

CS162 Operating Systems and Systems Programming Lecture 9

Address Translation

October 2, 2013

Anthony D. Joseph and John Canny

<http://inst.eecs.berkeley.edu/~cs162>

Goals for Today

- Review of Scheduling
- Address Translation Schemes
 - Segmentation
 - Paging
 - Multi-level translation
 - Paged page tables
 - Inverted page tables

Note: Some slides and/or pictures in the following are adapted from slides ©2005 Silberschatz, Galvin, and Gagne. Slides courtesy of Anthony D. Joseph, John Kubiawicz, AJ Shankar, George Necula, Alex Aiken, Eric Brewer, Ras Bodik, Ion Stoica, Doug Tygar, and David Wagner.

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.2

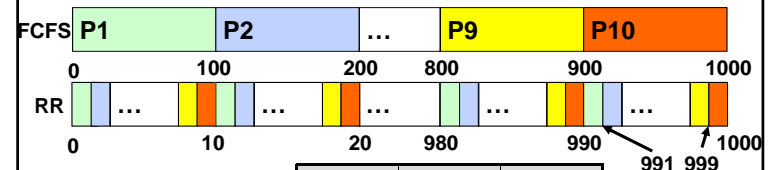
Lecture 8: Scheduling

- Please view lecture 8 on Youtube if you havent already:
- http://www.youtube.com/watch?v=Xr_5uYkWI2I&list=PL-XXv-cvA_iDrt_oPWfQ4-fjHm2KSSOPq&index=8
- Link is at top of class home page (and in Piazza).

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.3

FCFS and Round Robin

- Assuming zero-cost context-switching time, how does RR compare with FCFS?
- Simple example: 10 jobs, each takes 100s of CPU time
RR scheduler quantum of 1s
All jobs start at the same time



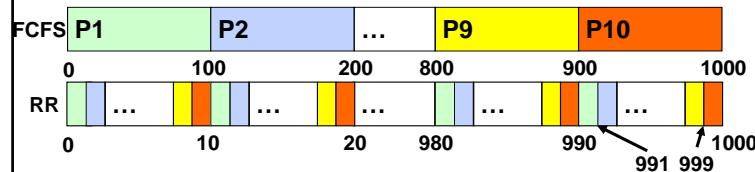
- Completion Times:
- FCFS average 550
- RR average 995.5!

Job #	FCFS	RR
1	100	991
2	200	992
...
9	900	999
10	1000	1000

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.4

Comparisons between FCFS and Round Robin

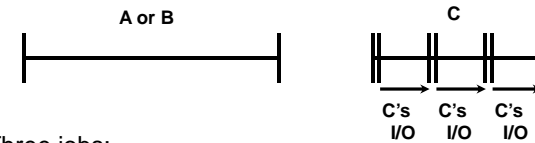
- Assuming zero-cost context-switching time, is RR always better than FCFS?
- Simple example: 10 jobs, each takes 100s of CPU time
RR scheduler quantum of 1s
All jobs start at the same time



- Both RR and FCFS finish at the same time
- Average response time is worse under RR!
 - Bad when all jobs same length, does better with short jobs
- Also: Cache state must be shared between all jobs with RR but can be devoted to each job with FCFS
 - Total time for RR longer even for zero-cost switch!

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.5

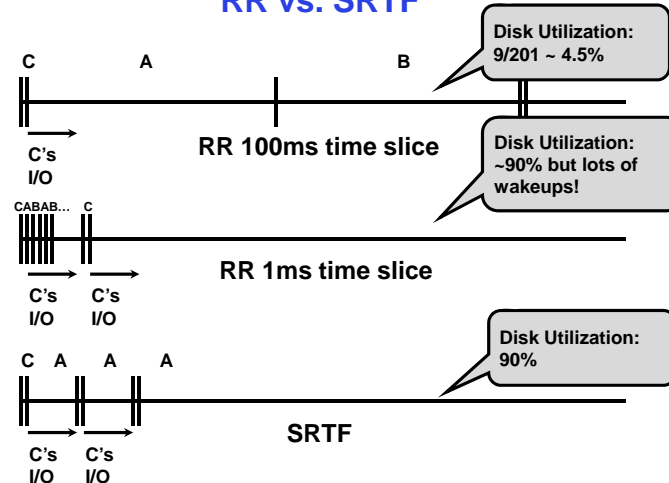
Example of Shortest Remaining Time First



- Three jobs:
 - A, B: CPU bound, each run for a week
 - C: I/O bound, loop 1ms CPU, 9ms disk I/O
 - If only one at a time, C uses 90% of the disk, A or B use 100% of the CPU
- With FCFS:
 - Once A or B get in, keep CPU for one week each
- What about RR or SRTF?
 - Easier to see with a timeline

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.6

RR vs. SRTF



10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.7

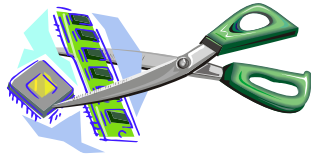
SRTF Tradeoffs

- Starvation
 - SRTF can lead to starvation if many small jobs!
 - Large jobs never get to run
- Somehow need to predict future
 - How can we do this?
 - Some systems ask the user
 - When you submit a job, have to say how long it will take
 - To stop cheating, system kills job if takes too long
 - But: even non-malicious users have trouble predicting runtime of their jobs
- Bottom line, can't really know how long job will take
 - However, can use SRTF as a yardstick for measuring other policies
 - Optimal, so can't do any better
- SRTF Pros & Cons
 - Optimal (average response time) (+)
 - Hard to predict future (-)
 - Unfair (-)



10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013

Virtualizing Resources



- Physical Reality: Processes/Threads share the same hardware
 - Need to multiplex CPU (CPU Scheduling)
 - Need to multiplex use of Memory (Today)
- Why worry about memory multiplexing?
 - The complete working state of a process and/or kernel is defined by its data in memory (and registers)
 - Consequently, cannot just let different processes use the same memory
 - Probably don't want different processes to even have access to each other's memory (protection)

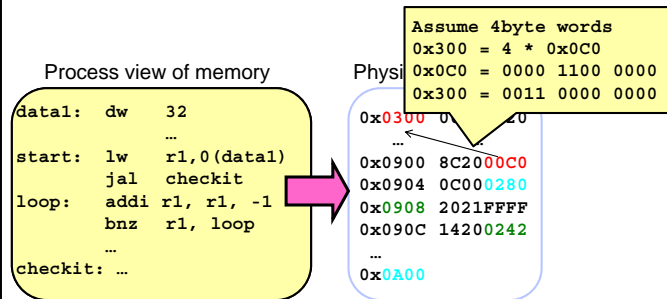
10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.9

Important Aspects of Memory Multiplexing

- **Controlled overlap:**
 - Processes should not collide in physical memory
 - Conversely, would like the ability to share memory when desired (for communication)
- **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
- **Translation:**
 - Ability to translate accesses from one address space (virtual) to a different one (physical)
 - When translation exists, process uses virtual addresses, physical memory uses physical addresses

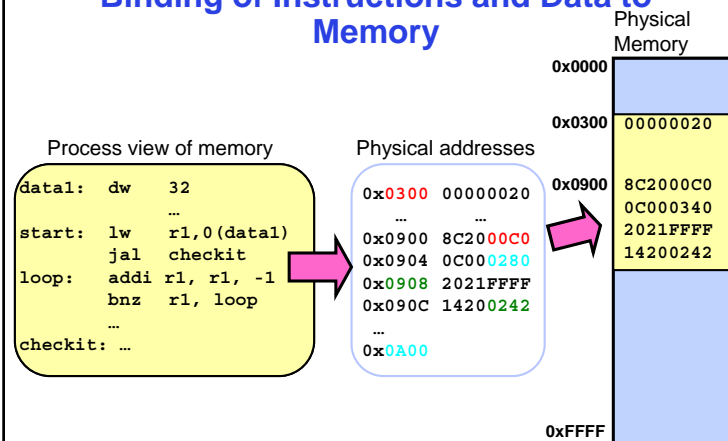
10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.10

Binding of Instructions and Data to Memory

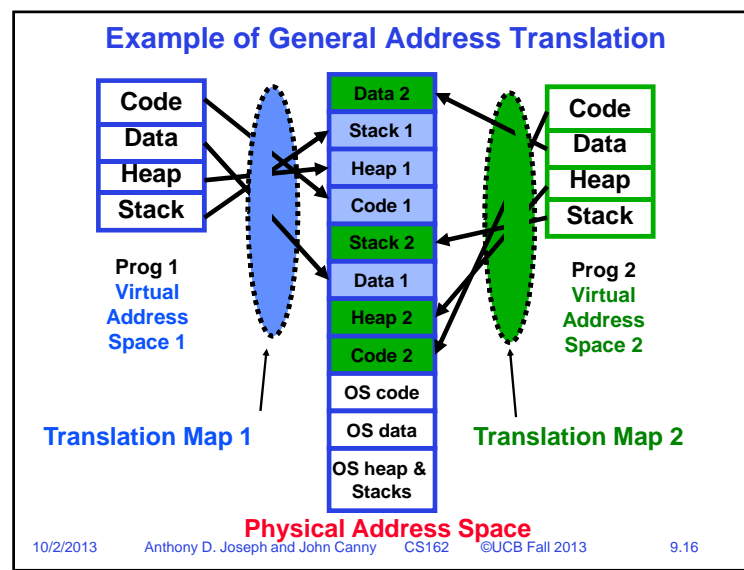
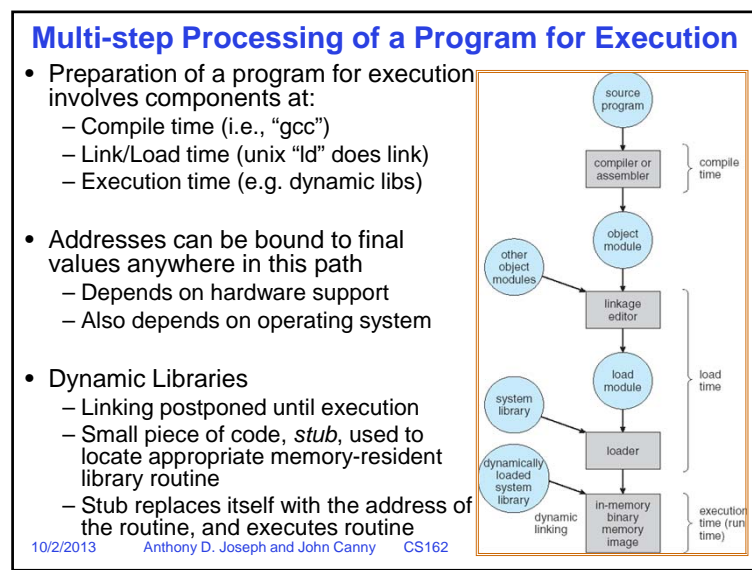
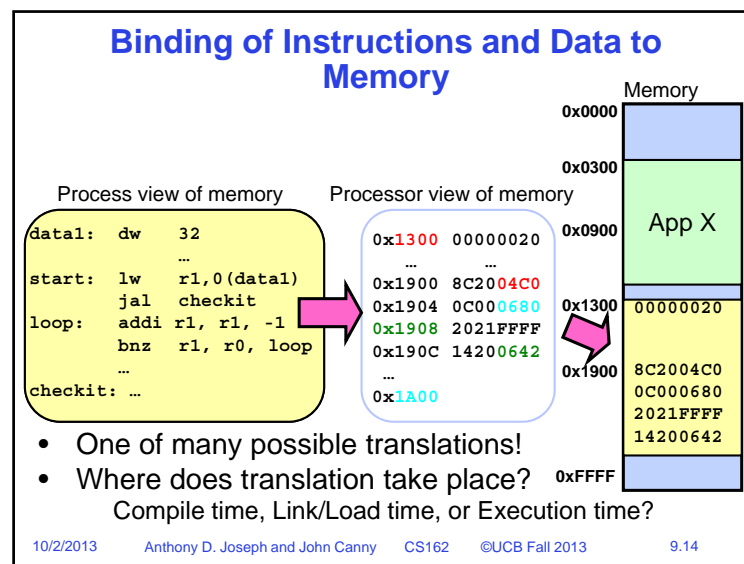
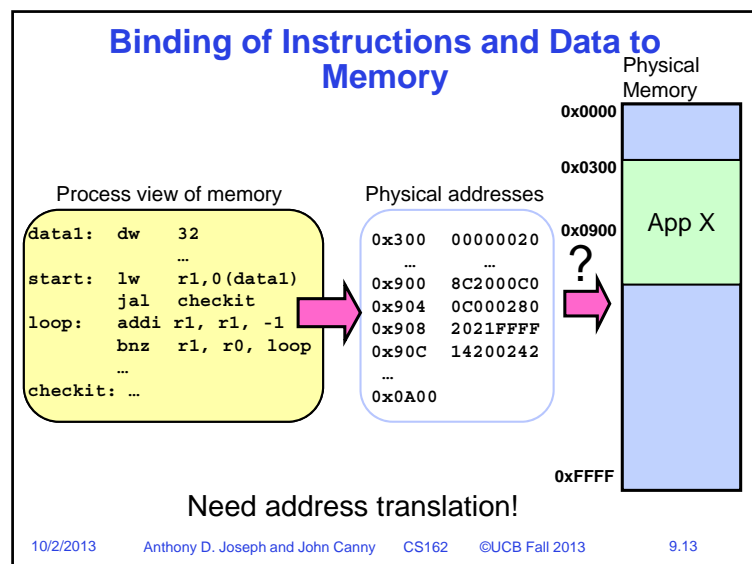


10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.11

Binding of Instructions and Data to Memory



10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.12



Two Views of Memory

- Address Space:
 - All the addresses and state a process can touch
 - Each process and kernel has different address space
- Consequently, two views of memory:
 - View from the CPU (what program sees, virtual memory)
 - View from memory (physical memory)
 - Translation box (MMU) converts between the two views
- Translation helps to implement protection
 - If task A cannot even gain access to task B's data, no way for A to adversely affect B
- With translation, every program can be linked/loaded into same region of user address space

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.17

Uniprogramming (MS-DOS)

Starting MS-DOS...
 C:\>_

- Uniprogramming (no Translation or Protection)
 - Application always runs at same place in physical memory since only one application at a time
 - Application can access any physical address

- Application given illusion of dedicated machine by giving it reality of a dedicated machine

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.18

Multiprogramming (First Version)

- Multiprogramming without Translation or Protection
 - Must somehow prevent address overlap between threads

- Trick: Use Loader/Linker: Adjust addresses while program loaded into memory (loads, stores, jumps)
 - » Everything adjusted to memory location of program
 - » Translation done by a linker-loader
 - » Was pretty common in early days
- With this solution, no protection: bugs in any program can cause other programs to crash or even the OS

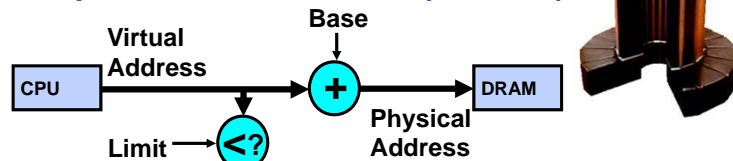
10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.19

Multiprogramming (Version with Protection)

- Can we protect programs from each other without translation?
 - Yes: use two special registers *BaseAddr* and *LimitAddr* to prevent user from straying outside designated area
 - » If user tries to access an illegal address, cause an error
 - During switch, kernel loads new base/limit from TCB (Thread Control Block)
 - » User not allowed to change base/limit registers

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.20

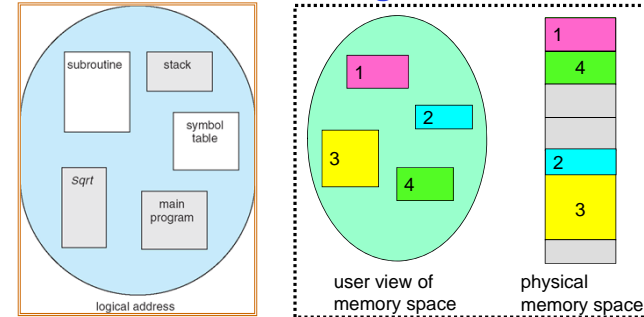
Simple Base and Bounds (CRAY-1)



- Could use base/limit for **dynamic address translation** (often called “segmentation”) – 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

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.21

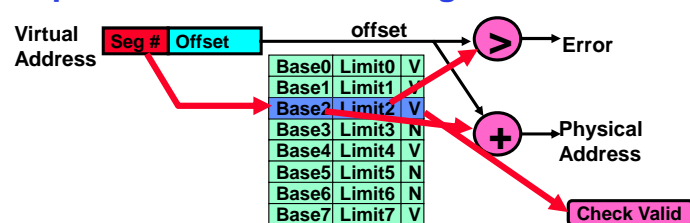
More Flexible Segmentation



- 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

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.22

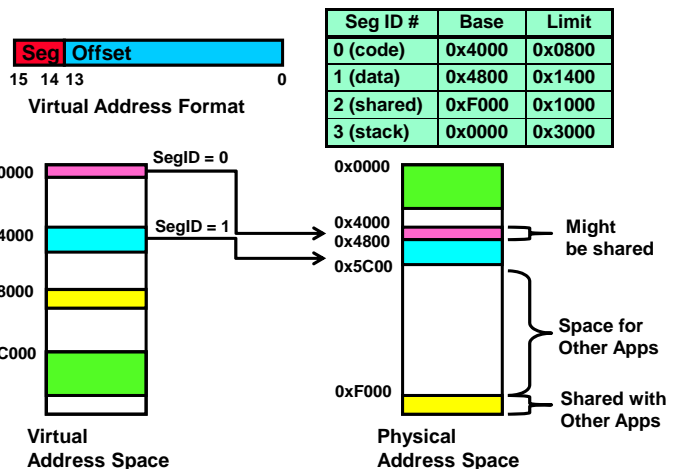
Implementation of Multi-Segment Model



- 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

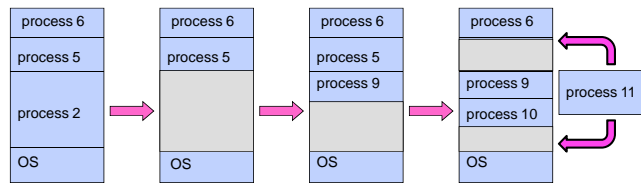
10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.23

Example: Four Segments (16 bit addresses)



10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.24

Issues with Simple Segmentation Method

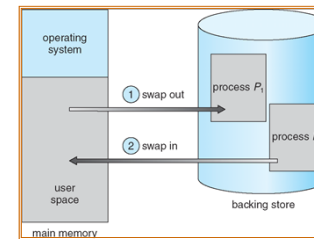


- Fragmentation problem
 - Not every process is the same size
 - Over time, memory space becomes fragmented
- Hard to do inter-process sharing
 - Want to share code segments when possible
 - Want to share memory between processes
 - Helped by providing multiple segments per process

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.25

Schematic View of Swapping

- Q: What if not all processes fit in memory?
- A: Swapping: Extreme form of Context Switch
 - In order to make room for next process, some or all of the previous process is moved to disk
 - 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

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.26

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

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.27

Administrivia

- Project #1 code due Tuesday Oct 8 by 11:59pm
- Design doc (submit proj1-final-design) and group evals (Google Docs form) due Wed 10/9 at 11:59PM
- Midterm #1 is Monday Oct 21 5:30-7pm in **145 Dwinelle (A-L)** and **2060 Valley LSB (M-Z)**
 - Closed book, one handwritten page of notes, no calculators, smartphones, Google glass etc.
 - Covers lectures #1-13, readings, handouts, and projects 1 and 2
 - Review session **390 HEARST MINING Friday, October 18, 2013; 5:00 PM-7:00 PM**

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.28

5min Break

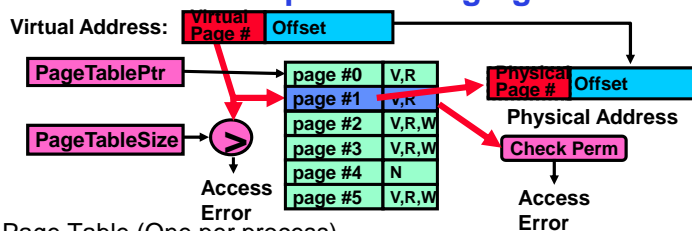
10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.29

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

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.30

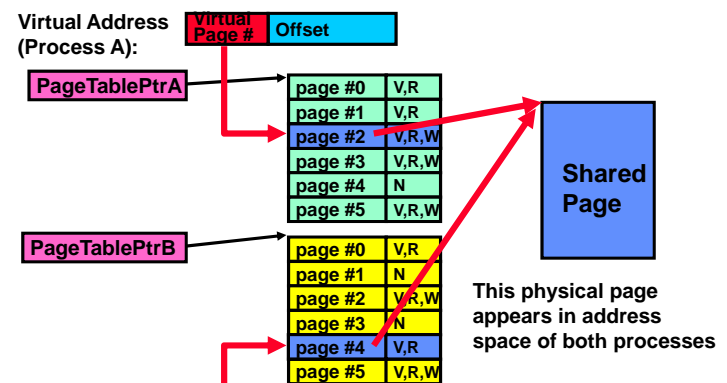
How to Implement Paging?



- 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

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.31

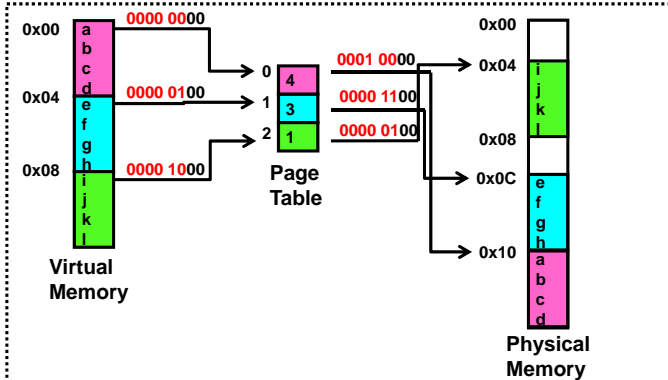
What about Sharing?



10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.32

Simple Page Table Example

Example (4 byte pages)



10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.33

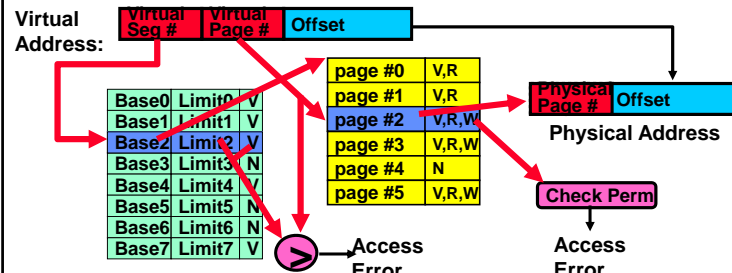
Page Table Discussion

- 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 4 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?

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.34

Multi-level Translation

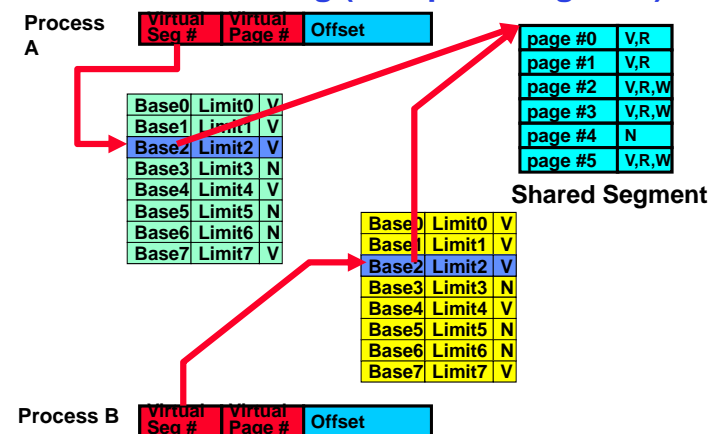
- What about a tree of tables?
 - Lowest level page table \Rightarrow memory still allocated with bitmap
 - Higher levels often segmented
- Could have any number of levels. Example (top segment):



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

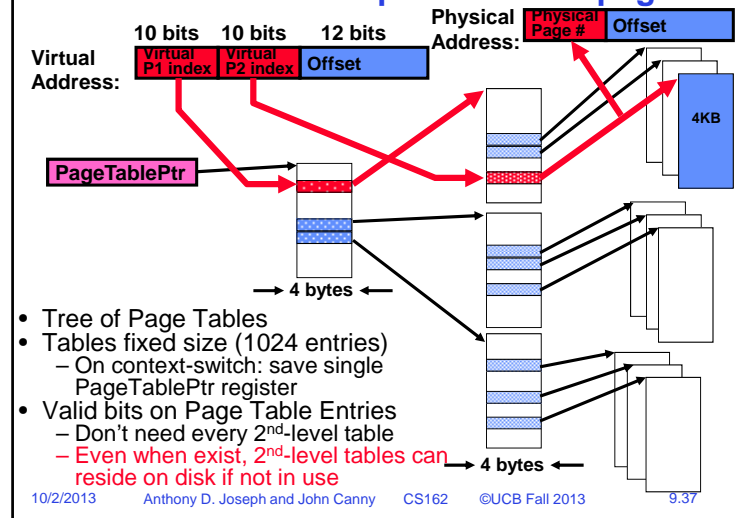
10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.35

What about Sharing (Complete Segment)?



10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.36

Another common example: two-level page table



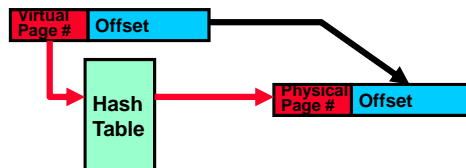
Multi-level Translation Analysis

- Pros:
 - Only need to allocate as many page table entries as we need for application – size is proportional to usage
 - 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!

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.38

Inverted Page Table

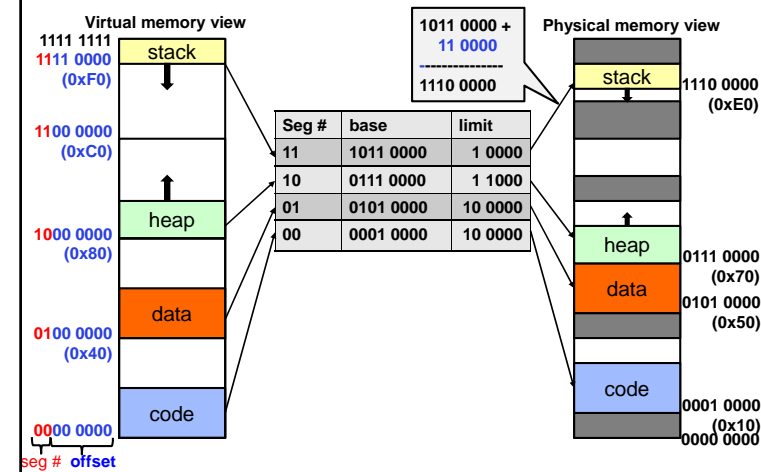
- With all previous examples (“Forward Page Tables”)
 - Size of page table is at least as large as amount of virtual memory allocated to processes
 - Physical memory may be much less
 - Much of process space may be out on disk or not in use



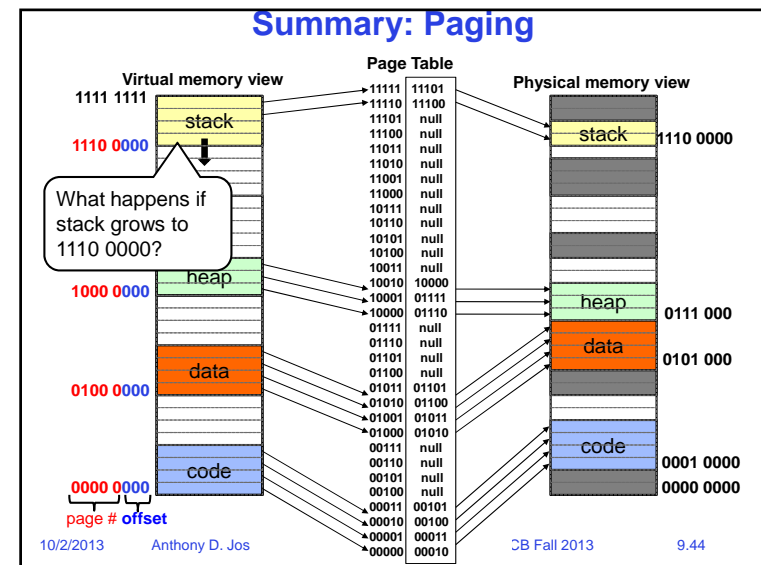
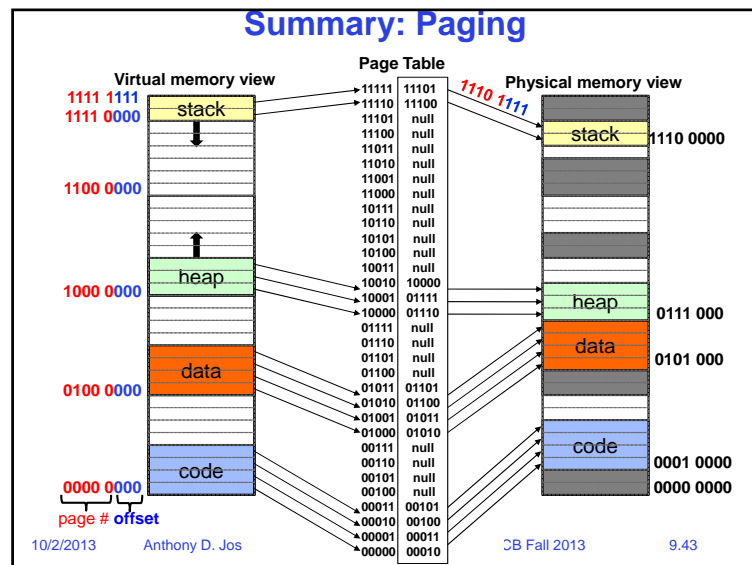
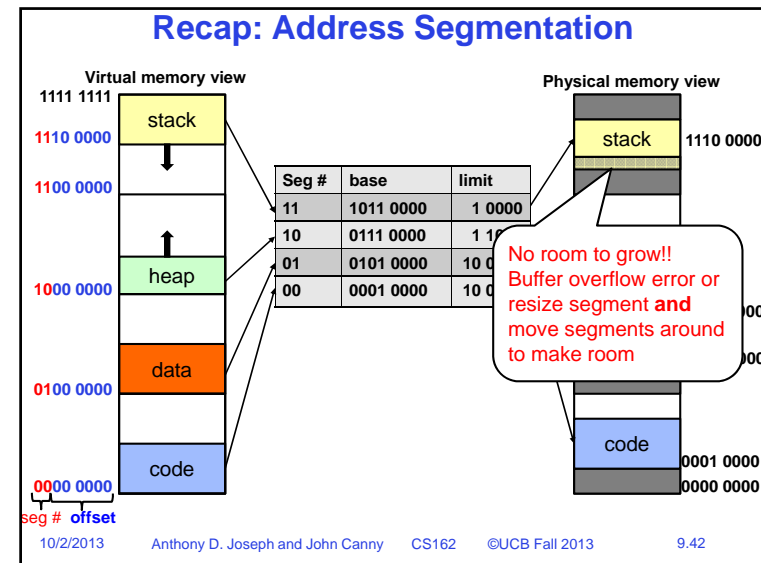
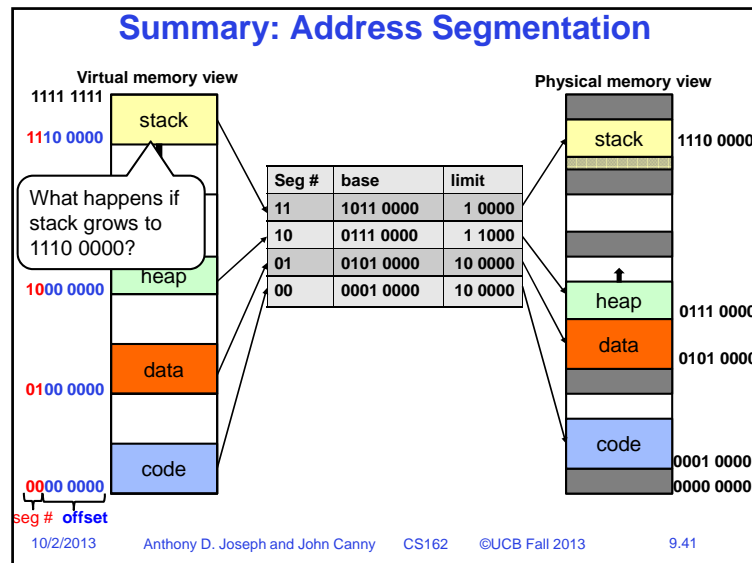
- Answer: use a hash table
 - Called an “Inverted Page Table”
 - Size is independent of virtual address space
 - Directly related to amount of physical memory
 - Very attractive option for 64-bit address spaces (IA64)
- Cons: Complexity of managing hash changes
 - Often in hardware!

10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.39

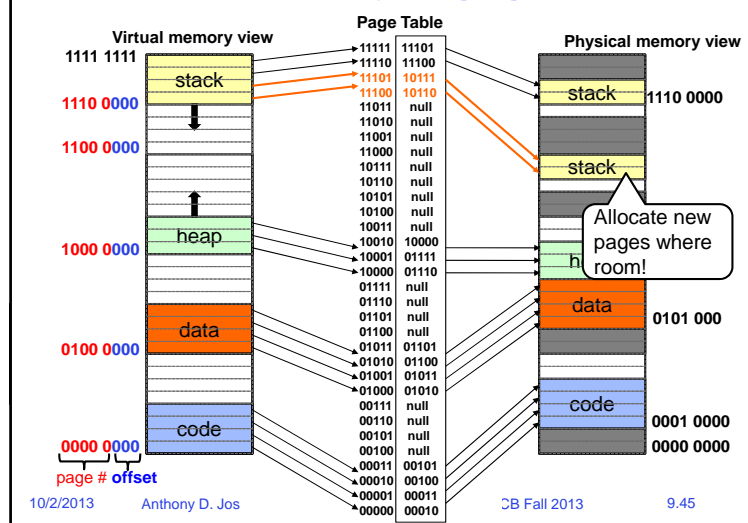
Summary: Address Segmentation



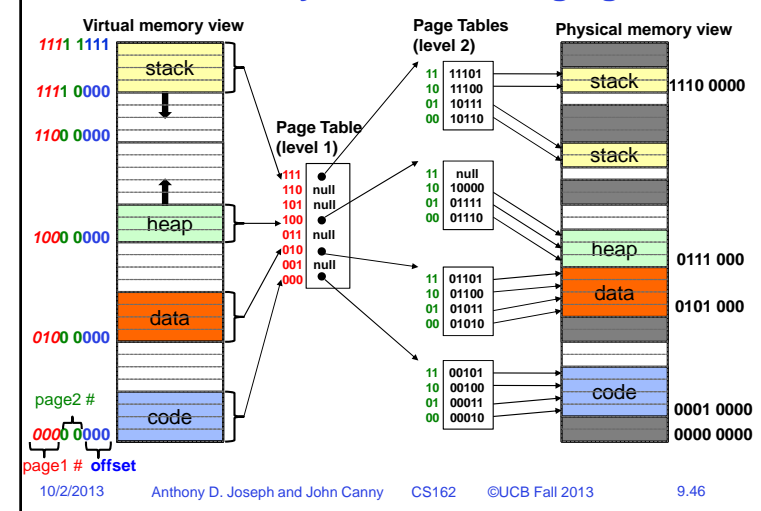
10/2/2013 Anthony D. Joseph and John Canny CS162 ©UCB Fall 2013 9.40



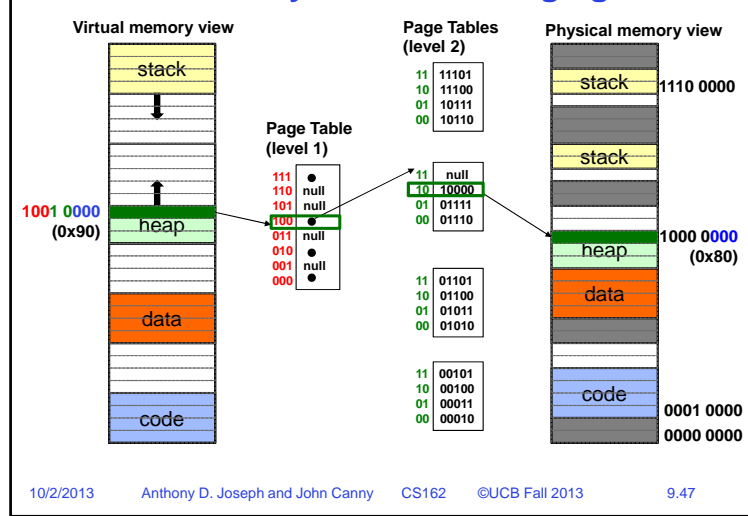
Summary: Paging



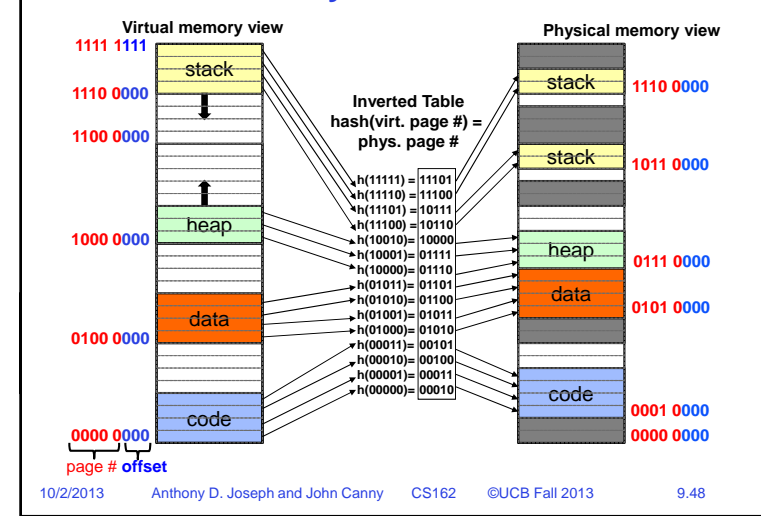
Summary: Two-Level Paging



Summary: Two-Level Paging



Summary: Inverted Table



Address Translation Comparison

	Advantages	Disadvantages
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
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

Summary

- Memory is a resource that must be multiplexed
 - Controlled Overlap: only shared when appropriate
 - Translation: Change virtual addresses into physical addresses
 - Protection: Prevent unauthorized sharing of resources
- Simple Protection through segmentation
 - Base+limit registers restrict memory accessible to user
 - Can be used to translate as well
- Page Tables
 - Memory divided into fixed-sized chunks of memory
 - Offset of virtual address same as physical address
- 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 mem. size

10/2/2013

Anthony D. Joseph and John Canny

CS162

©UCB Fall 2013

9.50