

CS162 Operating Systems and Systems Programming Lecture 12

Address Translation

March 10, 2008
Prof. Anthony D. Joseph
<http://inst.eecs.berkeley.edu/~cs162>

Review: Important Aspects of Memory Multiplexing

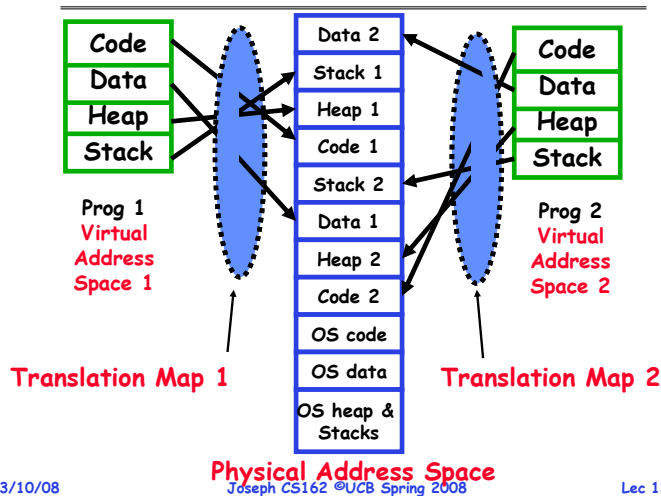
- **Controlled overlap:**
 - Separate state of threads should not collide in physical memory. Obviously, unexpected overlap causes chaos!
 - Conversely, would like the ability to overlap when desired (for communication)
- **Translation:**
 - Ability to translate accesses from one address space (virtual) to a different one (physical)
 - When translation exists, processor uses virtual addresses, physical memory uses physical addresses
 - Side effects:
 - » Can be used to avoid overlap
 - » Can be used to give uniform view of memory to programs
- **Protection:**
 - Prevent access to private memory of other processes
 - » Different pages of memory can be given special behavior (Read Only, Invisible to user programs, etc).
 - » Kernel data protected from User programs
 - » Programs protected from themselves

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.2

Review: General Address Translation



3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.3

Goals for Today

- **Address Translation Schemes**
 - Segmentation
 - Paging
 - Multi-level translation
 - Paged page tables
 - Inverted page tables
- **Comparison among options**

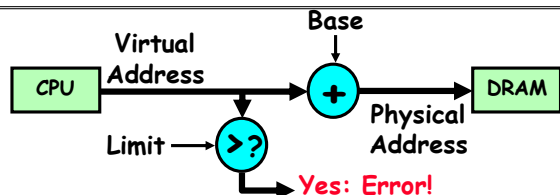
Note: Some slides and/or pictures in the following are adapted from slides ©2005 Silberschatz, Galvin, and Gagne. Many slides generated from my lecture notes by Kubiawicz.

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.4

Review: Simple Segmentation: Base and Bounds (CRAY-1)



- Can use base & bounds/limit for dynamic address translation (Simple form of "segmentation"):
 - Alter every address by adding "base"
 - Generate error if address bigger than limit
- This gives program the illusion that it is running on its own dedicated machine, with memory starting at 0
 - Program gets continuous region of memory
 - Addresses within program do not have to be relocated when program placed in different region of DRAM

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.5

Base and Limit segmentation discussion

- Provides level of indirection
 - OS can move bits around behind program's back
 - Can be used to correct if program needs to grow beyond its bounds or coalesce fragments of memory
- Only OS gets to change the base and limit!
 - Would defeat protection
- What gets saved/restored on a context switch?
 - Everything from before + base/limit values
 - Or: How about complete contents of memory (out to disk)?
 - » Called "Swapping"
- Hardware cost
 - 2 registers/Adder/Comparator
 - Slows down hardware because need to take time to do add/compare on every access
- Base and Limit Pros: Simple, relatively fast

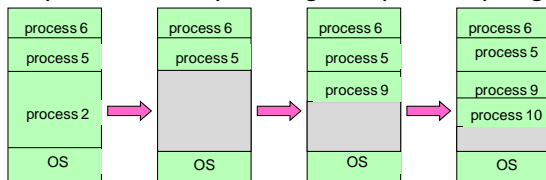
3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.6

Cons for Simple Segmentation Method

- Fragmentation problem (complex memory allocation)
 - Not every process is the same size
 - Over time, memory space becomes fragmented
 - Really bad if want space to grow dynamically (e.g. heap)



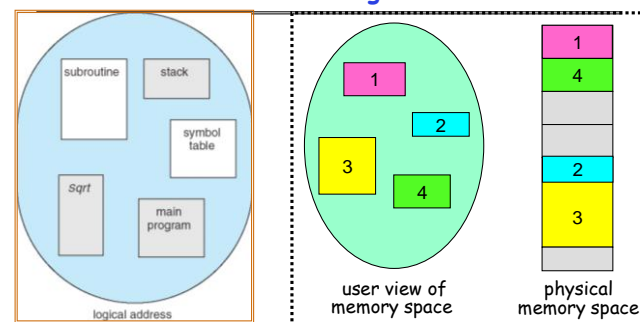
- Other problems for process maintenance
 - Doesn't allow heap and stack to grow independently
 - Want to put these as far apart in virtual memory space as possible so that they can grow as needed
- Hard to do inter-process sharing
 - Want to share code segments when possible
 - Want to share memory between processes

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.7

More Flexible Segmentation



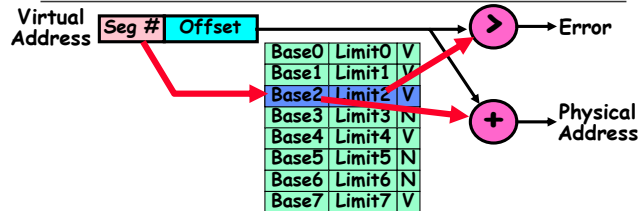
- Logical View: multiple separate segments
 - Typical: Code, Data, Stack
 - Others: memory sharing, etc
- Each segment is given region of contiguous memory
 - Has a base and limit
 - Can reside anywhere in physical memory

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.8

Implementation of Multi-Segment Model



- Segment map resides in processor
 - Segment number mapped into base/limit pair
 - Base added to offset to generate physical address
 - Error check catches offset out of range
- As many chunks of physical memory as entries
 - Segment addressed by portion of virtual address
 - However, could be included in instruction instead:
 - » x86 Example: `mov [es:bx], ax.`
- What is "V/N"?
 - Can mark segments as invalid; requires check as well

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.9

Administrivia

- Project 2
 - Initial Design Document due today (3/10) at 11:59pm
 - Look at the lecture schedule to keep up with due dates!
- Midterm #1 re-grade requests due by Friday at 5pm
 - Entire exam will be re-graded

3/10/08

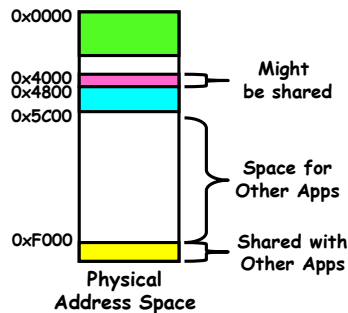
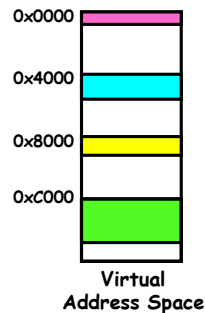
Joseph CS162 ©UCB Spring 2008

Lec 12.10

Example: Four Segments (16 bit addresses)



Seg ID #	Base	Limit
0 (code)	0x4000	0x0800
1 (data)	0x4800	0x1400
2 (shared)	0xF000	0x1000
3 (stack)	0x0000	0x3000



3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.11

Example of segment translation

0x240	main:	la \$a0, varx
0x244		jal strlen
...		...
0x360	strlen:	li \$v0, 0 ;count
0x364	loop:	lb \$t0, (\$a0)
0x368		beq \$r0,\$t1, done
...		...
0x4050	varx	dw 0x314159

Seg ID #	Base	Limit
0 (code)	0x4000	0x0800
1 (data)	0x4800	0x1400
2 (shared)	0xF000	0x1000
3 (stack)	0x0000	0x3000

Let's simulate a bit of this code to see what happens (PC=0x240):

- Fetch 0x240. Virtual segment #? 0; Offset? 0x240
Physical address? Base=0x4000, so physical addr=0x4240
Fetch instruction at 0x4240. Get "la \$a0, varx"
Move 0x4050 → \$a0, Move PC+4→PC
- Fetch 0x244. Translated to Physical=0x4244. Get "jal strlen"
Move 0x0248 → \$ra (return address!), Move 0x0360 → PC
- Fetch 0x360. Translated to Physical=0x4360. Get "li \$v0,0"
Move 0x0000 → \$v0, Move PC+4→PC
- Fetch 0x364. Translated to Physical=0x4364. Get "lb \$t0,(\$a0)"
Since \$a0 is 0x4050, try to load byte from 0x4050
Translate 0x4050. Virtual segment #? 1; Offset? 0x50
Physical address? Base=0x4800, Physical addr = 0x4850,
Load Byte from 0x4850→\$t0, Move PC+4→PC

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.12

Observations about Segmentation

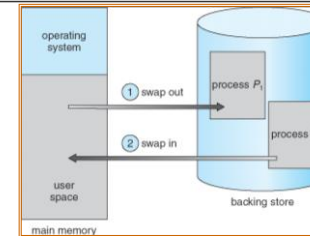
- **Virtual address space has holes**
 - Segmentation efficient for sparse address spaces
 - A correct program should never address gaps (except as mentioned in moment)
 - » If it does, trap to kernel and dump core
- **When it is OK to address outside valid range:**
 - This is how the stack and heap are allowed to grow
 - For instance, stack takes fault, system automatically increases size of stack
- **Need protection mode in segment table**
 - For example, code segment would be read-only
 - Data and stack would be read-write (stores allowed)
 - Shared segment could be read-only or read-write
- **What must be saved/restored on context switch?**
 - Segment table stored in CPU, not in memory (small)
 - Might store all of processes memory onto disk when switched (called "swapping")

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.13

Schematic View of Swapping



- **Extreme form of Context Switch: Swapping**
 - In order to make room for next process, some or all of the previous process is moved to disk
 - » Likely need to send out complete segments
 - This greatly increases the cost of context-switching
- **Desirable alternative?**
 - Some way to keep only active portions of a process in memory at any one time
 - Need finer granularity control over physical memory

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.14

Paging: Physical Memory in Fixed Size Chunks

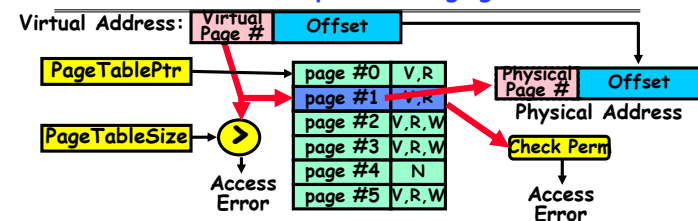
- **Problems with segmentation?**
 - Must fit variable-sized chunks into physical memory
 - May move processes multiple times to fit everything
 - Limited options for swapping to disk
- **Fragmentation: wasted space**
 - **External:** free gaps between allocated chunks
 - **Internal:** don't need all memory within allocated chunks
- **Solution to fragmentation from segments?**
 - Allocate physical memory in fixed size chunks ("pages")
 - Every chunk of physical memory is equivalent
 - » Can use simple vector of bits to handle allocation: 00110001110001101 ... 110010
 - » Each bit represents page of physical memory
1⇒allocated, 0⇒free
- **Should pages be as big as our previous segments?**
 - No: Can lead to lots of internal fragmentation
 - » Typically have small pages (1K-16K)
 - Consequently: need multiple pages/segment

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.15

How to Implement Paging?



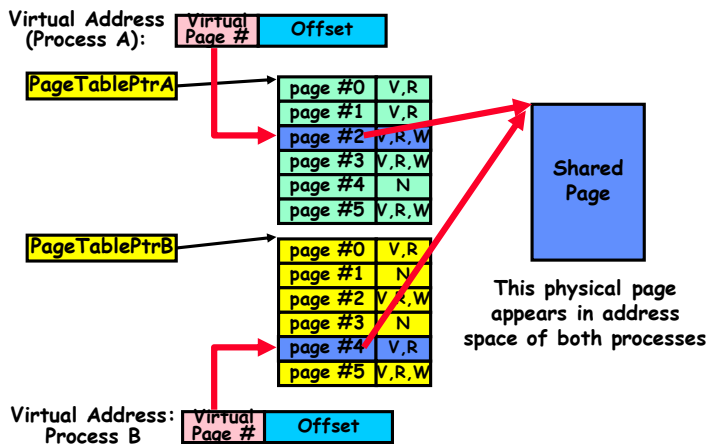
- **Page Table (One per process)**
 - Resides in physical memory
 - Contains physical page and permission for each virtual page
 - » Permissions include: Valid bits, Read, Write, etc
- **Virtual address mapping**
 - Offset from Virtual address copied to Physical Address
 - » Example: 10 bit offset ⇒ 1024-byte pages
 - Virtual page # is all remaining bits
 - » Example for 32-bits: 32-10 = 22 bits, i.e. 4 million entries
 - » Physical page # copied from table into physical address
 - Check Page Table bounds and permissions

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.16

What about Sharing?

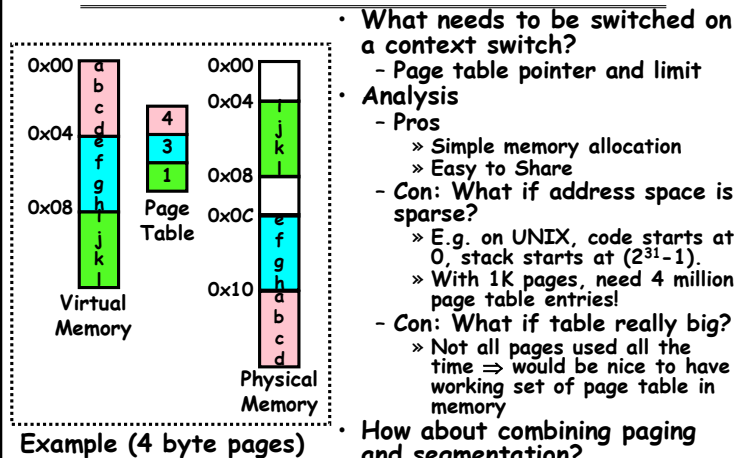


3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.17

Simple Page Table Discussion



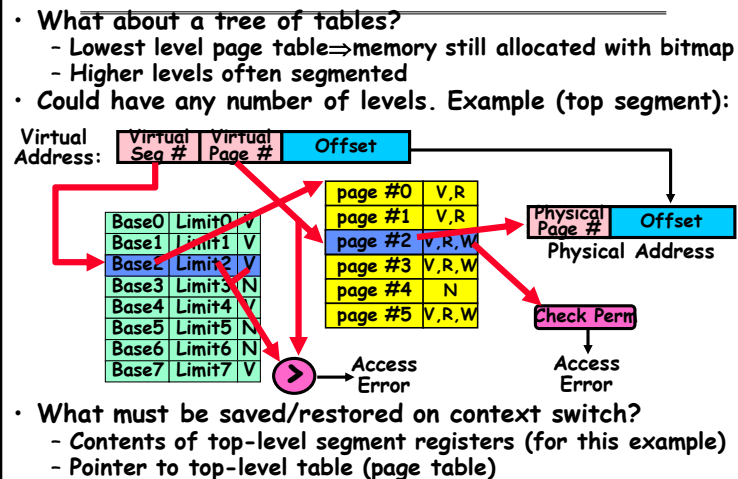
3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.18

BREAK

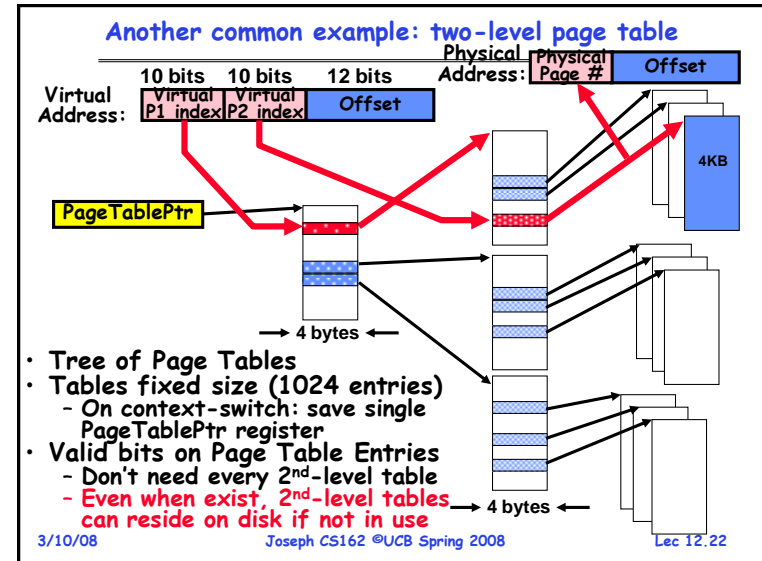
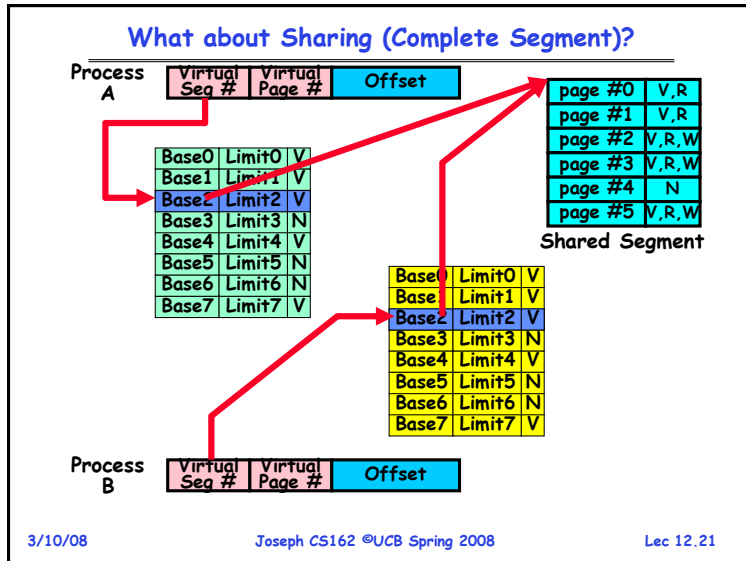
Multi-level Translation



3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.20

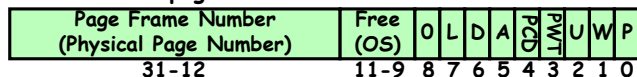


- ### Multi-level Translation Analysis
- Pros:
 - Only need to allocate as many page table entries as we need for application
 - » In other words, sparse address spaces are easy
 - Easy memory allocation
 - Easy Sharing
 - » Share at segment or page level (need additional reference counting)
 - Cons:
 - One pointer per page (typically 4K - 16K pages today)
 - Page tables need to be contiguous
 - » However, previous example keeps tables to exactly one page in size
 - Two (or more, if >2 levels) lookups per reference
 - » Seems very expensive!
- 3/10/08 Joseph CS162 ©UCB Spring 2008 Lec 12.23

- ### Inverted Page Table
- With all previous examples ("Forward Page Tables")
 - Size of page table is at least as large as amount of virtual memory allocated to processes
 - Physical memory may be much less
 - » Much of process space may be out on disk or not in use
-
- Virtual Page # | Offset
- Hash Table
- Physical Page # | Offset
- Answer: use a hash table
 - Called an "Inverted Page Table"
 - Size is independent of virtual address space
 - Directly related to amount of physical memory
 - Very attractive option for 64-bit address spaces
 - Cons: Complexity of managing hash changes
 - Often in hardware!
- 3/10/08 Joseph CS162 ©UCB Spring 2008 Lec 12.24

What is in a PTE?

- What is in a Page Table Entry (or PTE)?
 - Pointer to next-level page table or to actual page
 - Permission bits: valid, read-only, read-write, write-only
- Example: Intel x86 architecture PTE:
 - Address same format previous slide (10, 10, 12-bit offset)
 - Intermediate page tables called "Directories"



- P: Present (same as "valid" bit in other architectures)
- W: Writeable
- U: User accessible
- PWT: Page write transparent: external cache write-through
- PCD: Page cache disabled (page cannot be cached)
- A: Accessed: page has been accessed recently
- D: Dirty (PTE only): page has been modified recently
- L: L=1 ⇒ 4MB page (directory only).
Bottom 22 bits of virtual address serve as offset

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.25

Examples of how to use a PTE

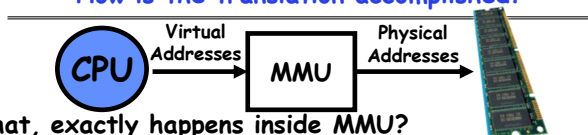
- How do we use the PTE?
 - Invalid PTE can imply different things:
 - » Region of address space is actually invalid or
 - » Page/directory is just somewhere else than memory
 - Validity checked first
 - » OS can use other (say) 31 bits for location info
- Usage Example: Demand Paging
 - Keep only active pages in memory
 - Place others on disk and mark their PTEs invalid
- Usage Example: Copy on Write
 - UNIX fork gives *copy* of parent address space to child
 - » Address spaces disconnected after child created
 - How to do this cheaply?
 - » Make copy of parent's page tables (point at same memory)
 - » Mark entries in both sets of page tables as read-only
 - » Page fault on write creates two copies
- Usage Example: Zero Fill On Demand
 - New data pages must carry no information (say be zeroed)
 - Mark PTEs as invalid; page fault on use gets zeroed page
 - Often, OS creates zeroed pages in background

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.26

How is the translation accomplished?



- What, exactly happens inside MMU?
- One possibility: Hardware Tree Traversal
 - For each virtual address, takes page table base pointer and traverses the page table in hardware
 - Generates a "Page Fault" if it encounters invalid PTE
 - » Fault handler will decide what to do
 - » More on this next lecture
 - Pros: Relatively fast (but still many memory accesses!)
 - Cons: Inflexible, Complex hardware
- Another possibility: Software
 - Each traversal done in software
 - Pros: Very flexible
 - Cons: Every translation must invoke Fault!
- In fact, need way to cache translations for either case!

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.27

Summary (1/2)

- Memory is a resource that must be shared
 - Controlled Overlap: only shared when appropriate
 - Translation: Change Virtual Addresses into Physical Addresses
 - Protection: Prevent unauthorized Sharing of resources
- Segment Mapping
 - Segment registers within processor
 - Segment ID associated with each access
 - » Often comes from portion of virtual address
 - » Can come from bits in instruction instead (x86)
 - Each segment contains base and limit information
 - » Offset (rest of address) adjusted by adding base

3/10/08

Joseph CS162 ©UCB Spring 2008

Lec 12.28

Summary (2/2)

- **Page Tables**
 - Memory divided into fixed-sized chunks of memory
 - Virtual page number from virtual address mapped through page table to physical page number
 - Offset of virtual address same as physical address
 - Large page tables can be placed into virtual memory
- **Multi-Level Tables**
 - Virtual address mapped to series of tables
 - Permit sparse population of address space
- **Inverted page table**
 - Size of page table related to physical memory size