

## TCCS 422: OPERATING SYSTEMS

### Limited Direct Execution, Introduction to CPU Scheduling



Wes J. Lloyd  
School of Engineering and Technology  
University of Washington - Tacoma

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

Tacoma

## FEEDBACK 1/14

- `fork()` – will a modification to the data of a forked process trigger Copy-and-Write to duplicate all process data?, or only that which was modified... ?
  - Only modified data – saves memory and time
  - How can this be tested?
- Example: `exec.c` revisited - example using `wait()`
  - Recompiled `exec.c`
  - Unable to reproduce bug (parent existing before child)
  - Old version was compiled with previous Linux kernel, suspect changes may have caused different behavior
  - Check out: `waitpid()` and capturing return code from `wait()`

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.2

## FEEDBACK - 2

- Why **`exec`**, instead of just calling another function like in Java?
  - `Fork()`, `exec()`, and `wait()` allow ***fine grained*** control over processes
  - Goal: remotely invoke another executable program, from an existing C program
  - Can `fork()` a new processes to delegate to run an external program
  - Can redirect input, output, stderr of processes
- `Exec..` is still a little complex
  - 6 versions: `execl()`, `execclp()`, `execle()`, `execv()`, `execvp()`, `execvpe()`
  - "p" – duplicates the PATH of the shell
  - "e" – allows the environment (variables) to be passed in as an array of pointers to NULL-terminated strings, list is NULL terminated
  - "v" – provide cmd & args as array (***vector***)
  - "l" – provide cmd & args as (***list***) of individual params sent to the function (***see execl.c example***)

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.3

## exec()

- Supports running an external program
- 6 types: `execl()`, `execclp()`, `execle()`, `execv()`, `execvp()`, `execvpe()`
- `execl()`, `execclp()`, `execle()`: `const char *arg` (***example: execl.c***)  
Provide cmd and args as individual params to the function  
Each arg is a pointer to a null-terminated string  
**ODD**: pass a variable number of args: (`arg0`, `arg1`, .. `argn`)
- `execv()`, `execvp()`, `execvpe()` (***example: exec.c***)  
Provide cmd and args as an Array of pointers to strings  
Strings are null-terminated  
First argument is name of command being executed  
Fixed number of args passed in

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.4

## FEEDBACK - 3

- What is the (too much / too little) control tradeoff for Operating System designs?
  - Who (what) has too much/little control?
  - What are consequences of too much control?
  - ... consequences of too little control?
- What is direct vs. limited direct execution with respect to operating systems?
  - ***Direct*** – code runs directly on HW, OS can not intervene, processes can "run away"...
  - ***Limited direct*** – code runs on HW, but can not directly access protected resources (memory, devices), OS can intervene and \*preempt\* a process



January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.5

## FEEDBACK - 4

- How do we identify system processes?
- All kernel processes (threads) are children of process (`pid==2`)
- ***"ps" command to display processes owned by the kernel:***  
`ps --ppid 2 -p 2 -o`  
`uname,pid,ppid,pcpu,pmem,vsz,rssize,start,time,cmd |`  
`more`
- ***"ps" command to display non-kernel processes:***  
`ps --ppid 2 -p 2 -o`  
`uname,pid,ppid,pcpu,pmem,vsz,rssize,start,time,cmd -`  
`-deselect | more`
- What is significant regarding the output?

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.6

FEEDBACK – 5

- Hard time following along. Overwhelming amount of information. Can't seem to find guideline/direction
  - AFTER CLASS - :
    - (1) Review Slides
    - (2) Read chapters (in this case Ch. 5 & 6)
    - (3) While reviewing/reading, **make list** of confusing topics
    - (4) **Formulate questions** about these topics
  - THEN
    - (1) Visit professor after class, during office hours, or make an appointment --or--
    - (2) Email list of questions to professor

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.7

OBJECTIVES

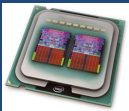
- Active Reading Quiz 1
- Assignment 0 / Linux Tutorial
- C Tutorial
- Chapter 6 – Limited Direct Execution
- **CPU Scheduling:**
- Chapter 7 – Scheduling Introduction
- Chapter 8 – Multi-level Feedback Queue (MLFQ)

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.8

CH. 6:  
LIMITED DIRECT  
EXECUTION



January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.9

VIRTUALIZING THE CPU

- How does the CPU support running so many jobs simultaneously?
- **Time Sharing**
- Tradeoffs:
  - Performance
    - Excessive overhead
  - Control
    - Fairness
    - Security
- Both HW and OS support is used

```
top - 09:21:07 up 45 days, 1:21, 1 user, load average: 0.00, 0.00, 0.00
task  PID   PPID  PSR   CPU%  MEM%  VSZ    RSS   TID    D    NAME
1      1      1   S     0.0    0.0    0K     0K    1    sshd: /usr/sbin/sshd -D
2      2      1   S     0.0    0.0    0K     0K    2    /usr/sbin/sshd -D
3      3      1   S     0.0    0.0    0K     0K    3    /usr/sbin/sshd -D
4      4      1   S     0.0    0.0    0K     0K    4    /usr/sbin/sshd -D
5      5      1   S     0.0    0.0    0K     0K    5    /usr/sbin/sshd -D
6      6      1   S     0.0    0.0    0K     0K    6    /usr/sbin/sshd -D
7      7      1   S     0.0    0.0    0K     0K    7    /usr/sbin/sshd -D
8      8      1   S     0.0    0.0    0K     0K    8    /usr/sbin/sshd -D
9      9      1   S     0.0    0.0    0K     0K    9    /usr/sbin/sshd -D
10     10     1   S     0.0    0.0    0K     0K   10    /usr/sbin/sshd -D
11     11     1   S     0.0    0.0    0K     0K   11    /usr/sbin/sshd -D
12     12     1   S     0.0    0.0    0K     0K   12    /usr/sbin/sshd -D
13     13     1   S     0.0    0.0    0K     0K   13    /usr/sbin/sshd -D
14     14     1   S     0.0    0.0    0K     0K   14    /usr/sbin/sshd -D
15     15     1   S     0.0    0.0    0K     0K   15    /usr/sbin/sshd -D
16     16     1   S     0.0    0.0    0K     0K   16    /usr/sbin/sshd -D
17     17     1   S     0.0    0.0    0K     0K   17    /usr/sbin/sshd -D
18     18     1   S     0.0    0.0    0K     0K   18    /usr/sbin/sshd -D
19     19     1   S     0.0    0.0    0K     0K   19    /usr/sbin/sshd -D
20     20     1   S     0.0    0.0    0K     0K   20    /usr/sbin/sshd -D
21     21     1   S     0.0    0.0    0K     0K   21    /usr/sbin/sshd -D
22     22     1   S     0.0    0.0    0K     0K   22    /usr/sbin/sshd -D
23     23     1   S     0.0    0.0    0K     0K   23    /usr/sbin/sshd -D
24     24     1   S     0.0    0.0    0K     0K   24    /usr/sbin/sshd -D
25     25     1   S     0.0    0.0    0K     0K   25    /usr/sbin/sshd -D
26     26     1   S     0.0    0.0    0K     0K   26    /usr/sbin/sshd -D
27     27     1   S     0.0    0.0    0K     0K   27    /usr/sbin/sshd -D
28     28     1   S     0.0    0.0    0K     0K   28    /usr/sbin/sshd -D
29     29     1   S     0.0    0.0    0K     0K   29    /usr/sbin/sshd -D
30     30     1   S     0.0    0.0    0K     0K   30    /usr/sbin/sshd -D
31     31     1   S     0.0    0.0    0K     0K   31    /usr/sbin/sshd -D
32     32     1   S     0.0    0.0    0K     0K   32    /usr/sbin/sshd -D
33     33     1   S     0.0    0.0    0K     0K   33    /usr/sbin/sshd -D
34     34     1   S     0.0    0.0    0K     0K   34    /usr/sbin/sshd -D
35     35     1   S     0.0    0.0    0K     0K   35    /usr/sbin/sshd -D
36     36     1   S     0.0    0.0    0K     0K   36    /usr/sbin/sshd -D
37     37     1   S     0.0    0.0    0K     0K   37    /usr/sbin/sshd -D
38     38     1   S     0.0    0.0    0K     0K   38    /usr/sbin/sshd -D
39     39     1   S     0.0    0.0    0K     0K   39    /usr/sbin/sshd -D
40     40     1   S     0.0    0.0    0K     0K   40    /usr/sbin/sshd -D
41     41     1   S     0.0    0.0    0K     0K   41    /usr/sbin/sshd -D
42     42     1   S     0.0    0.0    0K     0K   42    /usr/sbin/sshd -D
43     43     1   S     0.0    0.0    0K     0K   43    /usr/sbin/sshd -D
44     44     1   S     0.0    0.0    0K     0K   44    /usr/sbin/sshd -D
45     45     1   S     0.0    0.0    0K     0K   45    /usr/sbin/sshd -D
46     46     1   S     0.0    0.0    0K     0K   46    /usr/sbin/sshd -D
47     47     1   S     0.0    0.0    0K     0K   47    /usr/sbin/sshd -D
48     48     1   S     0.0    0.0    0K     0K   48    /usr/sbin/sshd -D
49     49     1   S     0.0    0.0    0K     0K   49    /usr/sbin/sshd -D
50     50     1   S     0.0    0.0    0K     0K   50    /usr/sbin/sshd -D
51     51     1   S     0.0    0.0    0K     0K   51    /usr/sbin/sshd -D
52     52     1   S     0.0    0.0    0K     0K   52    /usr/sbin/sshd -D
53     53     1   S     0.0    0.0    0K     0K   53    /usr/sbin/sshd -D
54     54     1   S     0.0    0.0    0K     0K   54    /usr/sbin/sshd -D
55     55     1   S     0.0    0.0    0K     0K   55    /usr/sbin/sshd -D
56     56     1   S     0.0    0.0    0K     0K   56    /usr/sbin/sshd -D
57     57     1   S     0.0    0.0    0K     0K   57    /usr/sbin/sshd -D
58     58     1   S     0.0    0.0    0K     0K   58    /usr/sbin/sshd -D
59     59     1   S     0.0    0.0    0K     0K   59    /usr/sbin/sshd -D
60     60     1   S     0.0    0.0    0K     0K   60    /usr/sbin/sshd -D
61     61     1   S     0.0    0.0    0K     0K   61    /usr/sbin/sshd -D
62     62     1   S     0.0    0.0    0K     0K   62    /usr/sbin/sshd -D
63     63     1   S     0.0    0.0    0K     0K   63    /usr/sbin/sshd -D
64     64     1   S     0.0    0.0    0K     0K   64    /usr/sbin/sshd -D
65     65     1   S     0.0    0.0    0K     0K   65    /usr/sbin/sshd -D
66     66     1   S     0.0    0.0    0K     0K   66    /usr/sbin/sshd -D
67     67     1   S     0.0    0.0    0K     0K   67    /usr/sbin/sshd -D
68     68     1   S     0.0    0.0    0K     0K   68    /usr/sbin/sshd -D
69     69     1   S     0.0    0.0    0K     0K   69    /usr/sbin/sshd -D
70     70     1   S     0.0    0.0    0K     0K   70    /usr/sbin/sshd -D
71     71     1   S     0.0    0.0    0K     0K   71    /usr/sbin/sshd -D
72     72     1   S     0.0    0.0    0K     0K   72    /usr/sbin/sshd -D
73     73     1   S     0.0    0.0    0K     0K   73    /usr/sbin/sshd -D
74     74     1   S     0.0    0.0    0K     0K   74    /usr/sbin/sshd -D
75     75     1   S     0.0    0.0    0K     0K   75    /usr/sbin/sshd -D
76     76     1   S     0.0    0.0    0K     0K   76    /usr/sbin/sshd -D
77     77     1   S     0.0    0.0    0K     0K   77    /usr/sbin/sshd -D
78     78     1   S     0.0    0.0    0K     0K   78    /usr/sbin/sshd -D
79     79     1   S     0.0    0.0    0K     0K   79    /usr/sbin/sshd -D
80     80     1   S     0.0    0.0    0K     0K   80    /usr/sbin/sshd -D
81     81     1   S     0.0    0.0    0K     0K   81    /usr/sbin/sshd -D
82     82     1   S     0.0    0.0    0K     0K   82    /usr/sbin/sshd -D
83     83     1   S     0.0    0.0    0K     0K   83    /usr/sbin/sshd -D
84     84     1   S     0.0    0.0    0K     0K   84    /usr/sbin/sshd -D
85     85     1   S     0.0    0.0    0K     0K   85    /usr/sbin/sshd -D
86     86     1   S     0.0    0.0    0K     0K   86    /usr/sbin/sshd -D
87     87     1   S     0.0    0.0    0K     0K   87    /usr/sbin/sshd -D
88     88     1   S     0.0    0.0    0K     0K   88    /usr/sbin/sshd -D
89     89     1   S     0.0    0.0    0K     0K   89    /usr/sbin/sshd -D
90     90     1   S     0.0    0.0    0K     0K   90    /usr/sbin/sshd -D
91     91     1   S     0.0    0.0    0K     0K   91    /usr/sbin/sshd -D
92     92     1   S     0.0    0.0    0K     0K   92    /usr/sbin/sshd -D
93     93     1   S     0.0    0.0    0K     0K   93    /usr/sbin/sshd -D
94     94     1   S     0.0    0.0    0K     0K   94    /usr/sbin/sshd -D
95     95     1   S     0.0    0.0    0K     0K   95    /usr/sbin/sshd -D
96     96     1   S     0.0    0.0    0K     0K   96    /usr/sbin/sshd -D
97     97     1   S     0.0    0.0    0K     0K   97    /usr/sbin/sshd -D
98     98     1   S     0.0    0.0    0K     0K   98    /usr/sbin/sshd -D
99     99     1   S     0.0    0.0    0K     0K   99    /usr/sbin/sshd -D
100    100     1   S     0.0    0.0    0K     0K  100    /usr/sbin/sshd -D
```

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.10

LIMITED DIRECT EXECUTION

- OS implements LDE to support time/resource sharing
- Limited direct execution means “only limited” processes can execute DIRECTLY on the CPU in **trusted** mode
- TRUSTED means the process is trusted, and it can do anything... (e.g. it is a system / kernel level process)
- Enabled by **protected (safe) control transfer**
- CPU supported context switch
- Provides data isolation

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.11

CPU MODES

- Utilize CPU Privilege Rings (Intel x86)
  - rings 0 (kernel), 1 (VM kernel), 2 (unused), 3 (user)

access ← no access

- **User mode:**  
Application is running, but w/o direct I/O access
- **Kernel mode:**  
OS kernel is running performing restricted operations

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.12

CPU MODES

- **User mode: ring 3 - untrusted**
  - Some instructions and registers are disabled by the CPU
  - Exception registers
  - HALT instruction
  - MMU instructions
  - OS memory access
  - I/O device access
- **Kernel mode: ring 0 – trusted**
  - All instructions and registers enabled

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.13

SYSTEM CALLS

- Implement restricted “OS” operations
- Kernel exposes key functions through an API:
  - Device I/O (e.g. file I/O)
  - Task swapping: context switching between processes
  - Memory management/allocation: malloc()
  - Creating/destroying processes

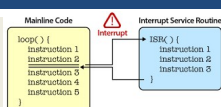
January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.14

TRAPS:  
SYSTEM CALLS, EXCEPTIONS, INTERRUPTS

- Trap: any transfer to kernel mode
- Three kinds of traps
  - **System call:** (planned) user → kernel
    - SYSCALL for I/O, etc.
  - **Exception:** (error) user → kernel
    - Div by zero, page fault, page protection error
  - **Interrupt:** (event) user → kernel
    - Non-maskable vs. maskable
    - Keyboard event, network packet arrival, timer ticks
    - Memory parity error (ECC), hard drive failure



January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.15

EXCEPTION TYPES

Exception type	Synchronous vs. asynchronous	User request vs. forced	User maskable vs. nonmaskable	Within vs. between instructions	Resume vs. terminate
I/O device request	Asynchronous	Coerced	Nonmaskable	Between	Resume
Invoke operating system	Synchronous	User request	Nonmaskable	Between	Resume
Trailing instruction execution	Synchronous	User request	User maskable	Between	Resume
Breakpoint	Synchronous	User request	User maskable	Between	Resume
Integer arithmetic overflow	Synchronous	Coerced	User maskable	Within	Resume
Floating-point arithmetic overflow or underflow	Synchronous	Coerced	User maskable	Within	Resume
Page fault	Synchronous	Coerced	Nonmaskable	Within	Resume
Unauthorized memory accesses	Synchronous	Coerced	User maskable	Within	Resume
Memory protection violation	Synchronous	Coerced	Nonmaskable	Within	Resume
Using undefined instruction	Synchronous	Coerced	Nonmaskable	Within	Terminate
Hardware malfunction	Asynchronous	Coerced	Nonmaskable	Within	Terminate
Power failure	Asynchronous	Coerced	Nonmaskable	Within	Terminate

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.16

OS @ boot (kernel mode)

Hardware

Program (user mode)

- Initialize trap table
  - remember address of ... syscall handler
- OS @ run (kernel mode)
  - Create entry for process list
  - Allocate memory for program
  - Load program into memory
  - Setup user stack with argv
  - Fill kernel stack with reg/PC
  - return-from-trap
    - restore regs from kernel stack
    - move to user mode
    - jump to main
    - Run main()
    - Call system trap into OS
- Handle trap
  - Do work of syscall
  - return-from-trap
    - save regs to kernel stack
    - move to kernel mode
    - jump to trap handler
    - restore regs from kernel stack
    - move to user mode
    - jump to PC after trap
    - return from main trap (via exit())
- Free memory of process
- Remove from process list

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.17

OS @ boot (kernel mode)

Hardware

Program (user mode)

- Initialize trap table
  - remember address of ... syscall handler
- OS @ run (kernel mode)
  - Create entry for process list
  - Allocate memory for program
  - Load program into memory
  - Setup user stack with argv
  - Fill kernel stack with reg/PC

Computer BOOT Sequence:  
OS with Limited Direct Execution

- move to kernel mode
- jump to trap handler
- Handle trap
  - Do work of syscall
  - return-from-trap
    - restore regs from kernel stack
    - move to user mode
    - jump to PC after trap
    - return from main trap (via exit())
- Free memory of process
- Remove from process list

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.18

MULTITASKING

- How/when should the OS regain control of the CPU to switch between processes?
- Cooperative multitasking (mostly pre 32-bit)
  - < Windows 95, Mac OSX
  - Opportunistic: running programs must give up control
    - User programs must call a special **yield** system call
    - When performing I/O
    - Illegal operations
  - (POLLEV)

What problems could you for see with this approach?

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.19

MULTITASKING

- How/when should the OS regain control of the CPU to switch between processes?
- Cooperative multitasking (mostly pre 32-bit)
  - < Windows 95, Mac OSX
  - Opportunistic: running programs must give up control
    - User programs must call a special **yield** system call
    - When performing I/O
    - Illegal operations
  - (POLLEV)

What problems could you for see with this approach?

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.20

A process gets stuck in an infinite loop.

→ Reboot the machine

What problems exist for regaining the control of the CPU with cooperative multitasking OSes?

Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

Total Recall

QUESTION: MULTITASKING

- What problems exist for regaining the control of the CPU with cooperative multitasking OSes?

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.22

MULTITASKING - 2

- Preemptive multitasking (32 & 64 bit OSes)
- >= Mac OSX, Windows 95+
- Timer interrupt
  - Raised at some regular interval (in ms)
  - Interrupt handling
    - Current program is halted
    - Program states are saved
    - OS Interrupt handler is run (kernel mode)
- (PollEV) What is a good interval for the timer interrupt?

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.23

MULTITASKING - 2

- Preemptive multitasking (32 & 64 bit OSes)
- >= Mac OSX, Windows 95+
- Timer interrupt
  - Raised at some regular interval (in ms)
  - Interrupt handling
    - Current program is halted
    - Program states are saved
    - OS Interrupt handler is run (kernel mode)
- (PollEV) What is a good interval for the timer interrupt?

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.24

A timer interrupt gives OS the ability to run again on a CPU.

W

For an OS that uses a system timer to force arbitrary context switches to share the CPU, what is a good value (in seconds) for the timer interrupt?

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

Total Rec: L4.2

QUESTION: TIME SLICE

■ For an OS that uses a system timer to force arbitrary context switches to share the CPU, what is a good value (in seconds) for the timer interrupt?

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.26

CONTEXT SWITCH

■ Preemptive multitasking initiates “trap” into the OS code to determine:

• Whether to continue running the **current process**, or switch to a **different one**.

• If the decision is made to switch, the OS performs a context switch swapping out the current process for a new one.

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.27

CONTEXT SWITCH - 2

1. Save register values of the current process to its kernel stack

■ General purpose registers

■ PC: program counter (instruction pointer)

■ kernel stack pointer

2. Restore soon-to-be-executing process from its kernel stack

3. Switch to the kernel stack for the soon-to-be-executing process

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.28

OS @ boot (kernel mode)

Hardware

Program (user mode)

initialize trap table

remember address of ...  
syscall handler  
timer handler

start interrupt timer

start timer  
interrupt CPU in X ms

OS @ run (kernel mode)

Hardware

Process A

timer interrupt

save regs(A) to k-stack(A)  
move to kernel mode  
jump to trap handler

Handle the trap

Call switch() routine

save regs(A) to proc-struct(A)  
restore regs(B) from proc-struct(B)  
switch to k-stack(B)  
return-from-trap (into B)

restore regs(B) from k-stack(B)  
move to user mode  
jump to B's PC

Process B

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.29

OS @ boot (kernel mode)

Hardware

Program (user mode)

initialize trap table

remember address of ...  
syscall handler  
timer handler

start interrupt timer

start timer  
interrupt CPU in X ms

OS @ run (kernel mode)

Hardware

Process A

timer interrupt

save regs(A) to k-stack(A)  
move to kernel mode  
jump to trap handler

Handle the trap

Call switch() routine

save regs(A) to proc-struct(A)  
restore regs(B) from proc-struct(B)  
switch to k-stack(B)  
return-from-trap (into B)

restore regs(B) from k-stack(B)  
move to user mode  
jump to B's PC

Process B

January 16, 2019

TCSS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.30

## INTERRUPTED INTERRUPTS

- What happens if during an interrupt (trap to kernel mode), another interrupt occurs?
- Linux
  - < 2.6 kernel: non-preemptive kernel
  - >= 2.6 kernel: preemptive kernel

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.31

## PREEMPTIVE KERNEL

- Use "locks" as markers of regions of non-preemptibility (non-maskable interrupt)
- Preemption counter (`preempt_count`)
  - begins at zero
  - increments for each lock acquired (not safe to preempt)
  - decrements when locks are released
- Interrupt can be interrupted when `preempt_count=0`
  - It is safe to preempt (maskable interrupt)
  - the interrupt is more important

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.32

## CHAPTER 7- SCHEDULING: INTRODUCTION



January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.33

## SCHEDULING METRICS

- Metrics:** A standard measure to quantify to what degree a system possesses some property. Metrics provide *repeatable* techniques to quantify and compare systems.
- Measurements** are the numbers derived from the application of metrics
- Scheduling Metric #1: **Turnaround time**
- The time at which the job completes minus the time at which the job arrived in the system

$$T_{\text{turnaround}} = T_{\text{completion}} - T_{\text{arrival}}$$

- How is turnaround time different than execution time?

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.34

## SCHEDULING METRICS - 2

- Scheduling Metric #2: **Fairness**
  - Jain's fairness index
  - Quantifies if jobs receive a fair share of system resources

$$\mathcal{J}(x_1, x_2, \dots, x_n) = \frac{(\sum_{i=1}^n x_i)^2}{n \cdot \sum_{i=1}^n x_i^2}$$

- n processes
- $x_i$  is time share of each process
- worst case =  $1/n$
- best case = 1

- Consider  $n=3$ , worst case = .333, best case=1
- With  $n=3$  and  $x_1=.2, x_2=.7, x_3=.1$ , fairness=.62
- With  $n=3$  and  $x_1=.33, x_2=.33, x_3=.33$ , fairness=1

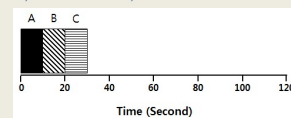
January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.35

## SCHEDULERS

- FIFO: first in, first out
  - Very simple, easy to implement
- Consider
  - 3 x 10sec jobs, arrival: A B C, duration 10 sec each



$$\text{Average turnaround time} = \frac{10 + 20 + 30}{3} = 20 \text{ sec}$$

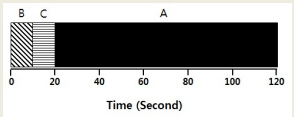
January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.36

SJF: SHORTEST JOB FIRST

- Given that we know execution times in advance:
  - Run in order of duration, shortest to longest
  - Non preemptive scheduler
  - This is not realistic
  - Arrival: A B C, duration a=100 sec, b/c=10sec



Average turnaround time =  $\frac{10 + 20 + 120}{3} = 50 \text{ sec}$

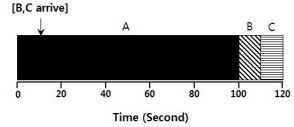
January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.37

SJF: WITH RANDOM ARRIVAL

- If jobs arrive at any time: duration a=100s, b/c=10s
- A @ t=0sec, B @ t=10sec, C @ t=10sec



Average turnaround time =  $\frac{100 + (110 - 10) + (120 - 10)}{3} = 107.33 \text{ sec}$

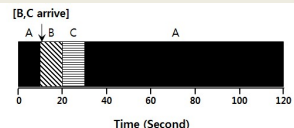
January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.38

STCF - 2

- Consider: duration a=100sec, b/c=10sec
- A<sub>len</sub>=100 A<sub>arrival</sub>=0
- B<sub>len</sub>=10, B<sub>arrival</sub>=10, C<sub>len</sub>=10, C<sub>arrival</sub>=10



Average turnaround time =  $\frac{(120 - 0) + (20 - 10) + (30 - 10)}{3} = 50 \text{ sec}$

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.39

SCHEDULING METRICS - 3

- Scheduling Metric #3: **Response Time**
- Time from when job arrives until it starts execution

$T_{\text{response}} = T_{\text{firstrun}} - T_{\text{arrival}}$

- STCF, SJF, FIFO
  - can perform poorly with respect to response time

What scheduling algorithm(s) can help minimize response time?

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.40

RR: ROUND ROBIN

- Run each job awhile, then switch to another distributing the CPU evenly (fairly)
- Scheduling Quantum is called a time slice
- Time a multimer interrupt period.

RR is fair, but performs poorly on metrics such as turnaround time

Process	Burst Time
P1	12
P5	5

Round Robin scheduling algorithm Gantt chart

P1	P2	P3	P4	P5	P1	P2	P4	P1	
0	5	10	14	19	24	29	32	37	39

Scheduling Quantum = 5 seconds

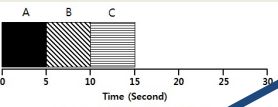
January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.41

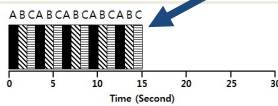
RR EXAMPLE

- ABC arrive at time=0, each run for 5 seconds



OVERHEAD not considered

$T_{\text{average response}} = \frac{0 + 5 + 10}{3} = 5 \text{ sec}$



RR with a time-slice of 1sec (Good for Response Time)

$T_{\text{average response}} = \frac{0 + 1 + 2}{3} = 1 \text{ sec}$

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.42

ROUND ROBIN: TRADEOFFS

Short Time Slice

Fast Response Time

High overhead from context switching

Long Time Slice

Slow Response Time

Low overhead from context switching

Time slice impact:

Turnaround time (for earlier example):  
ts(1,2,3,4,5)=14,14,13,14,10

Fairness: round robin is always fair, J=1

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.43

SCHEDULING WITH I/O

STCF scheduler

A: CPU=50ms, I/O=40ms, 10ms intervals

B: CPU=50ms, I/O=0ms

Consider A as 10ms subjobs (CPU, then I/O)

Without considering I/O:

CPU utilization= 100/140=71%

Poor Use of Resources

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.44

SCHEDULING WITH I/O - 2

When a job initiates an I/O request

A is blocked, waits for I/O to complete, frees CPU

STCF scheduler assigns B to CPU

When I/O completes → raise interrupt

Unblock A, STCF goes back to executing A: (10ms sub-job)

Cpu utilization = 100/100=100%

Overlap Allows Better Use of Resources

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.45

Which scheduler, thus far, best address fairness and average response time of jobs?

Respond at PollEv.com/wesleylloyd641

Text WESLEYLLOYD641 to 223333 once to join, then 1, 2, 3, 4, 5...

First In - First Out (FIFO)

1

Shortest Job First (SJF)

2

Shortest Time to Completion First (STCF)

3

Round Robin

4

None of the Above

5

All of the Above

6

Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

Total Results

CHAPTER 8 –  
MULTI-LEVEL FEEDBACK  
QUEUE (MLFQ) SCHEDULER

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.47

MULTI-LEVEL FEEDBACK QUEUE

Objectives:

Improve turnaround time:  
Run shorter jobs first

Minimize response time:  
Important for interactive jobs (UI)

Achieve without a priori knowledge of job length

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.48



MLFQ - 2

Round-Robin within a Queue

- Multiple job queues
- Adjust job priority based on observed behavior
  - Interactive Jobs
    - Frequent I/O → keep priority high
    - Interactive jobs require fast response time (GUI/UI)
  - Batch Jobs
    - Require long periods of CPU utilization
    - Keep priority low

[High Priority]

Q8 → A → B

Q7

Q6

Q5

Q4 → C

Q3

Q2

[Low Priority]

Q1 → D

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.49

MLFQ: DETERMINING JOB PRIORITY

- New arriving jobs are placed into highest priority queue
- If a job uses its entire time slice, priority is reduced (↓)
  - Jobs appears CPU-bound ( "batch" job), not interactive (GUI/UI)
- If a job relinquishes the CPU for I/O priority stays the same

MLFQ approximates SJF

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.50

MLFQ: LONG RUNNING JOB

- Three-queue scheduler, time slice=10ms

Priority

Q2

Q1

Q0

050100150200

Long-running Job Over Time (msec)

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.51

MLFQ: BATCH AND INTERACTIVE JOBS

- A<sub>arrival\_time</sub> = 0ms, A<sub>run\_time</sub> = 200ms,
- B<sub>run\_time</sub> = 20ms, B<sub>arrival\_time</sub> = 100ms

Priority

Q2

Q1

Q0

050100150200

Scheduling multiple jobs (ms)

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.52

MLFQ: BATCH AND INTERACTIVE - 2

- Continuous interactive job (B) with long running batch job (A)
  - Low response time is good for B
  - A continues to make progress

The MLFQ approach keeps interactive job(s) at the highest priority

Q2

Q1

Q0

050100150200

A Mixed I/O-intensive and CPU-intensive Workload (msec)

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.53

MLFQ: ISSUES

- Starvation

[High Priority]

Q8 → A → B → C → D → E → F

Q7

Q6

Q5

Q4

Q3

Q2

[Low Priority]

Q1 → G → H

CPU bound batch job(s)

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.54

MLFQ: ISSUES - 2

■ Gaming the scheduler

■ Issue I/O operation at 99% completion of the time slice

■ Keeps job priority fixed – never lowered

■ Job behavioral change

■ CPU/batch process becomes an interactive process

(High Priority) Q8 → (A) → (B) → (C) → (D) → (E) → (F)

Q7

Q6

Q5

Q4

Q3

Q2

(Low Priority) Q1 → (G) → (H) CPU bound batch jobs

Priority becomes stuck

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.55

RESPONDING TO BEHAVIOR CHANGE

↑

Q2

Q1

Q0

Without Priority Boost

Starvation

A: B: C:

■ Priority Boost

■ Reset all jobs to topmost queue after some time interval S

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.56

RESPONDING TO BEHAVIOR CHANGE - 2

■ With priority boost

■ Prevents starvation

↑

Q2

Q1

Q0

Without(Left) and With(Right) Priority Boost

A: B: C:

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.57

PREVENTING GAMING

■ Improved time accounting:

■ Track total job execution time in the queue

■ Each job receives a fixed time allotment

■ When allotment is exhausted, job priority is lowered

Q2

Q1

Q0

Without(Left) and With(Right) Gaming Tolerance

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.58

MLFQ: TUNING

■ Consider the tradeoffs:

■ How many queues?

■ What is a good time slice?

■ How often should we “Boost” priority of jobs?

■ What about different time slices to different queues?

Q2

Q1

Q0

Example) 10ms for the highest queue, 20ms for the middle, 40ms for the lowest

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.59

PRACTICAL EXAMPLE

■ Oracle Solaris MLFQ implementation

■ 60 Queues →  
w/ slowly increasing time slice (high to low priority)

■ Provides sys admins with set of editable table(s)

■ Supports adjusting time slices, boost intervals, priority changes, etc.

■ Advice

■ Provide OS with hints about the process

■ Nice command → Linux

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.60

MLFQ RULE SUMMARY

- The refined set of MLFQ rules:
  - Rule 1:** If  $\text{Priority}(A) > \text{Priority}(B)$ , A runs (B doesn't).
  - Rule 2:** If  $\text{Priority}(A) = \text{Priority}(B)$ , A & B run in RR.
  - Rule 3:** When a job enters the system, it is placed at the highest priority.
  - Rule 4:** Once a job uses up its time allotment at a given level (regardless of how many times it has given up the CPU), its priority is reduced(i.e., it moves down on queue).
  - Rule 5:** After some time period S, move all the jobs in the system to the topmost queue.

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

LA.61

Jackson deploys a 3-level MLFQ scheduler. The time slice is 1 for high priority jobs, 2 for medium priority, and 4 for low priority. This MLFQ scheduler performs a Priority Boost every 6 timer units. When the priority boost fires, the current job is preempted, and the next scheduled job is run in round-robin order.

Job	Arrival Time	Job Length
A	T=0	4
B	T=0	16
C	T=0	8

(11 points) Show a scheduling graph for the MLFQ scheduler for the jobs above. Draw vertical lines for key events and be sure to label the X-axis times as in the example. Please draw clearly. An unreadable graph will lose points.

HIGH

MED

LOW

0

QUESTIONS

COMPUTER BOOT SEQUENCE:  
OS WITH DIRECT EXECUTION

- What if programs could directly control the CPU / system?

OS	Program
1. Create entry for process list 2. Allocate memory for program 3. Load program into memory 4. Set up stack with <code>argc / argv</code> 5. Clear registers 6. Execute <code>call main()</code>	7. Run <code>main()</code> 8. Execute <code>return from main()</code>
9. Free memory of process 10. Remove from process list	

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

LA.64

COMPUTER BOOT SEQUENCE:  
OS WITH DIRECT EXECUTION

- What if programs could directly control the CPU / system?

OS	Program
1. Create entry for process list 2. Allocate memory for <div>Without <i>limits</i> on running programs, the OS wouldn't be in control of anything and would "just be a library"</div> <code>argv</code> 5. Clear registers 6. Execute <code>call main()</code>	7. Run <code>main()</code> 8. Execute <code>return from main()</code>
9. Free memory of process 10. Remove from process list	

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

LA.65

DIRECT EXECUTION - 2

- With direct execution:

How does the OS stop a program from running, and switch to another to support **time sharing**?

How do programs share disks and perform I/O if they are given direct control? Do they know about each other?

With direct execution, how can dynamic memory structures such as linked lists grow over time?

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

LA.66

CONTROL TRADEOFF

■ **Too little control:**

■ No security

■ No time sharing

■ **Too much control:**

■ Too much OS overhead

■ Poor performance for compute & I/O

■ Complex APIs (system calls), difficult to use

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.67

CONTEXT SWITCHING OVERHEAD

Context Switching

Multitasking

vs. Multitasking with context switching

Sequential

Overhead

Time

Total cost of context switching

January 16, 2019

TCCS422: Operating Systems [Winter 2019]  
School of Engineering and Technology, University of Washington - Tacoma

L4.68