# CS162
## Operating Systems and Systems Programming
## Lecture 12

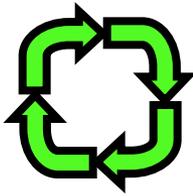## Address Translation

March 4th, 2015

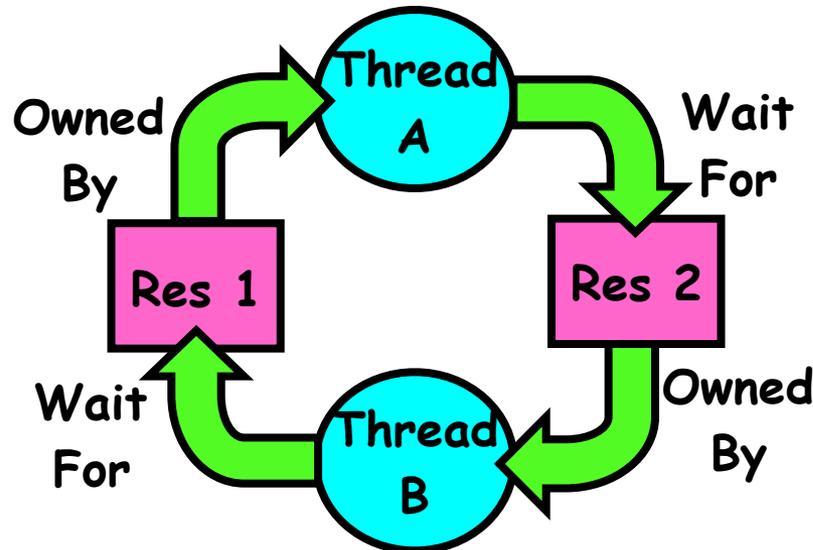Prof. John Kubiatowicz

http://cs162.eecs.Berkeley.edu

# Recall: Starvation vs Deadlock

- **Starvation vs. Deadlock**
  - Starvation: thread waits indefinitely
    - » Example, low-priority thread waiting for resources constantly in use by high-priority threads
  - Deadlock: circular waiting for resources
    - » Thread A owns Res 1 and is waiting for Res 2
      Thread B owns Res 2 and is waiting for Res 1



  - Deadlock $\Rightarrow$ Starvation but not vice versa
    - » Starvation can end (but doesn't have to)
    - » Deadlock can't end without external intervention

# Recall: Four requirements for Deadlock

- ## Mutual exclusion
  - Only one thread at a time can use a resource.
- ## Hold and wait
  - Thread holding at least one resource is waiting to acquire additional resources held by other threads
- ## No preemption
  - Resources are released only voluntarily by the thread holding the resource, after thread is finished with it
- ## Circular wait
  - There exists a set $\{T_1, \ldots, T_n\}$ of waiting threads
    - » $T_1$ is waiting for a resource that is held by $T_2$
    - » $T_2$ is waiting for a resource that is held by $T_3$
    - » …
    - » $T_n$ is waiting for a resource that is held by $T_1$

# Recall: Address translation

Virtual Addresses **CPU** → **MMU** → Physical Addresses

**Untranslated read or write**

- **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 essential to implementing 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

# Recall: General Address Translation



**Code**
**Data**
**Heap**
**Stack**

**Prog 1**
**Virtual**
**Address**
**Space 1**

**Translation Map 1**

**Data 2**
**Stack 1**
**Heap 1**
**Code 1**
**Stack 2**
**Data 1**
**Heap 2**
**Code 2**
**OS code**
**OS data**
**OS heap & Stacks**

**Code**
**Data**
**Heap**
**Stack**

**Prog 2**
**Virtual**
**Address**
**Space 2**

**Translation Map 2**

**Physical Address Space**

# Simple Base and Bounds (CRAY-1)



- **Could use base/limit 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

# Issues with Simple B&B Method



- **Fragmentation problem**
  - **Not every process is the same size**
  - **Over time, memory space becomes fragmented**
- **Missing support for sparse address space**
  - **Would like to have multiple chunks/program**
  - **E.g.: 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**

# More Flexible Segmentation



logical address

user view of memory space

physical memory space

- **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**

# Implementation of Multi-Segment Model

**Virtual Address**

| Seg # | Offset |
|-------|--------|

**offset** → **>** → **Error**

| | | |
|-------|--------|---|
| Base0 | Limit0 | V |
| Base1 | Limit1 | V |
| Base2 | Limit2 | V |
| Base3 | Limit3 | N |
| Base4 | Limit4 | V |
| Base5 | Limit5 | N |
| Base6 | Limit6 | N |
| Base7 | Limit7 | V |

**+** → **Physical Address**

**Check Valid** → **Access Error**

- **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

# Intel x86 Special Registers

## 80386 Special Registers

Segment registers

| | | |
|---|---|---|
| Code Seg. | | Data Seg. |
| 15 CS 0 | | 15 DS 0 |
| Stack Seg. | | Extra Seg. |
| 15 SS 0 | | 15 ES 0 |
| Extra Seg. | | Extra. Seg |
| 15 ES 0 | | 15 GS 0 |

| X | N T | IO PL | O F | D F | I F | T F | S F | Z F | x | A F | x | P F | x | C F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 15 | 14 | 13 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |

| P G | | E T | T S | T S | M P | P E | CR0 |
|---|---|---|---|---|---|---|---|
| 31 30 | | 5 | 4 | 3 | 2 | 1 0 | |

| Unused | CR1 |
|---|---|
| 31 | 0 Flags |

| Page Fault Linear Address | CR2 |
|---|---|
| 31 | 0 |

| Page Directory Base Register | Not Used | CR3 |
|---|---|---|
| 31 | 7 | 0 |

PG=Paging Enable
ET=Emulation Type
TS=Task Switched
EM=Emulate Coprocessor
MP=Math coprocessor present
PE=Protected Mode enable

X=Reserved
NT=Nested Task
IOPL=I/O Privilege Level
OF=Overflow Flag
DF=Direction Flag
IF=Interrupt Flag
TF=Trap Flag
SF=Sign Flag
ZF=Zero Flag
AF=Auxiliary Flag
PF=Parity Flag
CF=Carry Flag

| 15 | | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|
| Index | | | T I | RPL | |

RPL = Requestor Privilege Level
TI = Table Indicator
    (0 = GDT, 1 = LDT)
Index = Index into table

Protected Mode segment selector

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

# Example: Four Segments (16 bit addresses)

**Virtual Address Format**

| Seg | Offset |
|-----|--------|

15  14 13                                    0

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



Virtual Address Space

SegID = 0

SegID = 1

0x0000
0x4000
0x8000
0xC000

0x0000
0x4000
0x4800
0x5C00
0xF000

Might be shared

Space for Other Apps

Shared with Other Apps

Physical Address Space

# 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):**
1.  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
2.  Fetch 0x244. Translated to Physical=0x4244. Get "jal strlen"
    Move 0x0248 → $ra (return address!), Move 0x0360 → PC
3.  Fetch 0x360. Translated to Physical=0x4360. Get "li $v0,0"
    Move 0x0000 → $v0, Move PC+4→PC
4.  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

# 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")

# What if more segments than will fit into memory?



- **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

# 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**

# 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
      $1 \Rightarrow$ allocated, $0 \Rightarrow$ 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

# 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**

# Simple Page Table Example

**Example (4 byte pages)**



**0x00**
**a**
**b**
**c**
**d**
**0x04**
**e**
**f**
**0x06?**
**g**
**h**
**0x08**
**i**
**0x09?**
**j**
**k**
**l**

**0000 0000**
**0000 0100**
**0000 1000**

**Virtual Memory**

Page Table:
**0** **4**
**1** **3**
**2** **1**

**0001 0000**
**0000 1100**
**0000 0100**

**Page Table**

**0x00**
**0x04** i, j, k, l **0x05!**
**0x08**
**0x0C** e, f, g, h **0x0E!**
**0x10** a, b, c, d

**Physical Memory**

**0000 0110** ----> **0000 1110**
**0000 1001** ----> **0000 0101**

# What about Sharing?

**Virtual Address (Process A):**

| Virtual Page # | Offset |
|---|---|

**PageTablePtrA**

| page #0 | V,R |
|---|---|
| page #1 | V,R |
| page #2 | V,R,W |
| page #3 | V,R,W |
| page #4 | N |
| page #5 | V,R,W |

**PageTablePtrB**

| page #0 | V,R |
|---|---|
| page #1 | N |
| page #2 | V,R,W |
| page #3 | N |
| page #4 | V,R |
| page #5 | V,R,W |

**Virtual Address (Process B):**

| Virtual Page # | Offset |
|---|---|

**Shared Page**

This physical page appears in address space of both processes

# Memory Layout for Linux 32-bit



| | |
|---|---|
| **Kernel space** — User code CANNOT read from nor write to these addresses, doing so results in a Segmentation Fault | 0xc0000000 == TASK_SIZE |
| | Random stack offset |
| **Stack** (grows down) | RLIMIT_STACK (e.g., 8MB) |
| | Random mmap offset |
| **Memory Mapping Segment** — File mappings (including dynamic libraries) and anonymous mappings. Example: /lib/libc.so | |
| | program break / brk |
| **Heap** | start_brk |
| | Random brk offset |
| **BSS segment** — Uninitialized static variables, filled with zeros. Example: static char *userName; | |
| **Data segment** — Static variables initialized by the programmer. Example: static char *gonzo = "God's own prototype"; | end_data |
| **Text segment (ELF)** — Stores the binary image of the process (e.g., /bin/gonzo) | start_data / end_code / 0x08048000 / 0 |

1GB — Kernel space

3GB

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

# Summary: Simple Page Table

**Page Table**

**Virtual memory view**

**Physical memory view**

1111 1111
1111 0000 — stack

1100 0000

1000 0000 — heap

0100 0000 — data

0000 0000 — code

page # offset

| Page Table | |
|---|---|
| 11111 | 11101 |
| 11110 | 11100 |
| 11101 | null |
| 11100 | null |
| 11011 | null |
| 11010 | null |
| 11001 | null |
| 11000 | null |
| 10111 | null |
| 10110 | null |
| 10101 | null |
| 10100 | null |
| 10011 | null |
| 10010 | 10000 |
| 10001 | 01111 |
| 10000 | 01110 |
| 01111 | null |
| 01110 | null |
| 01101 | null |
| 01100 | null |
| 01011 | 01101 |
| 01010 | 01100 |
| 01001 | 01011 |
| 01000 | 01010 |
| 00111 | null |
| 00110 | null |
| 00101 | null |
| 00100 | null |
| 00011 | 00101 |
| 00010 | 00100 |
| 00001 | 00011 |
| 00000 | 00010 |

1110 1111

1110 0000 — stack

0111 000 — heap

0101 000 — data

0001 0000 — code

0000 0000

# Summary: Simple Page Table



**Page Table**

**Virtual memory view**

**Physical memory view**

1111 1111

stack

1110 0000

What happens if stack grows to 1110 0000?

heap

1000 0000

data

0100 0000

code

0000 0000

page # offset

stack 1110 0000

heap 0111 000

data 0101 000

code 0001 0000

0000 0000

| | |
|---|---|
| 11111 | 11101 |
| 11110 | 11100 |
| 11101 | null |
| 11100 | null |
| 11011 | null |
| 11010 | null |
| 11001 | null |
| 11000 | null |
| 10111 | null |
| 10110 | null |
| 10101 | null |
| 10100 | null |
| 10011 | null |
| 10010 | 10000 |
| 10001 | 01111 |
| 10000 | 01110 |
| 01111 | null |
| 01110 | null |
| 01101 | null |
| 01100 | null |
| 01011 | 01101 |
| 01010 | 01100 |
| 01001 | 01011 |
| 01000 | 01010 |
| 00111 | null |
| 00110 | null |
| 00101 | null |
| 00100 | null |
| 00011 | 00101 |
| 00010 | 00100 |
| 00001 | 00011 |
| 00000 | 00010 |

015

22

# Summary: Simple Page Table



Virtual memory view

Page Table

Physical memory view

Allocate new pages where room!

page # offset

# 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?**
  - Segments with pages inside them?
  - Need some sort of multi-level translation

# 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):

Virtual Address:

| Virtual Seg # | Virtual Page # | Offset |
|---|---|---|

| Base0 | Limit0 | V |
| Base1 | Limit1 | V |
| Base2 | Limit2 | V |
| Base3 | Limit3 | N |
| Base4 | Limit4 | V |
| Base5 | Limit5 | N |
| Base6 | Limit6 | N |
| Base7 | Limit7 | V |

| page #0 | V,R |
| page #1 | V,R |
| page #2 | V,R,W |
| page #3 | V,R,W |
| page #4 | N |
| page #5 | V,R,W |

| Physical Page # | Offset |
|---|---|

Physical Address

**Check Perm**

**>** → Access Error

Access Error

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

# What about Sharing (Complete Segment)?

**Process A**

| Virtual Seg # | Virtual Page # | Offset |
|---|---|---|

| Base0 | Limit0 | V |
|---|---|---|
| Base1 | Limit1 | V |
| Base2 | Limit2 | V |
| Base3 | Limit3 | N |
| Base4 | Limit4 | V |
| Base5 | Limit5 | N |
| Base6 | Limit6 | N |
| Base7 | Limit7 | V |

| page #0 | V,R |
|---|---|
| page #1 | V,R |
| page #2 | V,R,W |
| page #3 | V,R,W |
| page #4 | N |
| page #5 | V,R,W |

**Shared Segment**

| Base0 | Limit0 | V |
|---|---|---|
| Base1 | Limit1 | V |
| Base2 | Limit2 | V |
| Base3 | Limit3 | N |
| Base4 | Limit4 | V |
| Base5 | Limit5 | N |
| Base6 | Limit6 | N |
| Base7 | Limit7 | V |

**Process B**

| Virtual Seg # | Virtual Page # | Offset |
|---|---|---|

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

**Virtual Address:**

| 10 bits | 10 bits | 12 bits |
|---------|---------|---------|
| Virtual P1 index | Virtual P2 index | Offset |

**Physical Address:**

| Physical Page # | Offset |
|-----------------|--------|

**PageTablePtr**

→ 4 bytes ←

4KB

→ 4 bytes ←

**Tree of Page Tables**
**Tables fixed size (1024 entries)**
  – On context-switch: save single PageTablePtr register
**Valid bits on Page Table Entries**
  – Don't need every 2nd-level table
  – Even when exist, 2nd-level tables can reside on disk if not in use

# Summary: Two-Level Paging



**Virtual memory view**

*1111* *1111*
stack
*1111* *0000*

*1100* *0000*

heap
*1000* *0000*

data
*0100* *0000*

page2 #
code
*0000* *0000*

page1 # **offset**

**Page Table (level 1)**

*111*
*110* null
*101* null
*100*
*011* null
*010*
*001* null
*000*

**Page Tables (level 2)**

11 11101
10 11100
01 10111
00 10110

11 null
10 10000
01 01111
00 01110

11 01101
10 01100
01 01011
00 01010

11 00101
10 00100
01 00011
00 00010

**Physical memory view**

stack    1110 0000

stack

heap     0111 000

data     0101 000

code     0001 0000
         0000 0000

# Summary: Two-Level Paging



**Virtual memory view**

stack

1001 0000
(0x90)

heap

data

code

**Page Table (level 1)**

| | |
|---|---|
| 111 | ● |
| 110 | null |
| 101 | null |
| 100 | ● |
| 011 | null |
| 010 | ● |
| 001 | null |
| 000 | ● |

**Page Tables (level 2)**

| | |
|---|---|
| 11 | 11101 |
| 10 | 11100 |
| 01 | 10111 |
| 00 | 10110 |

| | |
|---|---|
| 11 | null |
| 10 | 10000 |
| 01 | 01111 |
| 00 | 01110 |

| | |
|---|---|
| 11 | 01101 |
| 10 | 01100 |
| 01 | 01011 |
| 00 | 01010 |

| | |
|---|---|
| 11 | 00101 |
| 10 | 00100 |
| 01 | 00011 |
| 00 | 00010 |

**Physical memory view**

stack — 1110 0000

stack

heap — 1000 0000 (0x80)

data

code

0001 0000

0000 0000

# Multi-level Translation Analysis

- Pros:
  - Only need to allocate as many page table entries as we need for application
    - » In other wards, sparse address spaces are easy
  - Easy memory allocation
  - Easy Sharing
    - » Share at segment or page level
- 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!

# 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