

15-213

P6/Linux Memory System April 5, 2001

Topics

- P6 address translation
- Linux memory management
- Linux page fault handling
- memory mapping

class21.ppt

Intel P6

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

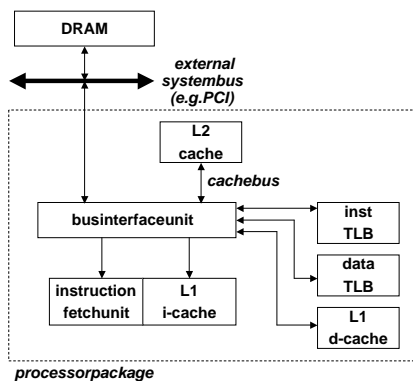
- Pentium Pro (1996)
- Pentium II (1997)
 - Incorporated MMX instructions
 - » special instructions for parallel processing
 - L2 cache on same chip
- Pentium III (1999)
 - Incorporated Streaming SIMD Extensions
 - » More instructions for parallel processing

class21.ppt

- 2 -

CS213S'01

P6 memory system



- 32bit address space
- 4KB 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
 - 16KB
 - 32B linesize
 - 128 sets
- L2 cache
 - unified
 - 128KB -- 2MB

class21.ppt

- 3 -

CS213S'01

Review of abbreviations

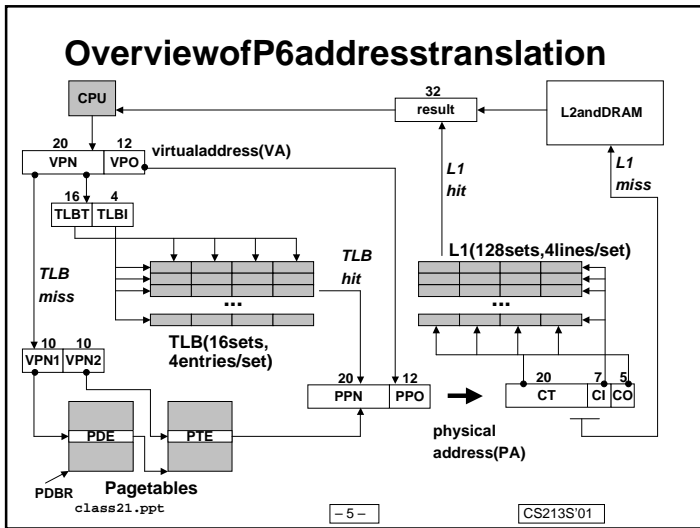
Symbols:

- Components of the virtual address (VA)
 - TLB: 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: cachetag
- Other
 - PDBR: Paged directory base register

class21.ppt

- 4 -

CS213S'01



P62 -level pagetable structure

Pagedirectory

- 1024 -byte pagedirectory entries (PDEs) that point to pagetables
- one pagedirectory per process.
- pagedirectory must be in memory when its process is running
- always pointed to by PDBR

Pagetables:

- 1024 -byte pagetable entries (PTEs) that point to pages.
- pagetables can be paged in and out.

Upto 1024 page tables

1024 PDEs

1024 PTEs

1024 PTEs

class21.ppt

- 6 -

CS213S'01

P6 pagedirectory entry (PDE)

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

Pagetable physical base address : 20 most significant bits of physical pagetable address (forces pagetable to be 4KB aligned)

Avail: 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: cached disabled (1) or enabled (0)

WT: write-through or write-back cache policy for this pagetable

U/S: user or supervisor mode access

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

P: pagetable is present in memory (1) or not (0)

31											1	0
Available for OS (pagetable location in secondary storage)												
											P=0	

class21.ppt

- 7 -

CS213S'01

P6 pagetable 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 page to be 4KB 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: cached disabled or enabled

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

U/S: user/supervisor

R/W: read/write

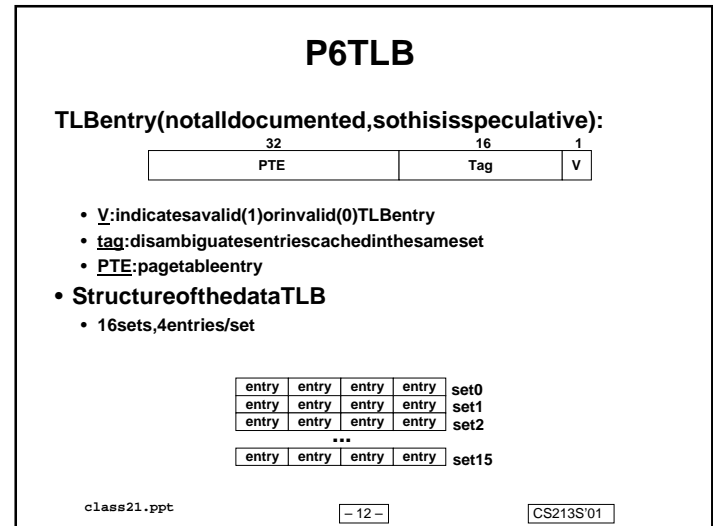
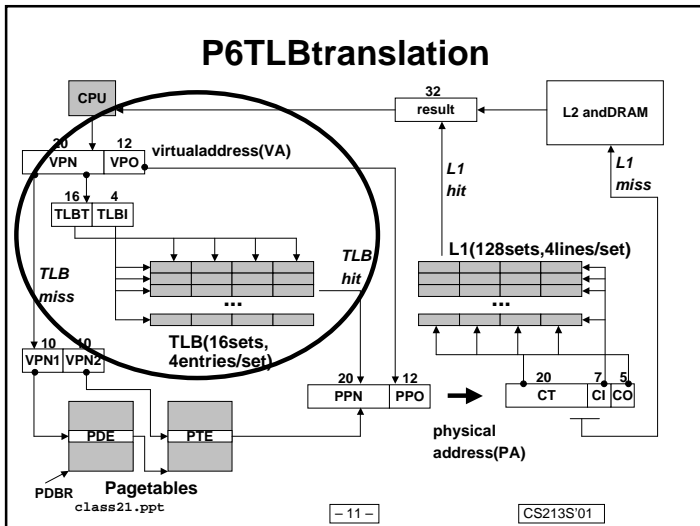
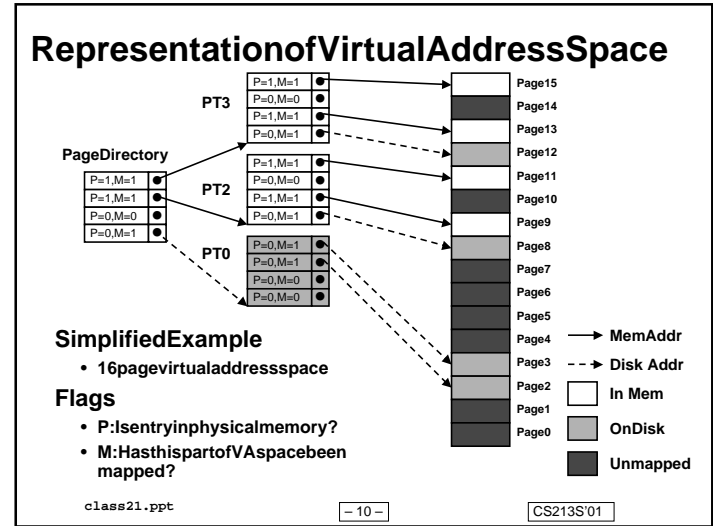
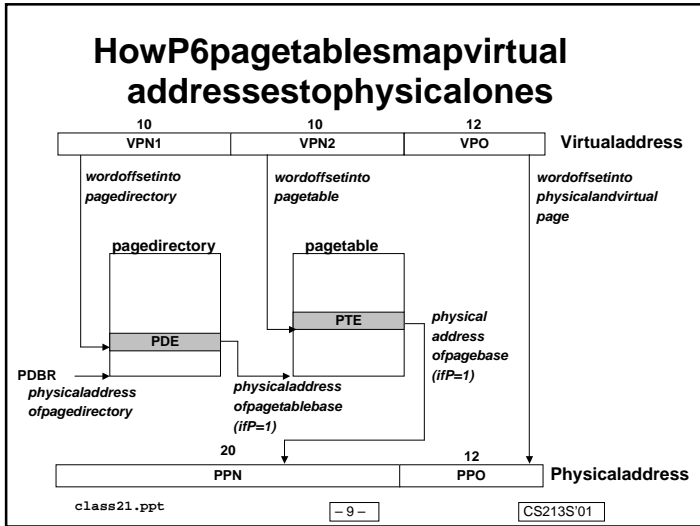
P: page is present in physical memory (1) or not (0)

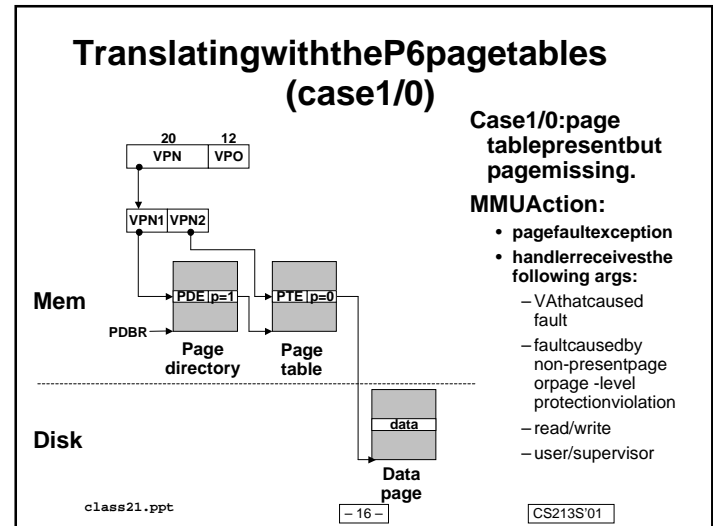
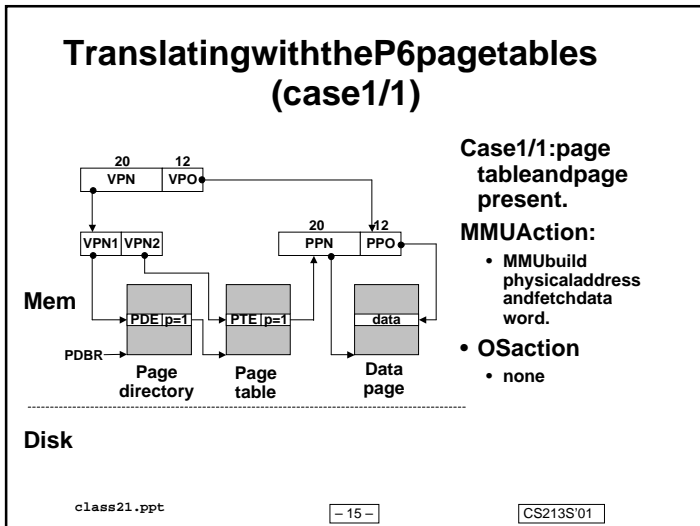
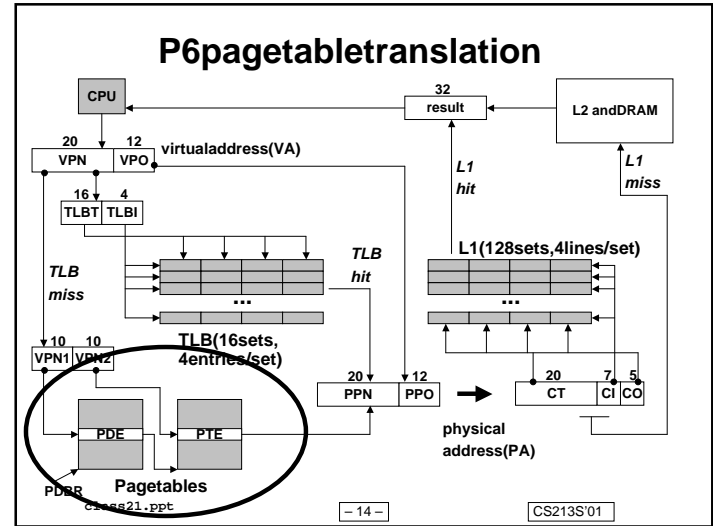
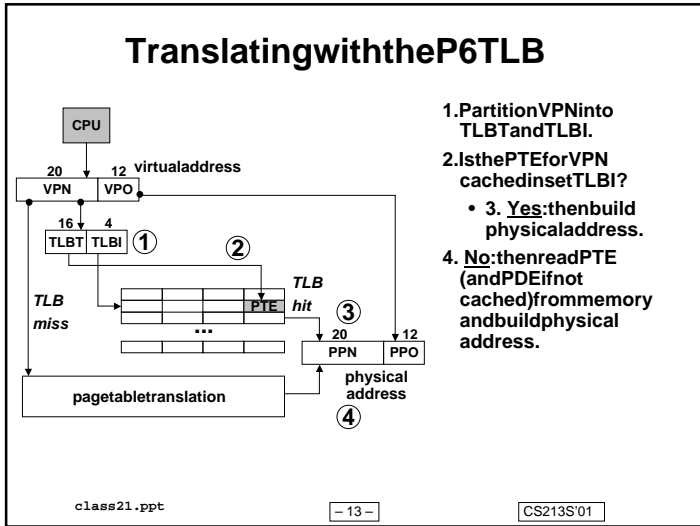
31											1	0
Available for OS (page location in secondary storage)												
											P=0	

class21.ppt

- 8 -

CS213S'01





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

OSAction:

- 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.

class21.ppt - 17 - CS213S'01

Translating with the P6 pagetables (case 0/1)

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

Introduces consistency issue.

- potentially every pageout requires update of disk page table.

Linux disallows this

- if a page table is swapped out, then swap out its data page too.

class21.ppt - 18 - CS213S'01

Translating with the P6 pagetables (case 0/0)

Case 0/0: page table and page missing.

MMUAction:

- page fault exception

class21.ppt - 19 - CS213S'01

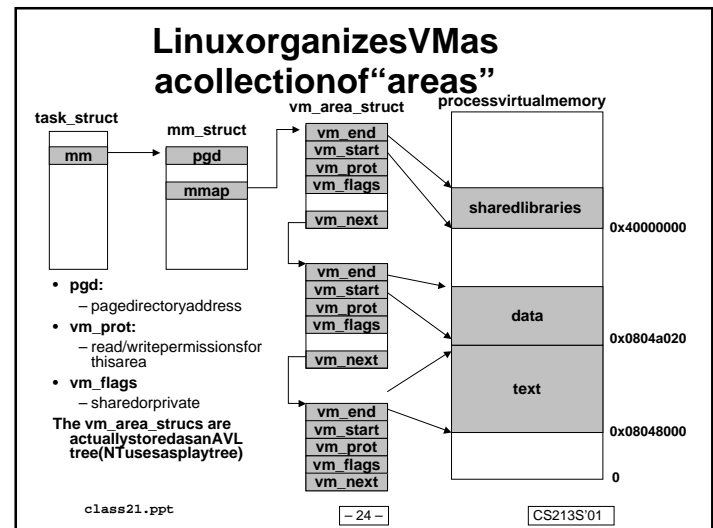
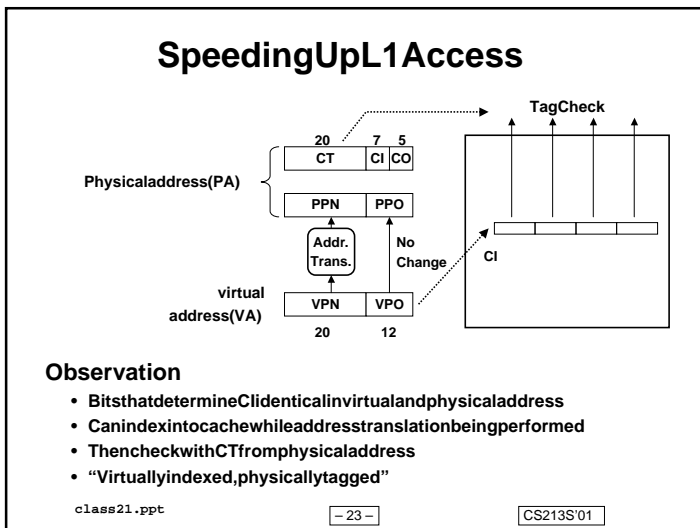
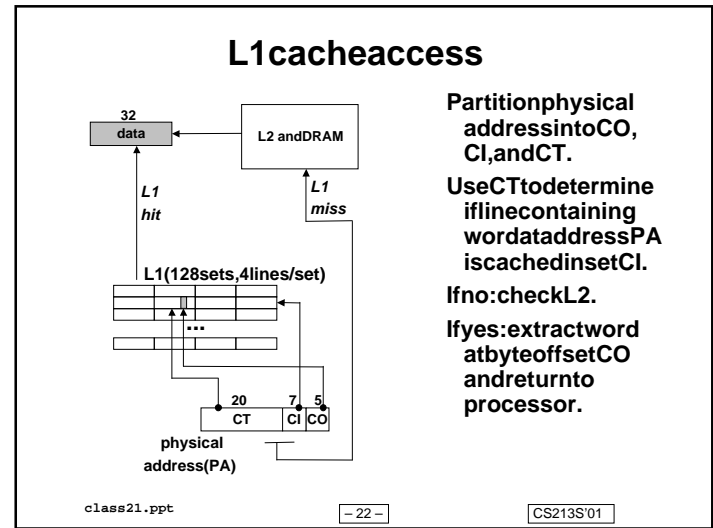
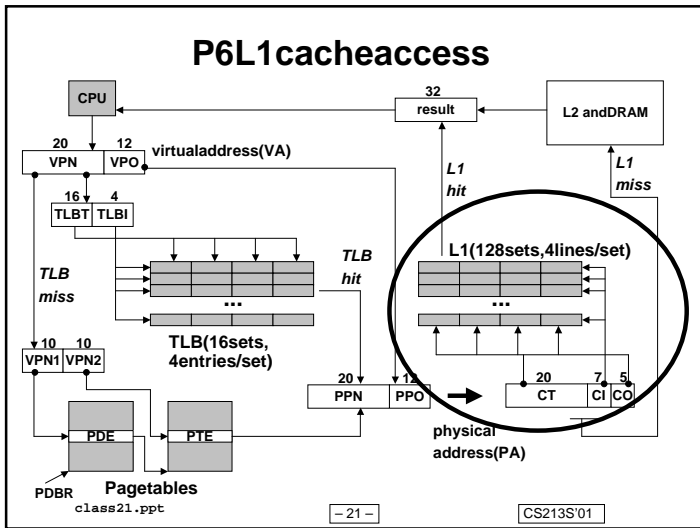
Translating with the P6 pagetables (case 0/0, cont)

OSAction:

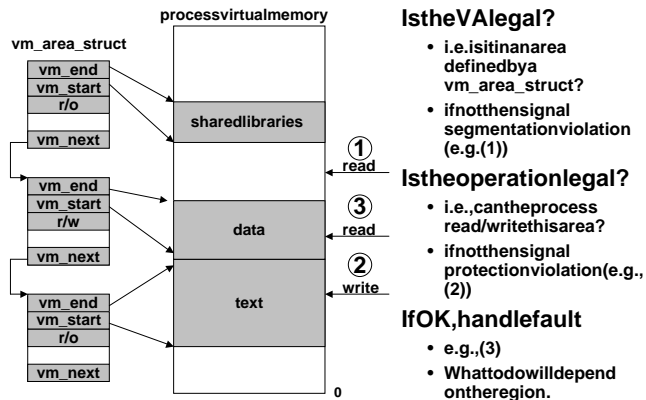
- swap in page table.
- restart faulting instruction by returning from handler.

Like case 0/1 from here on.

class21.ppt - 20 - CS213S'01



Linux page fault handling



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)
- What to do will depend on the region.

class21.ppt

- 25 -

CS213S'01

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. i.e.

Keypoint: no virtual pages are copied into physical memory until they are referenced!

- known as "demand paging"
- crucial for time and space efficiency

class21.ppt

- 26 -

CS213S'01

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 if 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 userspace.

class21.ppt

- 27 -

CS213S'01

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 and 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);

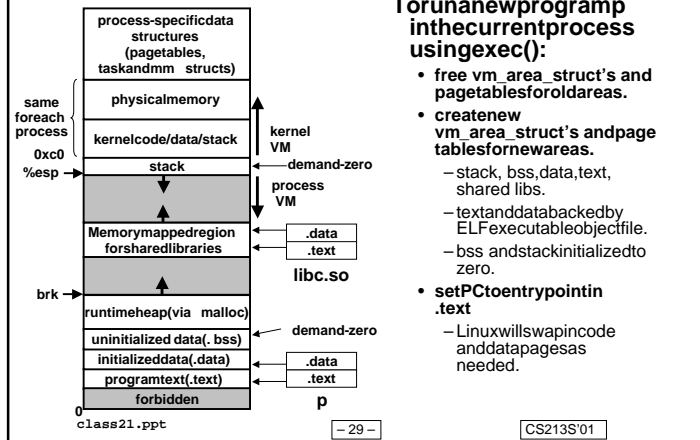
    /* write the VM area to stdout */
    write(1, bufp, size);
}
```

class21.ppt

- 28 -

CS213S'01

Exec() revisited



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

- free `vm_area_struct`'s and pagetables for old areas.
- create new `vm_area_struct`'s and pagetables 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.

Fork() revisited

To create a new process using fork:

- make copies of the old process's `smm_struct`, `vm_area_struct`'s, and pagetables.
 - 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).

class21.ppt

- 30 -

CS213S'01

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 to tables
 - Hardware caching of table entries (TLB)

class21.ppt

- 31 -

CS213S'01