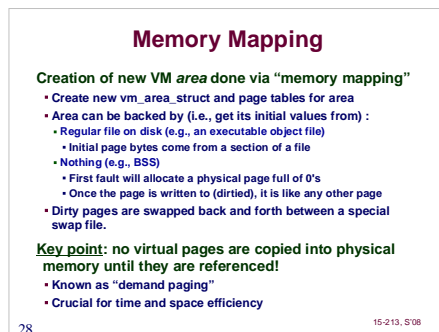
25



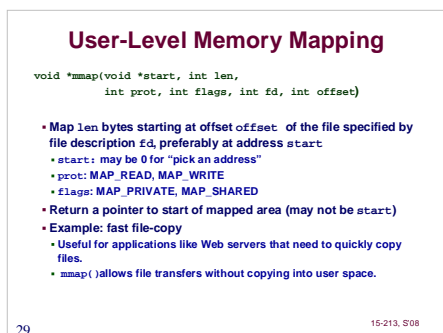
26



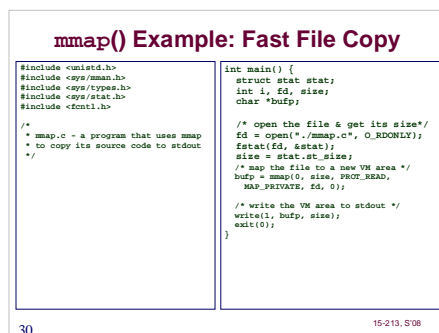
27



28

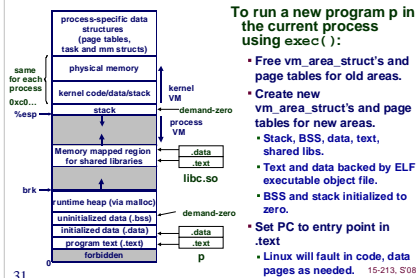


29



30

Exec() Revisited



31

15-213, S'08

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 share all of their pages.
- How to get separate spaces without copying all the virtual pages from one space to another?
 - "Copy on Write" (COW) technique.
- Copy-on-write
 - Mark PTE's of writeable areas as read-only
 - Writes by either process to these pages will cause page faults
 - Flag `vm_area_struct`'s for these areas as private "copy-on-write"
 - 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).

32

15-213, S'08

Memory System Summary

L1/L2 Memory Cache

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

Virtual Memory

- Supports many OS-related functions
 - Process creation, task switching, protection
- Software
 - Allocates/shares physical memory among processes
 - Maintains high-level tables tracking memory type, source, sharing
 - Handles exceptions, fills in hardware-defined mapping tables
- Hardware
 - Translates $V \rightarrow P$ via mapping tables, enforcing permissions
 - Accelerates mapping via translation cache (TLB)

33

15-213, S'08

Further Reading

Intel TLBs

Application Note: "TLBs, Paging-Structure Caches, and Their Invalidation", April 2007

34

15-213, S'08