

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

**P6 / Linux Memory System**  
**March 23, 2004**

**Topics**

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

class19.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

- 2 - 15-213, S'04

**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'04

**Review of Abbreviations**

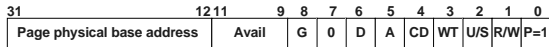
**Symbols:**

- 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'04



## P6 Page Table Entry (PTE)



**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)

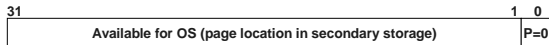
**CD:** cache disabled or enabled

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

**U/S:** user/supervisor

**R/W:** read/write

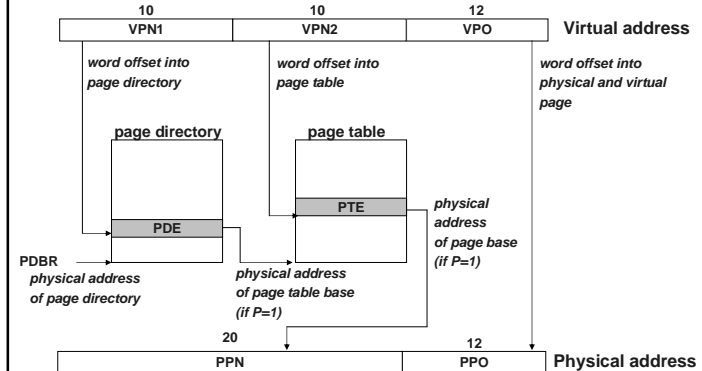
**P:** page is present in physical memory (1) or not (0)



- 9 -

15-213, S'04

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

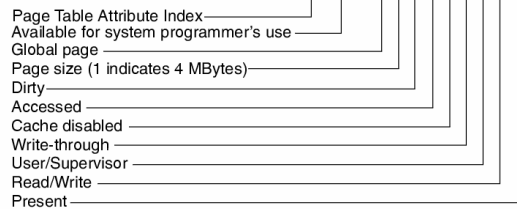


- 10 -

15-213, S'04

## 4Mbyte PDE's

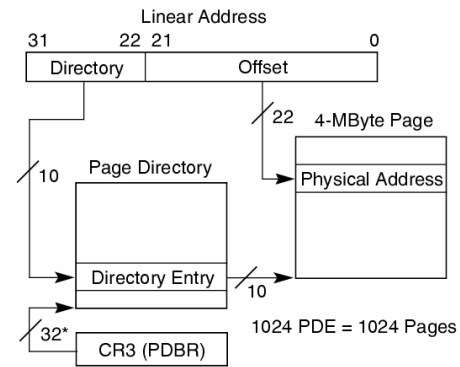
### Page-Directory Entry (4-MByte Page)



- 11 -

15-213, S'04

## Support for 4Mbyte Pages

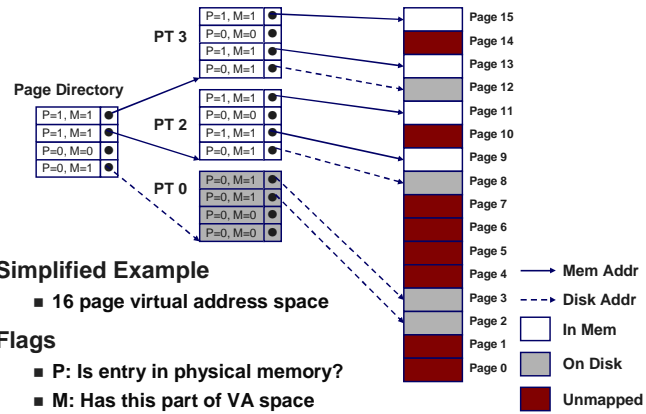


\*32 bits aligned onto a 4-KByte boundary.

- 12 -

15-213, S'04

## Representation of VM Address Space

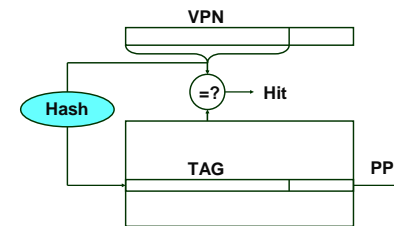


- 13 -

15-213, S'04

## Aside: Inverted Page Tables

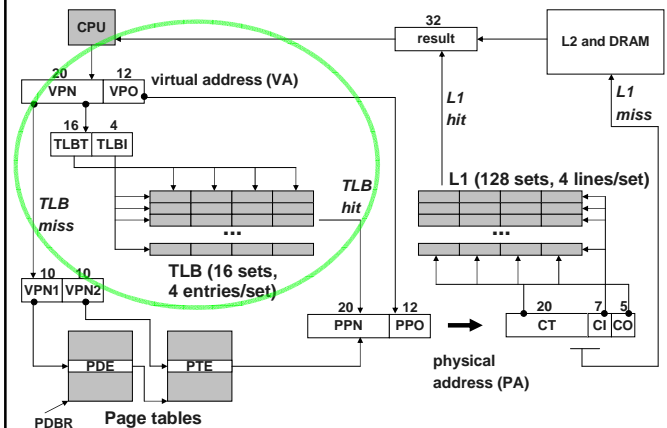
Problem: Page-table storage overhead is a function of the virtual address space size.



- 14 -

15-213, S'04

## P6 TLB Translation

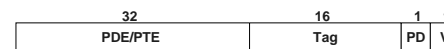


- 15 -

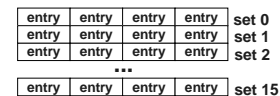
15-213, S'04

## P6 TLB

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



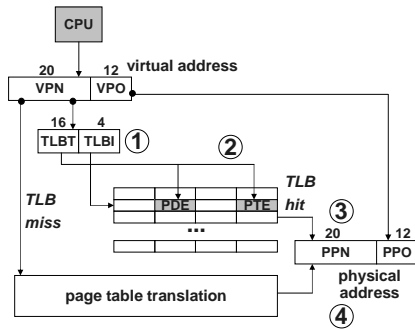
- **V**: indicates a valid (1) or invalid (0) TLB entry
- **PD**: is this entry a PDE (1) or a PTE (0)?
- **tag**: disambiguates entries cached in the same set
- **PDE/PTE**: page directory or page table entry
- **Structure of the data TLB:**
  - 16 sets, 4 entries/set



- 16 -

15-213, S'04

## Translating with the P6 TLB

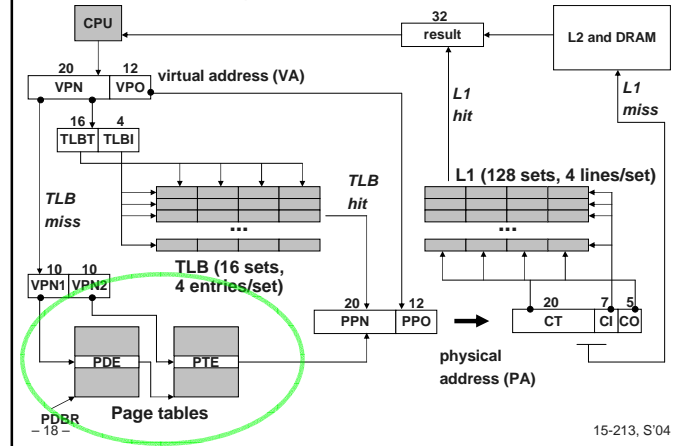


1. Partition VPN into TLBT and TLBI.
2. Is the PTE for VPN cached in set TLBI?
  - **3. Yes:** then build physical address.
4. **No:** then read PTE (and PDE if not cached) from memory and build physical address.

- 17 -

15-213, S'04

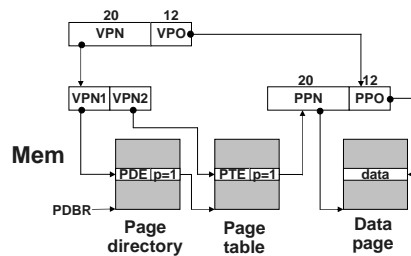
## P6 Page Table Translation



- 18 -

15-213, S'04

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



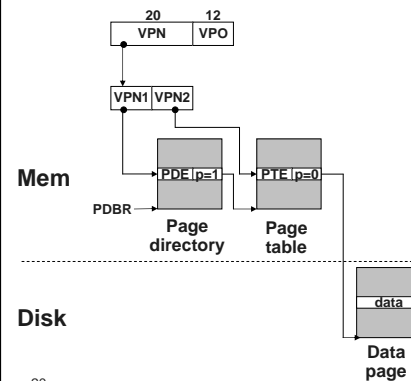
Case 1/1: page table and page present.

- MMU Action:
- MMU builds physical address and fetches data word.
- OS action
- none

- 19 -

15-213, S'04

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



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

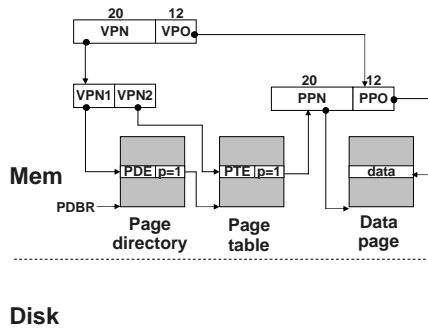
MMU Action:

- 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

- 20 -

15-213, S'04

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



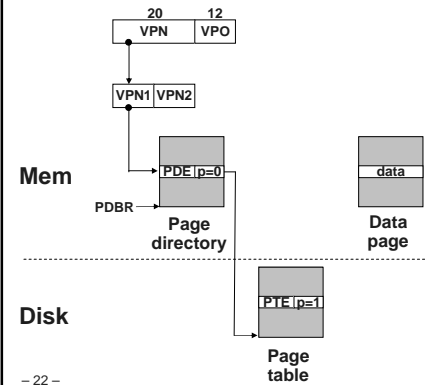
OS Action:

- 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.

- 21 -

15-213, S'04

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



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

Introduces consistency issue.

- potentially every page out requires update of disk page table.

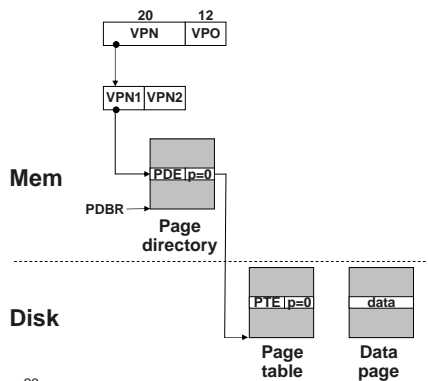
Linux disallows this

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

- 22 -

15-213, S'04

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



Case 0/0: page table and page missing.

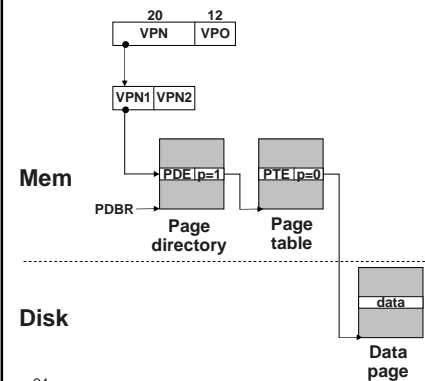
MMU Action:

- page fault exception

- 23 -

15-213, S'04

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



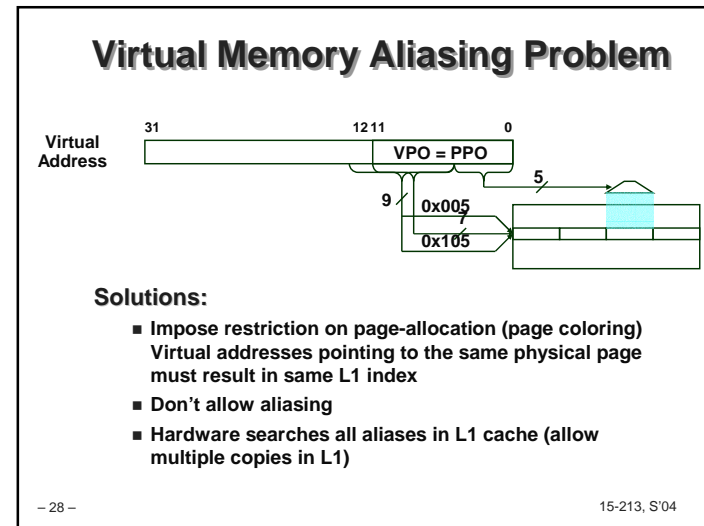
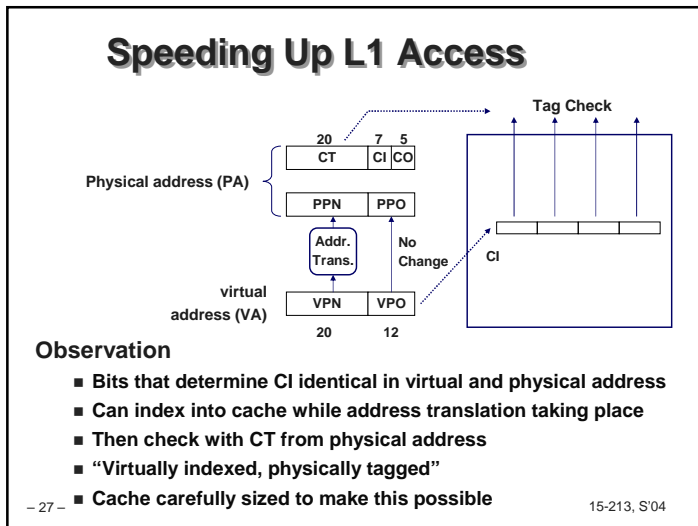
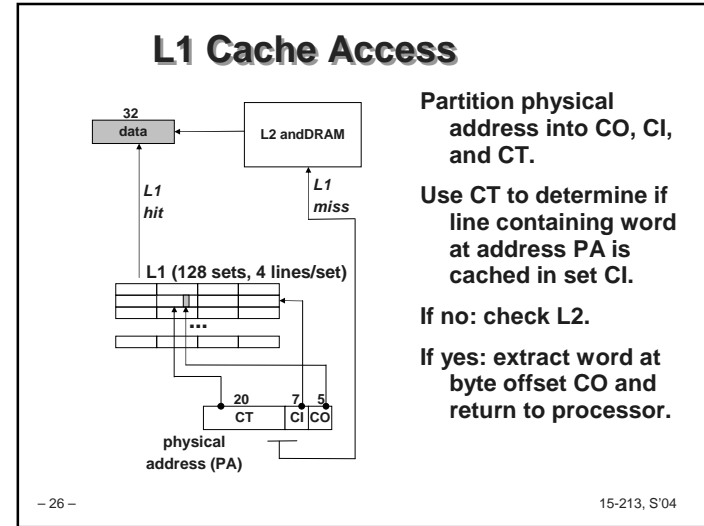
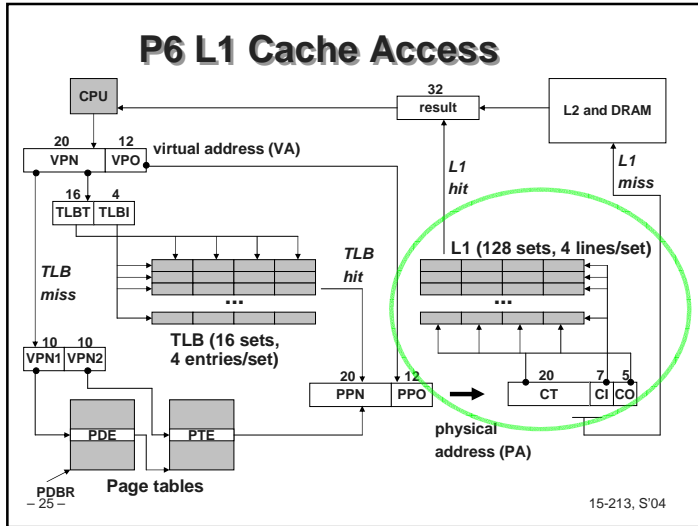
OS action:

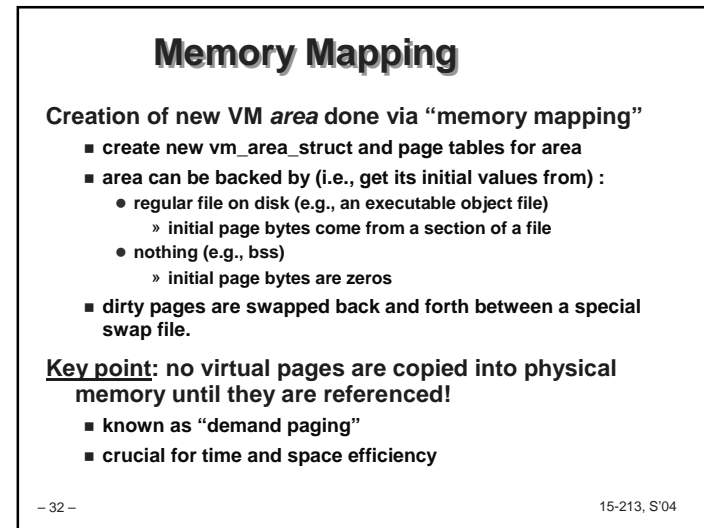
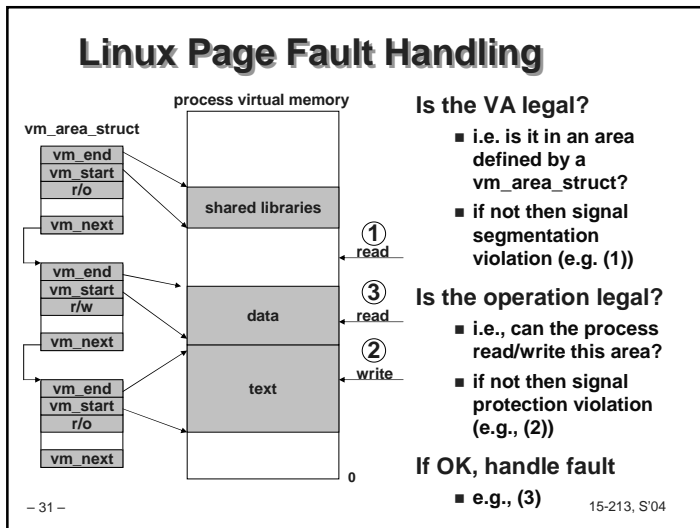
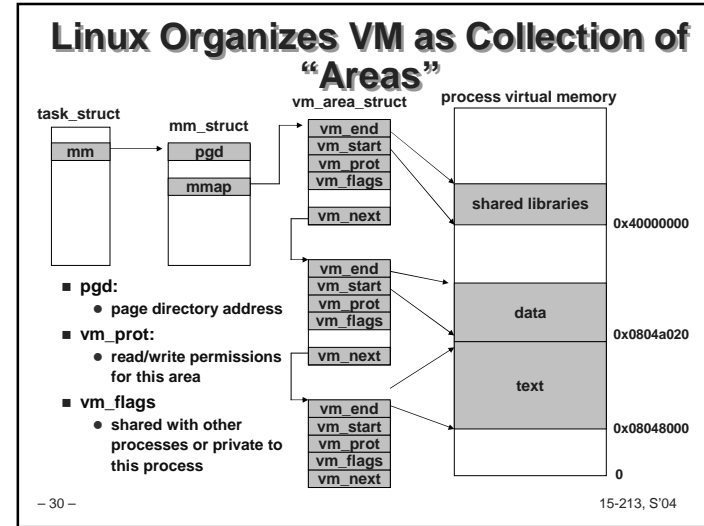
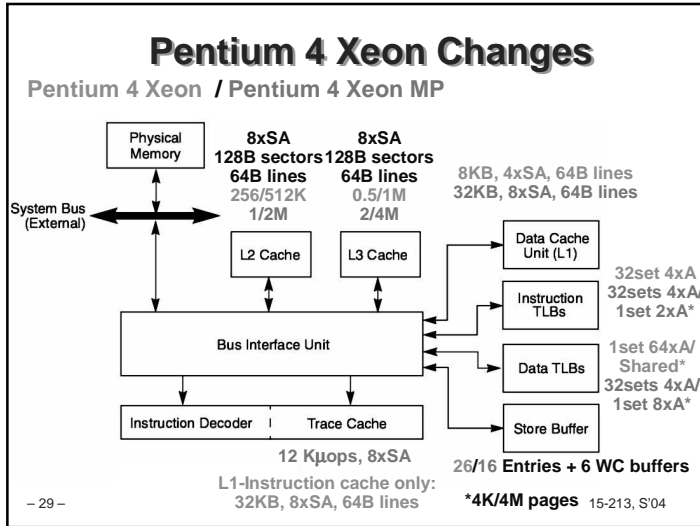
- swap in page table.
- restart faulting instruction by returning from handler.

Like case 0/1 from here on.

- 24 -

15-213, S'04







## 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.

- 33 -

15-213, S'04

## 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;

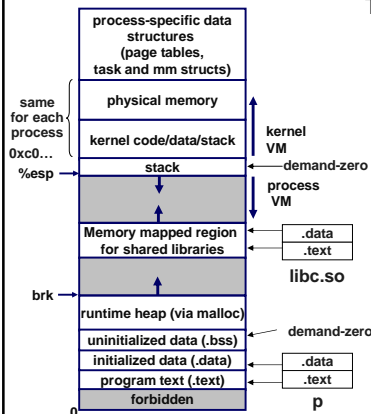
    /* open the file & get its size*/
    fd = open("./mmap.c", O_RDONLY);
    fstat(fd, &stat);
    size = stat.st_size;
    /* map the file to a new VM area */
    bufp = mmap(0, size, PROT_READ,
               MAP_PRIVATE, fd, 0);

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

- 34 -

15-213, S'04

## Exec() Revisited



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.

- 35 -

15-213, S'04

## 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).

- 36 -

15-213, S'04

## **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)