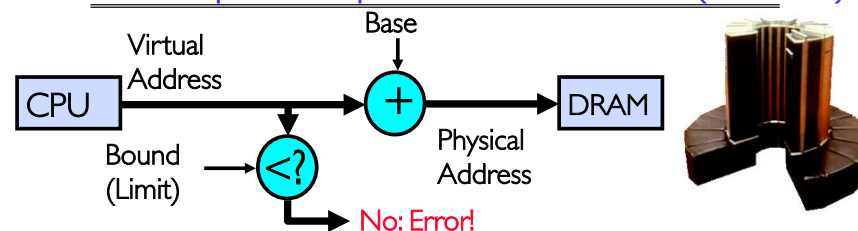


CSI 62
Operating Systems and
Systems Programming
Lecture 12

Address Translation

March 2nd, 2016
Prof. Anthony D. Joseph
<http://cs162.eecs.Berkeley.edu>

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



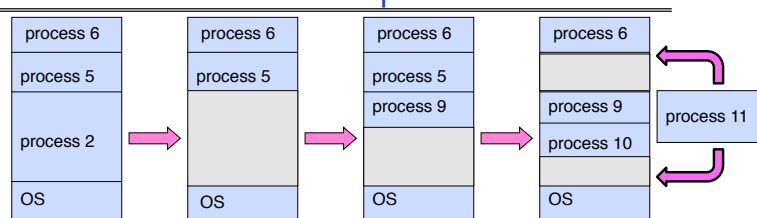
- Could use base/bounds for **dynamic address translation** – translation happens at execution:
 - Alter address of every load/store by adding “base”
 - Generate error if address bigger than limit
- This gives program the illusion that it is running on its own dedicated machine, with memory starting 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

3/2/16

Joseph CSI 62 ©UCB Spring 2016

Lec 12.2

Review: Issues with Simple B&B Method



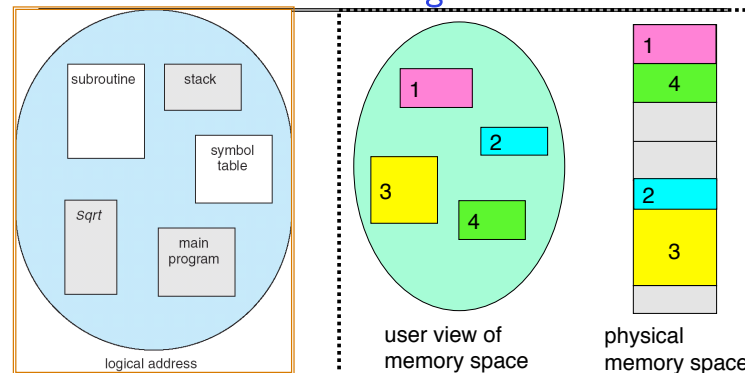
- Fragmentation problem over time
 - Not every process is same size → memory becomes fragmented
- Missing support for sparse address space
 - Would like to have multiple chunks/program (Code, Data, Stack)
- Hard to do inter-process sharing
 - Want to share code segments when possible
 - Want to share memory between processes
 - Helped by providing multiple segments per process

3/2/16

Joseph CSI 62 ©UCB Spring 2016

Lec 12.3

More Flexible Segmentation



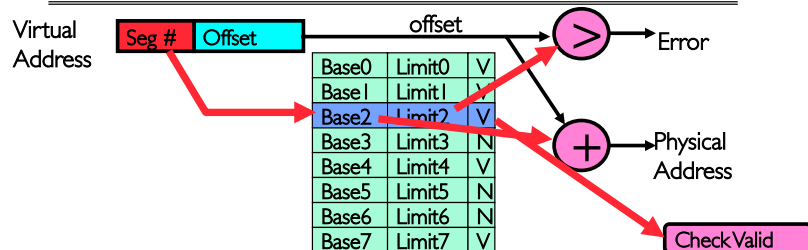
- Logical View: multiple separate segments
 - Typical: Code, Data, Stack
 - Others: memory sharing, etc
- Each segment is given region of contiguous memory
 - Has a base and limit
 - Can reside anywhere in physical memory

3/2/16

Joseph CSI 62 ©UCB Spring 2016

Lec 12.4

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" (valid / not valid)?
 - Can mark segments as invalid; requires check as well

3/2/16

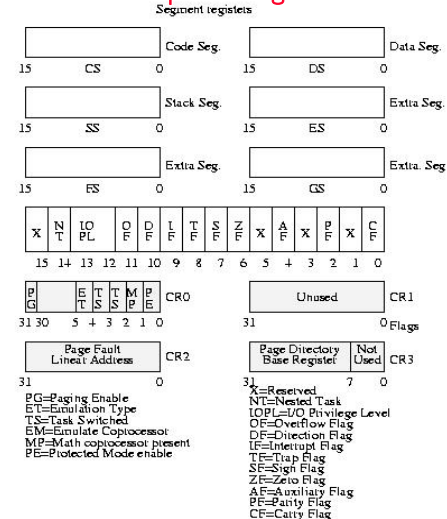
Joseph CS162 ©UCB Spring 2016

Lec 12.5

Intel x86 Special Registers



80386 Special Registers



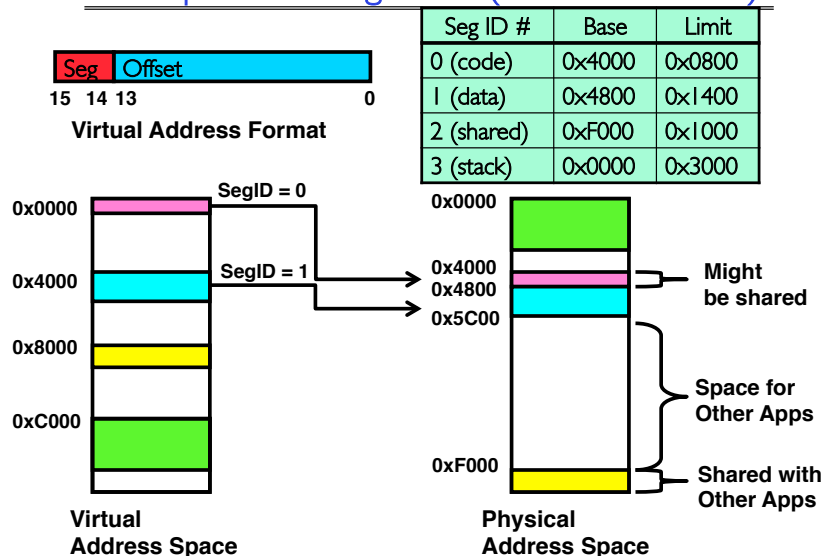
Typical Segment Register
Current Priority is RPL
Of Code Segment (CS)

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.6

Example: Four Segments (16 bit addresses)



3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.7

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

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

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.8

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”)

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.9

Problems with Segmentation

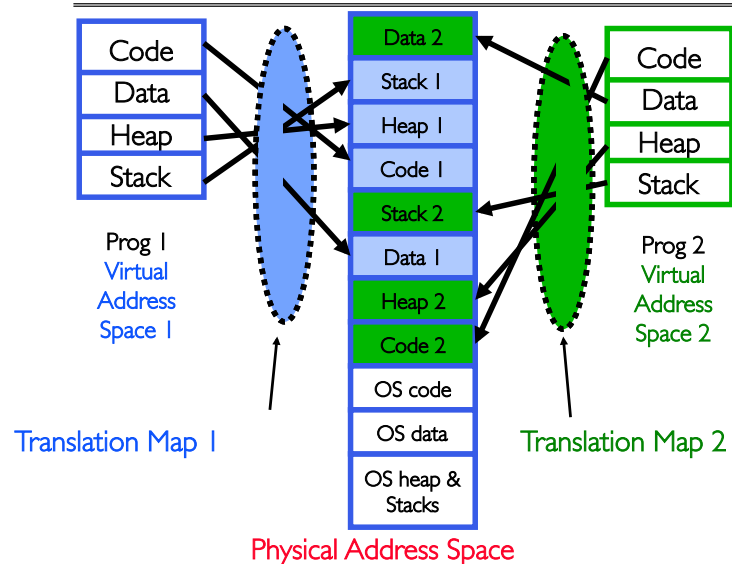
- Must fit variable-sized chunks into physical memory
- May move processes multiple times to fit everything
- Limited options for swapping to disk
- **Fragmentation**: wasted space
 - **External**: free gaps between allocated chunks
 - **Internal**: don't need all memory within allocated chunks

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.10

Recall: General Address Translation



3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.11

Paging: Physical Memory in Fixed Size Chunks

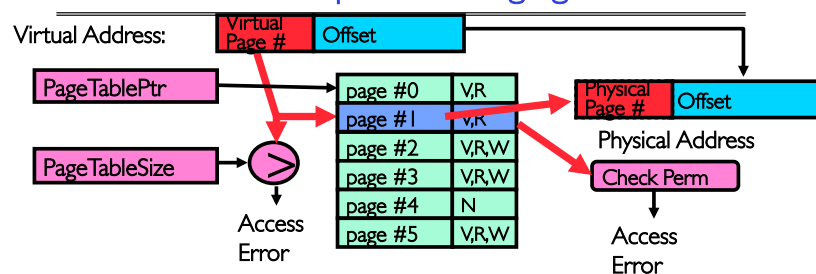
- Solution to fragmentation from segments?
 - Allocate physical memory in fixed size chunks (“pages”)
 - Every chunk of physical memory is equivalent
 - » Can use simple vector of bits to handle allocation: 00110001110001101 ... 110010
 - » Each bit represents page of physical memory
1⇒allocated, 0⇒free
- 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

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.12

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 \Rightarrow 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

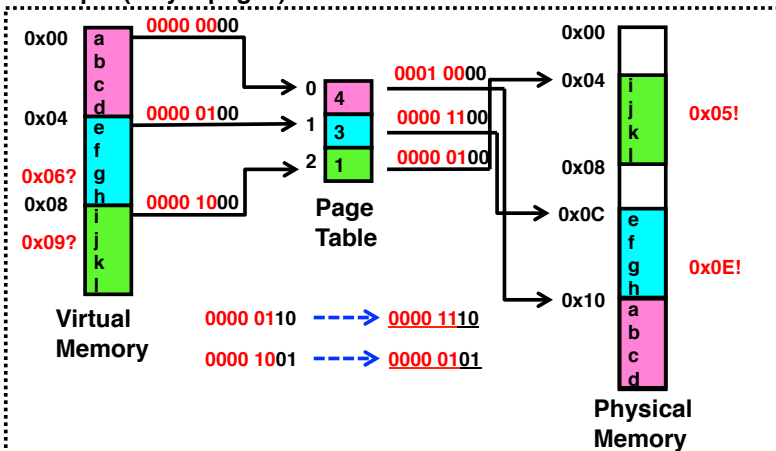
3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.13

Simple Page Table Example

Example (4 byte pages)

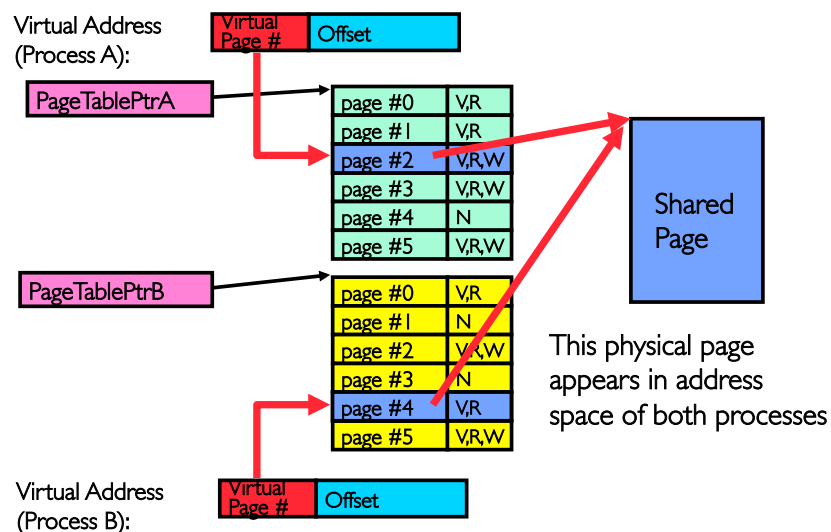


3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.14

What about Sharing?



3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.15

Administrivia

- Upcoming deadlines:
 - Project I final code due Fri 3/4, final report due Mon 3/7
- Midterm next week Wed 3/9 6-7:30 10 Evans and 155 Dwinelle
 - Midterm review session: Sun 3/6 2-5PM at 2060 VLSB
 - Rooms assignment: aa-eh 10 Evans, ej-oa 155 Dwinelle
 - Lectures (including #12), project, homeworks readings, textbook
 - No books, no calculators, one double-side page handwritten notes
 - No class that day, extra office hours

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.16

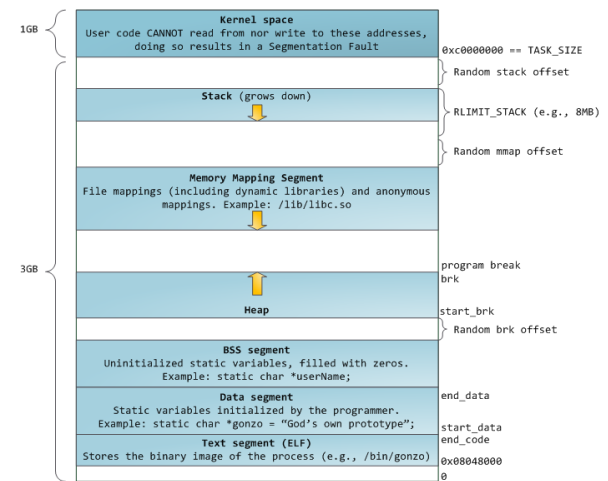
BREAK

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.17

Memory Layout for Linux 32-bit



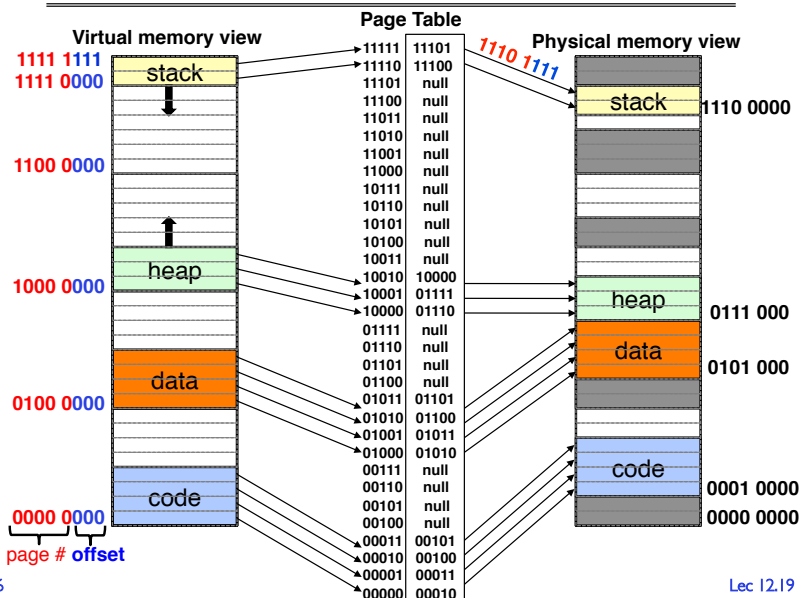
<http://static.duartes.org/img/blogPosts/linuxFlexibleAddressSpaceLayout.png>

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.18

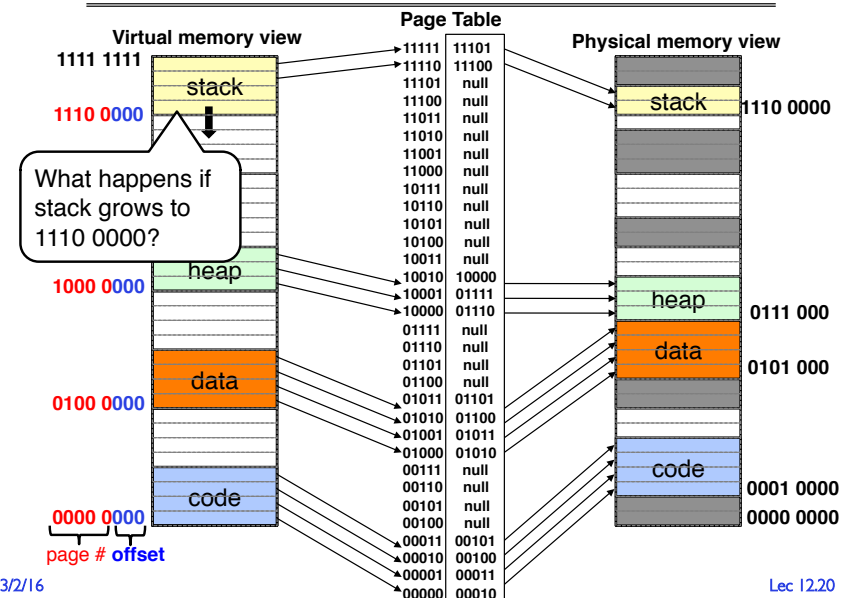
Summary: Paging



3/2/16

Lec 12.19

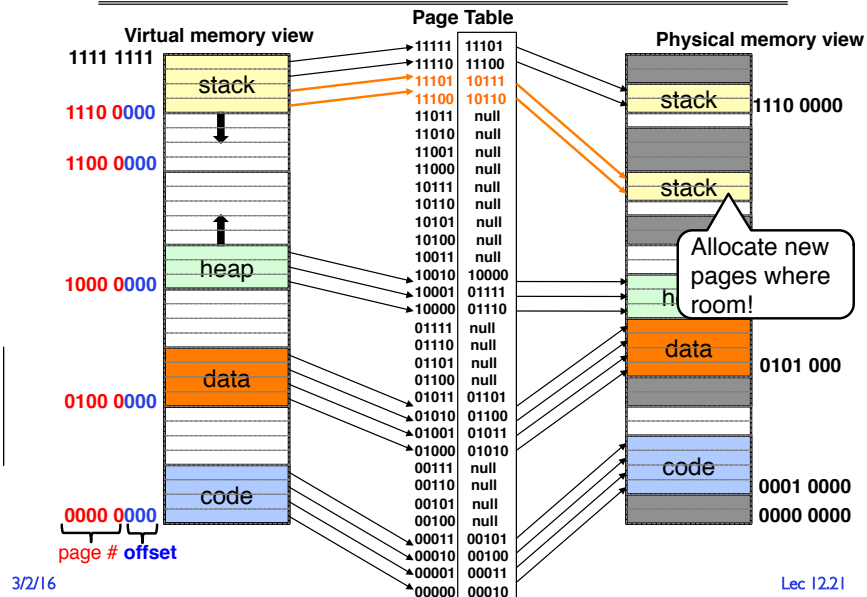
Summary: Paging



3/2/16

Lec 12.20

Summary: Paging



Page Table Discussion

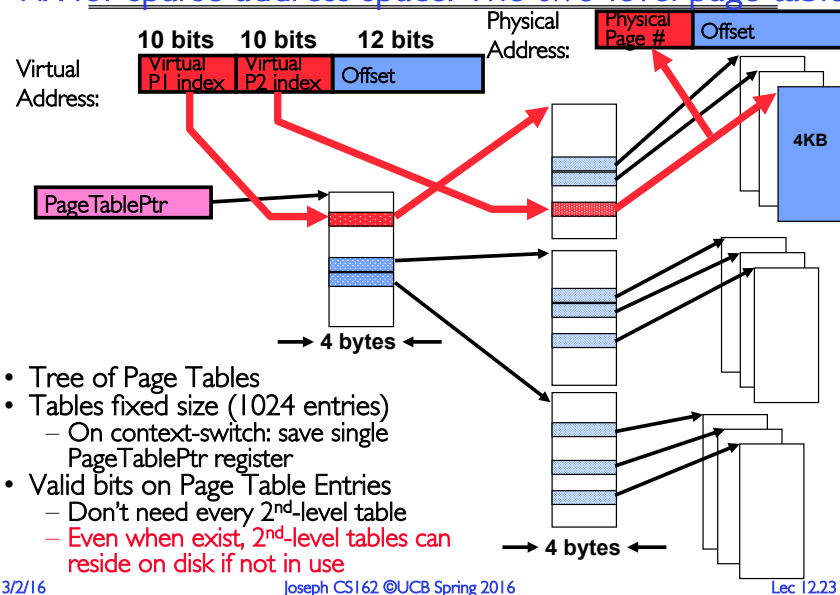
- What needs to be switched on a context switch?
 - Page table pointer and limit
- Analysis
 - Pros
 - Simple memory allocation
 - Easy to Share
 - Con: What if address space is sparse?
 - E.g. on UNIX, code starts at 0, stack starts at $(2^{31}-1)$.
 - With 1K pages, need 2 million page table entries!
 - Con: What if table really big?
 - Not all pages used all the time \Rightarrow would be nice to have working set of page table in memory
- How about combining paging and segmentation?

3/2/16

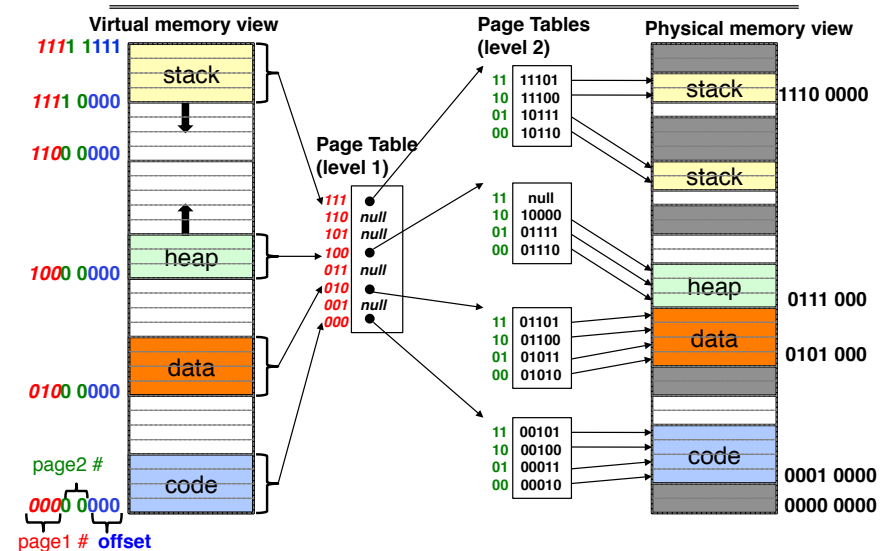
Joseph CS162 ©UCB Spring 2016

Lec 12.22

Fix for sparse address space: The two-level page table



Summary: Two-Level Paging

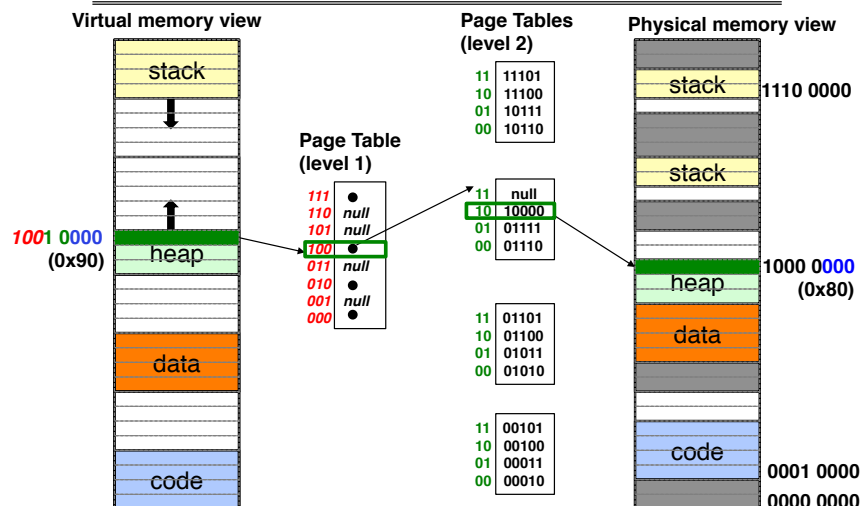


3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.24

Summary: Two-Level Paging



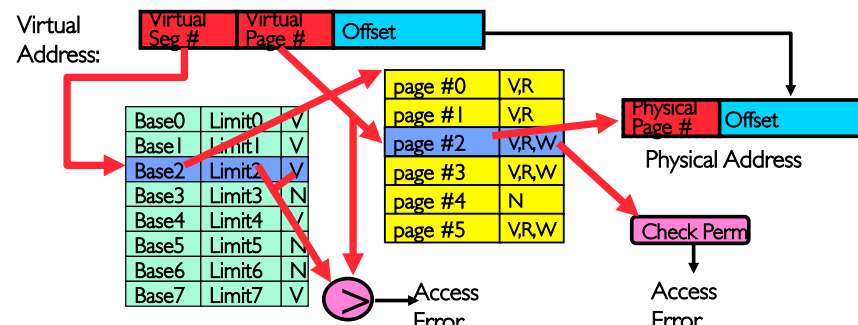
3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.25

Multi-level Translation: Segments + Pages

- What about a tree of tables?
 - Lowest level page table ⇒ memory still allocated with bitmap
 - Higher levels often segmented
- Could have any number of levels. Example (top segment):



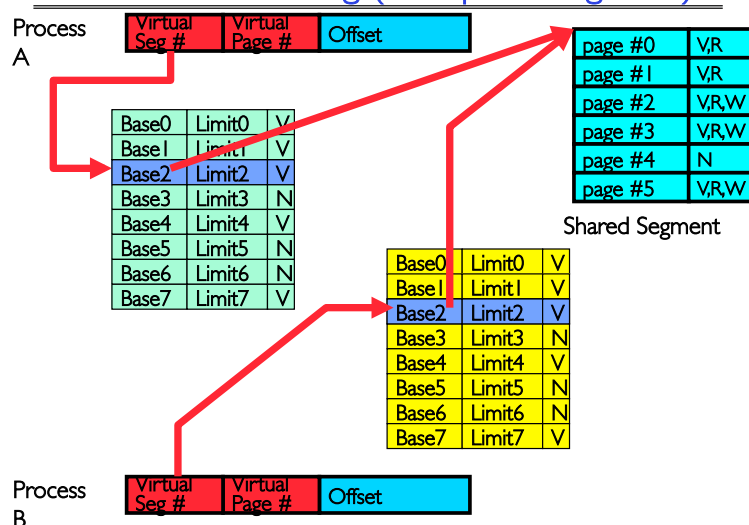
- What must be saved/restored on context switch?
 - Contents of top-level segment registers (for this example)
 - Pointer to top-level table (page table)

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.26

What about Sharing (Complete Segment)?



3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.27

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!

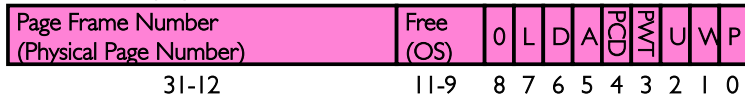
3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.28

What is in a Page Table Entry?

- 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

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.29

Examples of how to use a PTE

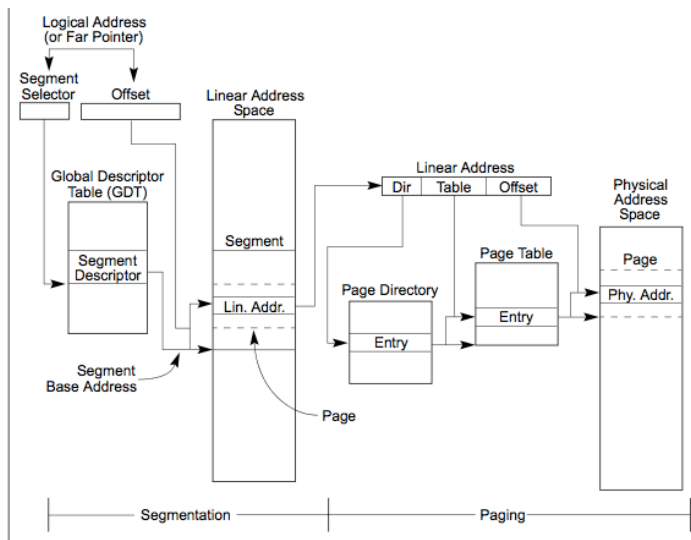
- How do we use the PTE?
 - Invalid PTE can imply different things:
 - Region of address space is actually invalid or
 - Page/directory is just somewhere else than memory
 - Validity checked first
 - OS can use other (say) 31 bits for location info
- Usage Example: Demand Paging
 - Keep only active pages 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

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.30

Making it real: X86 Memory model with segmentation (16/32-bit)



3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.31

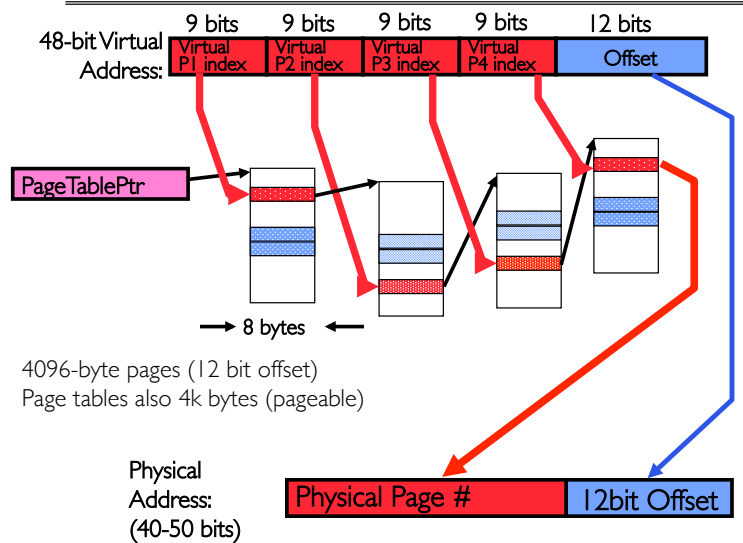
BREAK

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.34

X86_64: Four-level page table!

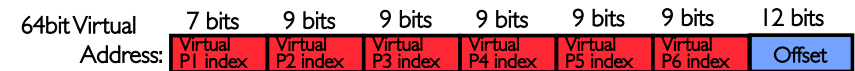


3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.35

IA64: 64bit addresses: Six-level page table!?



No!

Too slow
Too many almost-empty tables

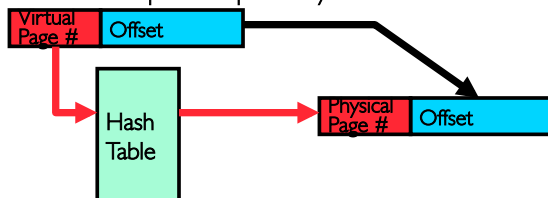
3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.36

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



- 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!

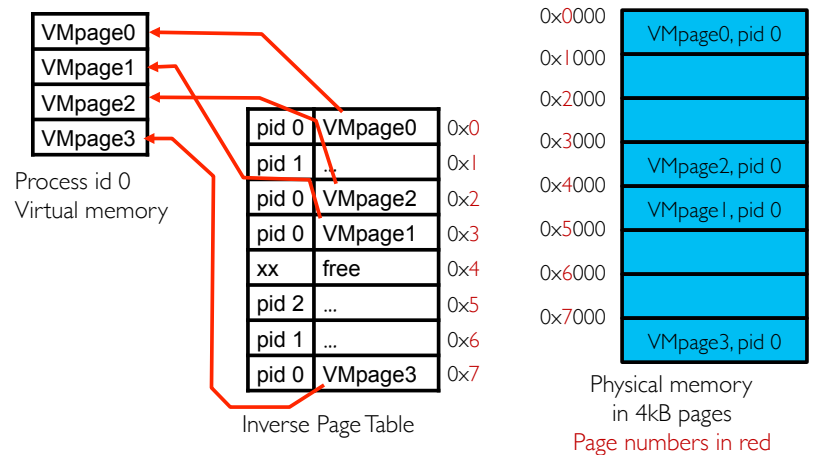
3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.37

IA64: Inverse Page Table (IPT)

Idea: index page table by physical pages instead of VM



3/2/16

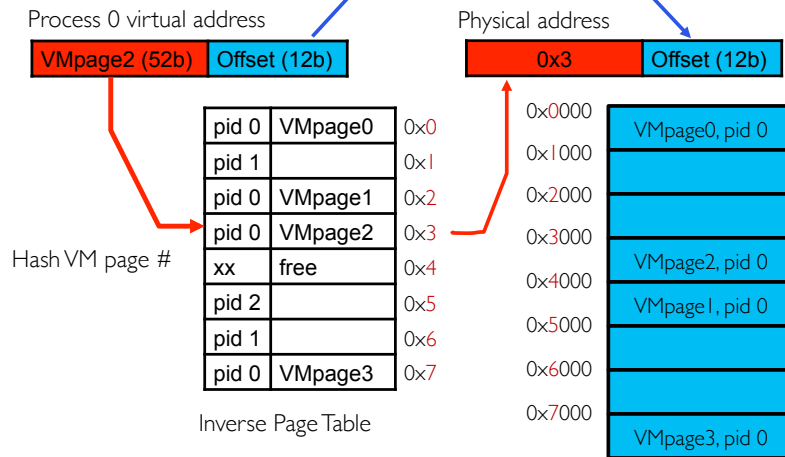
Joseph CS162 ©UCB Spring 2016

Lec 12.38

IPT address translation

- Need an associative map from VM page to IPT address:

– Use a hash map

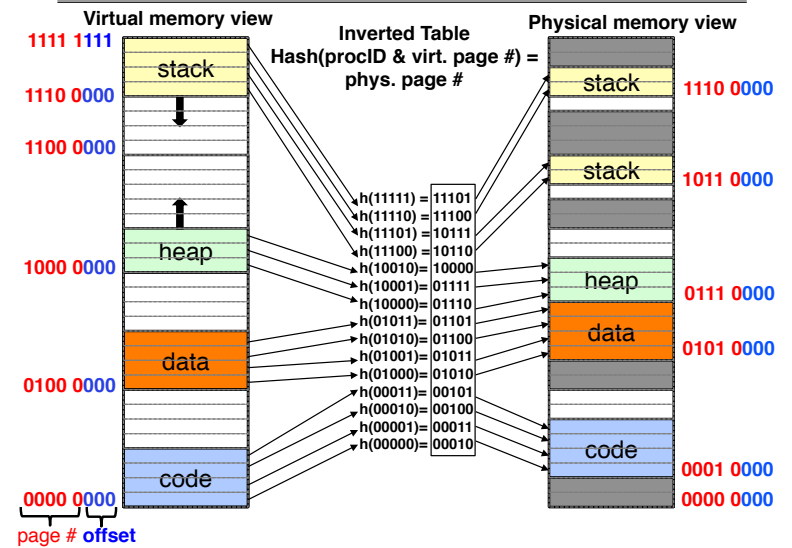


3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.39

Summary: Inverted Table



3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.40

Address Translation Comparison

	Advantages	Disadvantages
Simple Segmentation	Fast context switching: Segment mapping maintained by CPU	External fragmentation
Paging (single-level page)	No external fragmentation, fast easy allocation	Large table size ~ virtual memory Internal fragmentation
Paged segmentation	Table size ~ # of pages in virtual memory , fast easy allocation	Multiple memory references per page access
Two-level pages		
Inverted Table	Table size ~ # of pages in physical memory	Hash function more complex

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.41

Summary

- 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
- 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

3/2/16

Joseph CS162 ©UCB Spring 2016

Lec 12.42