## CS162
## Operating Systems and
## Systems Programming
## Lecture 8

## CPU Scheduling, Protection Address Spaces

February 14, 2011

Ion Stoica

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

---

## Review: Last Time

- Scheduling: selecting a waiting process from the ready queue and allocating the CPU to it

- FCFS Scheduling:
  - Run threads to completion in order of submission
  - Pros: Simple (+)
  - Cons: Short jobs get stuck behind long ones (-)

- Round-Robin Scheduling:
  - Give each thread a small amount of CPU time when it executes; cycle between all ready threads
  - Pros: Better for short jobs (+)
  - Cons: Poor when jobs are same length (-)

---

## Goals for Today

- Finish discussion of Scheduling
- Kernel vs User Mode
- What is an Address Space?
- How is it Implemented?

**Note: Some slides and/or pictures in the following are
adapted from slides ©2005 Silberschatz, Galvin, and Gagne**

---

## Round-Robin Discussion

- How do you choose time slice?
  - What if too big?
    - » Response time suffers
  - What if infinite ($\infty$)?
    - » Get back FIFO
  - What if time slice too small?
    - » Throughput suffers!
- Actual choices of timeslice:
  - Initially, UNIX timeslice one second:
    - » Worked ok when UNIX was used by one or two people.
    - » What if three compilations going on? 3 seconds to echo each keystroke!
  - In practice, need to balance short-job performance and long-job throughput:
    - » Typical time slice today is between 10ms – 100ms
    - » Typical context-switching overhead is 0.1ms – 1ms
    - » Roughly 1% overhead due to context-switching

---

Page 1

## 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

- Completion Times:

| Job # | FIFO | RR |
|-------|------|------|
| 1 | 100 | 991 |
| 2 | 200 | 992 |
| ... | ... | ... |
| 9 | 900 | 999 |
| 10 | 1000 | 1000 |

  – Both RR and FCFS finish at the same time
  – Average response time is much worse under RR!
    » Bad when all jobs same length
- 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!

---

## Earlier Example with Different Time Quantum

Best FCFS:

| $P_2$ [8] | $P_4$ [24] | $P_1$ [53] | $P_3$ [68] |
|---|---|---|---|

0    8    32    85    153

| Quantum | $P_1$ | $P_2$ | $P_3$ | $P_4$ | Average |
|---------|-------|-------|-------|-------|---------|

| $P_1$ | $P_2$ | $P_3$ | $P_4$ | $P_1$ | $P_3$ | $P_4$ | $P_1$ | $P_3$ | $P_4$ | $P_1$ | $P_3$ | $P_1$ | $P_3$ | $P_1$ | $P_3$ | $P_1$ | $P_3$ | $P_3$ | $P_3$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

0  8  **16**  24  32  40  48  56  64  72  **80**  88  96  104  112  120  128  **133**  141  149  **153**

| | Quantum | $P_1$ | $P_2$ | $P_3$ | $P_4$ | Average |
|---|---|---|---|---|---|---|
| Wait Time | Q = 8 | 80 | 8 | 85 | 56 | 57¼ |
| | Q = 10 | 82 | 10 | 85 | 68 | 61¼ |
| | Q = 20 | 72 | 20 | 85 | 88 | 66¼ |
| | Worst FCFS | 68 | 145 | 0 | 121 | 83½ |
| Completion Time | Best FCFS | 85 | 8 | 153 | 32 | 69½ |
| | Q = 1 | 137 | 30 | 153 | 81 | 100½ |
| | Q = 5 | 135 | 28 | 153 | 82 | 99½ |
| | Q = 8 | 133 | 16 | 153 | 80 | 95½ |
| | Q = 10 | 135 | 18 | 153 | 92 | 99½ |
| | Q = 20 | 125 | 28 | 153 | 112 | 104½ |
| | Worst FCFS | 121 | 153 | 68 | 145 | 121¾ |

---

## What if we Knew the Future?

- Could we always mirror best FCFS?
- Shortest Job First (SJF):
  – Run whatever job has the least amount of computation to do

- Shortest Remaining Time First (SRTF):
  – Preemptive version of SJF: if job arrives and has a shorter time to completion than the remaining time on the current job, immediately preempt CPU

- These can be applied either to a whole program or the current CPU burst of each program
  – Idea is to get short jobs out of the system
  – Big effect on short jobs, only small effect on long ones
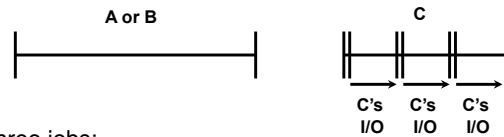  – Result is better average response time

---

## Discussion

- SJF/SRTF are the best you can do at minimizing average response time
  – Provably optimal (SJF among non-preemptive, SRTF among preemptive)
  – Since SRTF is always at least as good as SJF, focus on SRTF

- Comparison of SRTF with FCFS and RR
  – What if all jobs the same length?
    » SRTF becomes the same as FCFS (i.e. FCFS is best can do if all jobs the same length)
  – What if jobs have varying length?
    » SRTF (and RR): short jobs not stuck behind long ones

Page 2

## Example to illustrate benefits of SRTF

**A or B**

**C**

C's I/O     C's I/O     C's I/O

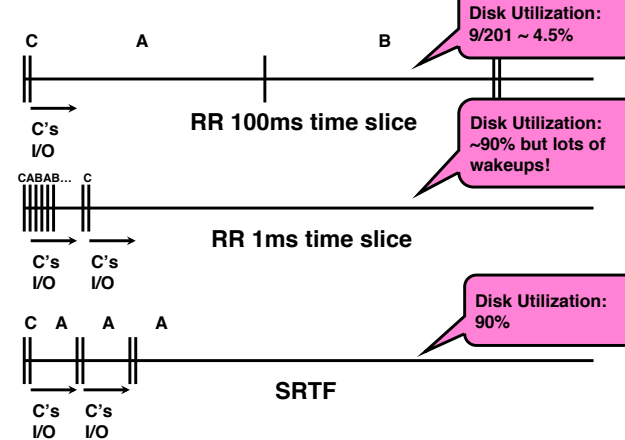- Three jobs:
  - A,B: CPU bound, each run for a week
    C: I/O bound, loop 1ms CPU, 9ms disk I/O
  - If only one at a time, C uses 90% of the disk, A or B could use 100% of the CPU
- With FIFO:
  - Once A or B get in, keep CPU for one week each
- What about RR or SRTF?
  - Easier to see with a timeline

---

## RR vs. SRTF

**C          A                              B**

**Disk Utilization: 9/201 ~ 4.5%**

C's I/O

**RR 100ms time slice**

**Disk Utilization: ~90% but lots of wakeups!**

CABAB... C

C's I/O     C's I/O

**RR 1ms time slice**

**Disk Utilization: 90%**

**C    A     A      A**

C's I/O     C's I/O

**SRTF**

---

## SRTF Further discussion

- 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) (+)
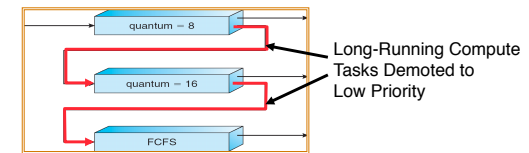  - Hard to predict future (-)
  - Unfair (-)

RESEARCH DEPT.

---

## Multi-Level Feedback Scheduling

quantum = 8

quantum = 16

FCFS

Long-Running Compute Tasks Demoted to Low Priority

- Multiple queues, each with different priority
  - Higher priority queues often considered "foreground" tasks

- Each queue has its own scheduling algorithm
  - e.g. foreground – RR, background – FCFS

- Adjust each job's priority as follows (details vary)
  - Job starts in highest priority queue
  - If it doesn't finish in its time quantum, drop one level
  - If it finishes, push up one level (or to top)

## Scheduling Fairness

- What about fairness?
  - Strict fixed-priority scheduling between queues is unfair (run highest, then next, etc):
    - » long running jobs may never get CPU
    - » In Multics, shut down machine, found 10-year-old job
  - Must give long-running jobs a fraction of the CPU even when there are shorter jobs to run
  - Tradeoff: fairness gained by hurting avg response time!

- How to implement fairness?
  - Could give each queue some fraction of the CPU
    - » What if one long-running job and 100 short-running ones?
    - » Like express lanes in a supermarket—sometimes express lanes get so long, get better service by going into one of the other lines
  - Could increase priority of jobs that don't get service
    - » What is done in UNIX
    - » This is ad hoc—what rate should you increase priorities?

## Lottery Scheduling

- Yet another alternative: Lottery Scheduling
  - Give each job some number of lottery tickets
  - On each time slice, randomly pick a winning ticket
  - On average, CPU time is proportional to number of tickets given to each job

- How to assign tickets?
  - To approximate SRTF, short running jobs get more, long running jobs get fewer
  - To avoid starvation, every job gets at least one ticket (everyone makes progress)

- Advantage over strict priority scheduling: behaves gracefully as load changes
  - Adding or deleting a job affects all jobs proportionally, independent of how many tickets each job possesses

## Lottery Scheduling Example

- Lottery Scheduling Example
  - Assume short jobs get 10 tickets, long jobs get 1 ticket

| # short jobs/ # long jobs | % of CPU each short jobs gets | % of CPU each long jobs gets |
|---|---|---|
| 1/1 | 91% | 9% |
| 0/2 | N/A | 50% |
| 2/0 | 50% | N/A |
| 10/1 | 9.9% | 0.99% |
| 1/10 | 50% | 5% |

  - What if too many short jobs to give reasonable response time?
    - » In UNIX, if load average is 100, hard to make progress
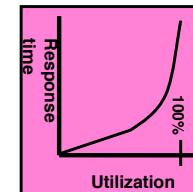    - » One approach: log some user out

## A Final Word On Scheduling

- When do the details of the scheduling policy and fairness really matter?
  - When there aren't enough resources to go around
- When should you simply buy a faster computer?
  - (Or network link, or expanded highway, or …)
  - One approach: Buy it when it will pay for itself in improved response time
    - » Might think that you should buy a faster X when X is utilized 100%, but usually, response time goes to infinity as utilization⇒100%



- An interesting implication of this curve:
  - Most scheduling algorithms work fine in the "linear" portion of the load curve, fail otherwise
  - Argues for buying a faster X when hit "knee" of curve

Page 4

## Administrivia

- Deadlines, project 1:
  - Design: Tomorrow, February 15th
  - Code: March 1st
  - Submitting instructions posted on Piazzza

---

## 5min Break

---

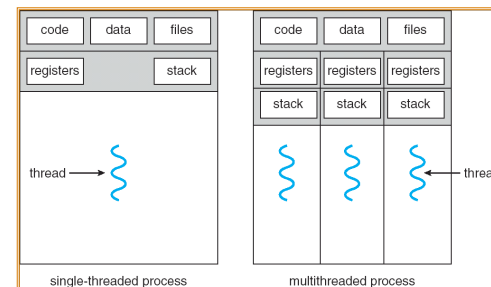## Virtualizing Resources



- Physical Reality:
  Different Processes/Threads share the same hardware
  - Need to multiplex CPU (Just finished: scheduling)
  - Need to multiplex use of Memory (Today)
  - Need to multiplex disk and devices (later in term)

- Why worry about memory sharing?
  - The complete working state of a process and/or kernel is defined by its data in memory (and registers)
  - Consequently, cannot just let different threads of control use the same memory
  - Probably don't want different threads to even have access to each other's memory (protection)

---

## Recall: Single and Multithreaded Processes



| code | data | files |
| registers | | stack |

single-threaded process            multithreaded process

- Threads encapsulate execution (concurrency)
  - "Active" component of a process
- Address spaces encapsulate protection
  - Keeps buggy program from trashing the system
  - "Passive" component of a process

Page 5

## 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
    » Programs protected from themselves
- 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

## Binding of Instructions and Data to Memory

- Binding of instructions and data to add
  - Choose addresses for instructions and standpoint of the processor

```
Assume 4byte words
0x300 = 4 * 0x0C0
0x0C0 = 0000 1100 0000
0x300 = 0011 0000 0000
```

```
data1:  dw    32                      0x300        000020
        …                             …            …
start:  lw    r1,0(data1)             0x900        8C2000C0
        jal   checkit                 0x904        0C000340
loop:   addi  r1, r1, -1              0x908        2021FFFF
        bnz   r1, r0, loop            0x90C        1420FFFF
        …                             …
checkit: …                            0xD00        …
```

  - Could we place `data1`, `start`, and/or `checkit` at different addresses?
    » Yes
    » When? Compile time/Load time/Execution time
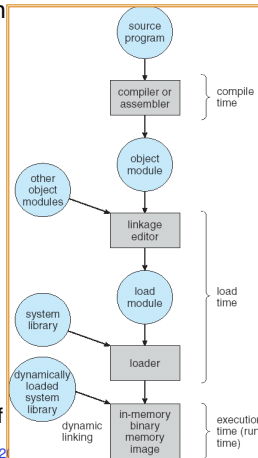  - Related: which physical memory locations hold particular instructions or data?

## Multi-step Processing of a Program for Execution

- Preparation of a program for execution involves components at:
  - Compile time (i.e., "gcc")
  - Link/Load time (unix "ld" does link)
  - Execution time (e.g. dynamic libs)

- Addresses can be bound to final values anywhere in this path
  - Depends on hardware support
  - Also depends on operating system

- Dynamic Libraries
  - Linking postponed until execution
  - Small piece of code, *stub*, used to locate appropriate memory-resident library routine
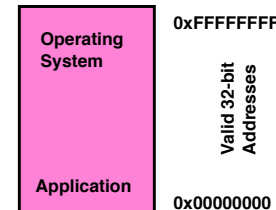  - Stub replaces itself with the address of the routine, and executes routine

## Uniprogramming

- 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
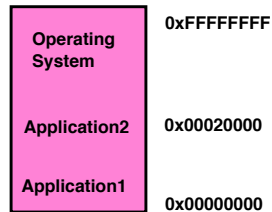- Of course, this doesn't help us with multithreading

Page 6

## Multiprogramming (First Version)

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

| | |
|---|---|
| **Operating System** | 0xFFFFFFFF |
| **Application2** | 0x00020000 |
| **Application1** | 0x00000000 |

  - 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

## Multiprogramming (Version with Protection)

- Can we protect programs from each other without translation?

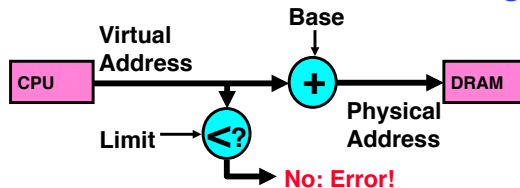| | |
|---|---|
| **Operating System** | 0xFFFFFFFF |
| | LimitAddr=0x10000 |
| **Application2** | 0x00020000 ← BaseAddr=0x20000 |
| **Application1** | 0x00000000 |

  - 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
    » User not allowed to change base/limit registers

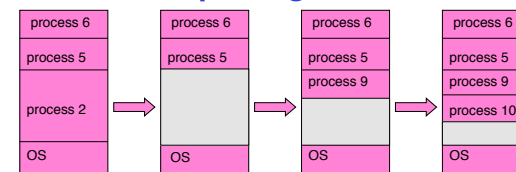## Segmentation with Base and Limit Registers



**No: Error!**

- Could use base/limit for dynamic address translation (often called "segmentation"):
  - Alter address of every load/store by adding "base"
  - User allowed to read/write within segment
    » Accesses are relative to segment so don't have to be relocated when program moved to different segment
  - User may have multiple segments available (e.g x86)
    » Loads and stores include segment ID in opcode:
      x86 Example: `mov [es:bx],ax.`
    » Operating system moves around segment base pointers as necessary

## 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 by providing multiple segments per process
- Need enough physical memory for every process

Page 7

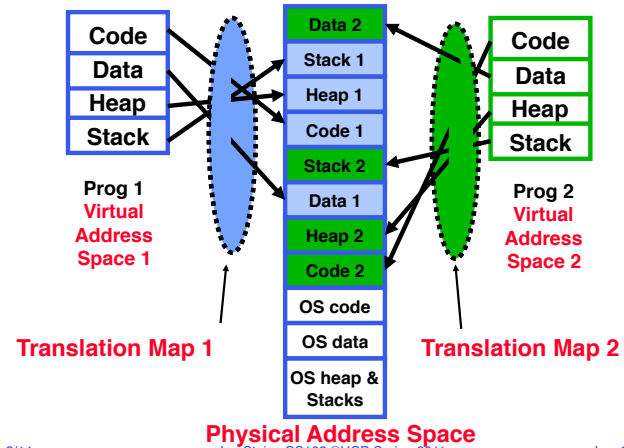## Multiprogramming (Translation and Protection v2)

- Problem: Run multiple applications in such a way that they are protected from one another
- Goals:
  - Isolate processes and kernel from one another
  - Allow flexible translation that:
    - » Doesn't lead to fragmentation
    - » Allows easy sharing between processes
    - » Allows only part of process to be resident in physical memory
- (Some of the required) Hardware Mechanisms:
  - General Address Translation
    - » Flexible: Can fit physical chunks of memory into arbitrary places in users address space
    - » Not limited to small number of segments
    - » Think of this as providing a large number (thousands) of fixed-sized segments (called "pages")
  - Dual Mode Operation
    - » Protection base involving kernel/user distinction

## Example of General Address Translation



Prog 1
**Virtual Address Space 1**

Prog 2
**Virtual Address Space 2**

**Translation Map 1**     **Translation Map 2**

**Physical Address Space**

## Two Views of Memory



- Recall: 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 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
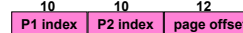- With translation, every program can be linked/loaded into same region of user address space

## Example of Translation Table Format

**Two-level Page Tables
32-bit address:**

| 10 | 10 | 12 |
|---|---|---|
| P1 index | P2 index | page offset |



- Page: a unit of memory translatable by memory management unit (MMU)
  - Typically 1K – 8K
- Page table structure in memory
  - Each user has different page table
- Address Space switch: change pointer to base of table (hardware register)
  - Hardware traverses page table (for many architectures)
  - MIPS uses software to traverse table

Page 8

## Dual-Mode Operation

- Can Application Modify its own translation tables?
  - If it could, could get access to all of physical memory
  - Has to be restricted somehow

- To Assist with Protection, hardware provides at least two modes (Dual-Mode Operation):
  - "Kernel" mode (or "supervisor" or "protected")
  - "User" mode (Normal program mode)
  - Mode set with bits in special control register only accessible in kernel-mode

- Intel processor actually has four "rings" of protection:
  - PL (Priviledge Level) from 0 – 3
    » PL0 has full access, PL3 has least
  - Typical OS kernels on Intel processors only use PL0 ("user") and PL3 ("kernel")

## For Protection, Lock User-Programs in Asylum

- Idea: Lock user programs in padded cell with no exit or sharp objects
  - Cannot change mode to kernel mode
  - User cannot modify page table mapping
  - Limited access to memory: cannot adversely effect other processes
    » Side-effect: Limited access to memory-mapped I/O operations
  - What else needs to be protected?

- A couple of issues
  - How to share CPU between kernel and user programs?
    » Kinda like both the inmates and the warden in asylum are the same person. How do you manage this???
  - How does one switch between kernel and user modes?
    » OS → user (kernel → user mode): getting into cell
    » User→ OS (user → kernel mode): getting out of cell

## How to get from Kernel→User

- What does the kernel do to create a new user process?
  - Allocate and initialize process control block
  - Read program off disk and store in memory
  - Allocate and initialize translation table
    » Point at code in memory so program can execute
    » Possibly point at statically initialized data
  - Run Program:
    » Set machine registers
    » Set hardware pointer to translation table
    » Set processor status word for user mode
    » Jump to start of program
- How does kernel switch between processes?
  - Same saving/restoring of registers as before
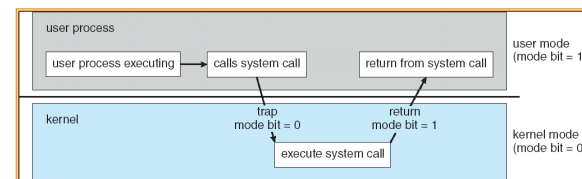  - Save/restore hardware pointer to translation table

## User→Kernel (System Call)

- Can't let inmate (user) get out of padded cell on own
  - Would defeat purpose of protection!
  - So, how does the user program get back into kernel?



- System call: Voluntary procedure call into kernel
  - Hardware for controlled User→Kernel transition
  - Can any kernel routine be called?
    » No! Only specific ones.
  - System call ID encoded into system call instruction
    » Index forces well-defined interface with kernel

## Summary (1)

- Shortest Job First (SJF)/Shortest Remaining Time First (SRTF):
  - Run whatever job has the least amount of computation to do/least remaining amount of computation to do
  - Pros: Optimal (average response time)
  - Cons: Hard to predict future, Unfair
- Multi-Level Feedback Scheduling:
  - Multiple queues of different priorities
  - Automatic promotion/demotion of process priority in order to approximate SJF/SRTF
- Lottery Scheduling:
  - Give each thread a priority-dependent number of tokens (short tasks⇒more tokens)
  - Reserve a minimum number of tokens for every thread to ensure forward progress/fairness
- Evaluation of mechanisms:
  - Analytical, Queuing Theory, Simulation

## Summary (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
- Simple Protection through segmentation
  - Base+limit registers restrict memory accessible to user
  - Can be used to translate as well
- Full translation of addresses through Memory Management Unit (MMU)
  - Every Access translated through page table
  - Changing of page tables only available to user
- Dual-Mode
  - Kernel/User distinction: User restricted
  - User→Kernel: System calls, Traps, or Interrupts