# TCSS 422: OPERATING SYSTEMS

## CPU Schedulers: MLFQ, Proportional Share Schedulers, Linux Completely Fair Scheduler

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

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington Tacoma

---

## OBJECTIVES – 4/16

- **Questions from 4/14**
- **C Tutorial**
- **Active Reading Quiz – Ch. 7**
- **Assignment 0**
- **Chapter 8: Multi-level Feedback Queue**
  - Examples
- **Chapter 9: Proportional Share Schedulers**
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- **Chapter 26: Concurrency: An Introduction**

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma
L6.2

---

## MATERIAL / PACE

- Please classify your perspective on material covered in today's class (45 respondents):
- 1-mostly review, 5-equal new/review, 10-mostly new
- **Average – 7.6 (↓ from 7.875)**

- Please rate the pace of today's class:
- 1-slow, 5-just right, 10-fast
- **Average – 5.45 (↓ from 5.93)**

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma
L6.3

---

## FEEDBACK FROM 4/14

- No survey questions from 4/14

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma
L6.4

---

## OBJECTIVES – 4/16

- Questions from 4/14
- **C Tutorial**
- Active Reading Quiz – Ch. 7
- Assignment 0
- Chapter 8: Multi-level Feedback Queue
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma
L6.5

---

## OBJECTIVES – 4/16

- Questions from 4/14
- C Tutorial
- **Active Reading Quiz – Ch. 7**
- Assignment 0
- Chapter 8: Multi-level Feedback Queue
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma
L6.6

## OBJECTIVES – 4/16

- Questions from 4/14
- C Tutorial
- Active Reading Quiz – Ch. 7
- Assignment 0
- Chapter 8: Multi-level Feedback Queue
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction

April 16, 2020 — TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma — L6.7

---

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

April 16, 2020 — TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma — L6.8

---

## OBJECTIVES – 4/16

- Questions from 4/14
- C Tutorial
- Active Reading Quiz – Ch. 7
- Assignment 0
- Chapter 8: Multi-level Feedback Queue
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction

April 16, 2020 — TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma — L6.9

---

## 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

April 16, 2020 — TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma — L6.10

---

## 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

April 16, 2020 — TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma — L6.11

---

## RESPONDING TO BEHAVIOR CHANGE



Without Priority Boost — A: B: C:

- Priority Boost
  - Reset all jobs to topmost queue after some time interval S

April 16, 2020 — TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma — L6.12

---

## RESPONDING TO BEHAVIOR CHANGE - 2

- With priority boost
  - Prevents starvation



With Priority Boost    A: ■    B: ▨    C: ▤

## 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



Without(Left) and With(Right) Gaming Tolerance

## 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?



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

## MLFQ RULE SUMMARY

- The refined set of MLFQ rules:
- **Rule 1:** If Priority(A) > Priority(B), A runs (B doesn't).
- **Rule 2:** If Priority(A) = 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.

---

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 loose points.

HIGH

MED

LOW

0

---

## EXAMPLE

- Question:
- Given a system with a quantum length of 10 ms in its highest queue, how often would you have to boost jobs back to the highest priority level to guarantee that a single long-running (and potentially starving) job gets at least 5% of the CPU?

- Some combination of n short jobs runs for a total of 10 ms per cycle without relinquishing the CPU
  - E.g. 2 jobs = 5 ms ea; 3 jobs = 3.33 ms ea, 10 jobs = 1 ms ea
  - n jobs always uses full time quantum (10 ms)
  - Batch jobs starts, runs for full quantum of 10ms
  - All other jobs run and context switch totaling the quantum per cycle
  - If 10ms is 5% of the CPU, when must the priority boost be ???
  - **ANSWER → _Priority boost should occur every 200ms_**

# CHAPTER 9 - PROPORTIONAL SHARE SCHEDULER

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma

L6.19

---

## OBJECTIVES – 4/16

- **Questions from 4/14**
- **C Tutorial**
- **Active Reading Quiz – Ch. 7**
- **Assignment 0**
- **Chapter 8: Multi-level Feedback Queue**
  - Examples
- **Chapter 9: Proportional Share Schedulers**
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- **Chapter 26: Concurrency: An Introduction**

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma

L6.20

---

## PROPORTIONAL SHARE SCHEDULER

- **Also called fair-share scheduler or lottery scheduler**

- **Guarantees each job receives some percentage of CPU time based on share of "tickets"**

- **Each job receives an allotment of tickets**

- **% of tickets corresponds to potential share of a resource**

- **Can conceptually schedule any resource this way**
  - CPU, disk I/O, memory

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma

L6.21

---

## LOTTERY SCHEDULER

- **Simple implementation**

  - **Just need a random number generator**
    - Picks the winning ticket

  - **Maintain a data structure of jobs and tickets (list)**

  - **Traverse list to find the owner of the ticket**

  - **Consider sorting the list for speed**

April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma

L6.22

---

## LOTTERY SCHEDULER IMPLEMENTATION

head → Job:A Tix:100 → Job:B Tix:50 → Job:C Tix:250 → NULL

```
1    // counter: used to track if we've found the winner yet
2    int counter = 0;
3
4    // winner: use some call to a random number generator to
5    // get a value, between 0 and the total # of tickets
6    int winner = getrandom(0, totaltickets);
7
8    // current: use this to walk through the list of jobs
9    node_t *current = head;
10
11   // loop until the sum of ticket values is > the winner
12   while (current) {
13       counter = counter + current->tickets;
14       if (counter > winner)
15           break; // found the winner
16       current = current->next;
17   }
18   // 'current' is the winner: schedule it...
```
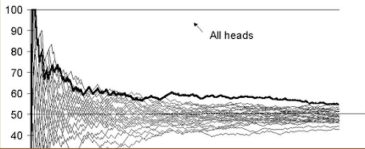
April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma

L6.23

---

## TICKET MECHANISMS

- **Ticket currency / exchange**
  - **User allocates tickets in any desired way**
  - **OS converts user currency into global currency**

- **Example:**
  - **There are 200 global tickets assigned by the OS**

    User A  → *500* (A's currency) to A1 → *50* (global currency)
            → *500* (A's currency) to A2 → *50* (global currency)

    User B  → *10* (B's currency) to B1 → *100* (global currency)
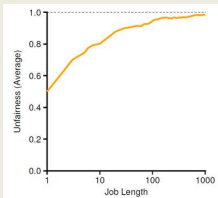
April 16, 2020
TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma

L6.24

---

## TICKET MECHANISMS - 2

- Ticket transfer
  - Temporarily hand off tickets to another process

- Ticket inflation
  - Process can temporarily raise or lower the number of tickets it owns
  - If a process needs more CPU time, it can boost tickets.

## LOTTERY SCHEDULING

- Scheduler picks a **winning** ticket
  - Load the job with the winning ticket and run it

- Example:
  - Given 100 tickets in the pool
    - Job A has 75 tickets: 0 - 74
    - Job B has 25 tickets: 75 – 99

| Scheduler's winning tickets: | 63 85 70 39 76 17 29 41 36 39 10 99 68 83 63 |
| Scheduled job: | A B A A B A A A A A B A B A |

- But what do we know about probability of a coin flip?

## COIN FLIPPING

- Equality of distribution (fairness) requires a lot of flips!

> Similarly,
> Lottery scheduling requires lots of "rounds" to achieve fairness.

## LOTTERY FAIRNESS

- With two jobs
  - Each with the same number of tickets (t=100)

> When the job length is not very long,
> average unfairness can be _quite severe_.

## LOTTERY SCHEDULING CHALLENGES

- What is the best approach to assign tickets to jobs?
  - Typical approach is to assume users know best
  - Users are provided with tickets, which they allocate as desired

- How should the OS automatically distribute tickets upon job arrival?
  - What do we know about incoming jobs a priori ?
  - Ticket assignment is really an open problem…

# TCSS 422 WILL RETURN AT ~2:40PM

## STRIDE SCHEDULER

- Addresses statistical probability issues with lottery scheduling

- Instead of guessing a random number to select a job, simply count…

## STRIDE SCHEDULER - 2

- Jobs have a "stride" value
  - A stride value describes the counter pace when the job should give up the CPU
  - Stride value is **inverse in proportion** to the job's number of tickets (more tickets = smaller stride)

- Total system tickets = 10,000
  - Job A has 100 tickets → $A_{stride}$ = 10000/100 = 100 stride
  - Job B has 50 tickets → $B_{stride}$ = 10000/50 = 200 stride
  - Job C has 250 tickets → $C_{stride}$ = 10000/250 = 40 stride

- Stride scheduler tracks "pass" values for each job (A, B, C)

## STRIDE SCHEDULER - 3

- Basic algorithm:
  1. Stride scheduler picks job with the lowest pass value
  2. Scheduler increments job's pass value by its stride and starts running
  3. Stride scheduler increments a counter
  4. When counter exceeds pass value of current job, pick a new job (go to 1)

- <u>KEY:</u> When the counter reaches a job's "PASS" value, the scheduler <u>passes</u> on to the next job…

## STRIDE SCHEDULER - EXAMPLE

- Stride values
  - Tickets = priority to select job
  - Stride is inverse to tickets
  - Lower stride = more chances to run <u>(higher priority)</u>

<u>Priority</u>
C stride = 40
A stride = 100
B stride = 200

## STRIDE SCHEDULER EXAMPLE - 2

- <u>Three-way tie</u>: randomly pick job A (all pass values=0)
- Set A's pass value to A's stride = 100
- Increment counter until > 100
- Pick a new job: <u>two-way tie</u>

Tickets
C = 250
A = 100
B = 50

| Pass(A) (stride=100) | Pass(B) (stride=200) | Pass(C) (stride=40) | Who Runs? |
|---|---|---|---|
| 0 | 0 | 0 | A |
| 100 | 0 | 0 | B |
| 100 | 200 | 0 | C |
| 100 | 200 | 40 | C |
| 100 | 200 | 80 | C |
| 100 | 200 | 120 | A |
| 200 | 200 | 120 | C |
| 200 | 200 | 160 | C |
| 200 | 200 | 200 | … |

Initial job selection is random. All @ 0

C has the most tickets and receives a lot of opportunities to run…

## STRIDE SCHEDULER EXAMPLE - 3

- We set A's counter (pass value) to A's stride = 100
- Next scheduling decision between B (pass=0) and C (pass=0)
  - Randomly choose B
- C has the lowest counter for next 3 rounds

Tickets
C = 250
A = 100
B = 50

| Pass(A) (stride=100) | Pass(B) (stride=200) | Pass(C) (stride=40) | Who Runs? |
|---|---|---|---|
| 0 | 0 | 0 | A |
| 100 | 0 | 0 | B |
| 100 | 200 | 0 | C |
| 100 | 200 | 40 | C |
| 100 | 200 | 80 | C |
| 100 | 200 | 120 | A |
| 200 | 200 | 120 | C |
| 200 | 200 | 160 | C |
| 200 | 200 | 200 | … |

C has the most tickets and is selected to run more often …

## STRIDE SCHEDULER EXAMPLE - 4

- Job counters support determining which job to run next
- Over time jobs are scheduled to run based on their priority represented as their **share of tickets...**
- **Tickets are analogous to job priority**

| Tickets |
|---|
| C = 250 |
| A = 100 |
| B = 50 |

| Pass(A)<br>(stride=100) | Pass(B)<br>(stride=200) | Pass(C)<br>(stride=40) | Who Runs? |
|---|---|---|---|
| 0 | 0 | 0 | A |
| 100 | 0 | 0 | B |
| 100 | 200 | 0 | C |
| 100 | 200 | 40 | C |
| 100 | 200 | 80 | C |
| 100 | 200 | 120 | A |
| 200 | 200 | 120 | C |
| 200 | 200 | 160 | C |
| 200 | 200 | 200 | ... |

| April 16, 2020 | TCSS422: Operating Systems [Spring 2020]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.37 |
|---|---|---|

## LINUX: COMPLETELY FAIR SCHEDULER (CFS)

- Large Google datacenter study:
  *"Profiling a Warehouse-scale Computer"* (Kanev et al.)
- Monitored 20,000 servers over 3 years
- Found 20% of CPU time spent in the Linux kernel
- 5% of CPU time spent in the CPU scheduler!
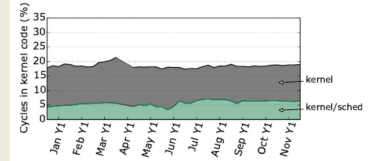- Study highlights importance for high performance OS kernels and CPU schedulers !



Figure 5: Kernel time, especially time spent in the scheduler, is a significant fraction of WSC cycles.

See: https://dl.acm.org/doi/pdf/10.1145/2749469.2750392

| April 16, 2020 | TCSS422: Operating Systems [Spring 2020]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.38 |
|---|---|---|

## LINUX: COMPLETELY FAIR SCHEDULER (CFS)

- Loosely based on the stride scheduler

- CFS models system as a Perfect Multi-Tasking System
  - In perfect system every process of the same priority (class) receive exactly $1/n^{th}$ of the CPU time

- Each scheduling class has a runqueue
  - Groups process of same class
  - In class, scheduler picks task w/ lowest `vruntime` to run
  - Time slice varies based on how many jobs in shared runqueue
  - Minimum time slice prevents too many context switches (e.g. 3 ms)

| April 16, 2020 | TCSS422: Operating Systems [Spring 2020]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.39 |
|---|---|---|

## COMPLETELY FAIR SCHEDULER - 2

- Every thread/process has a scheduling class (policy):
- **Normal classes**: SCHED_OTHER (TS), SCHED_IDLE, SCHED_BATCH
  - TS = Time Sharing
- **Real-time classes**: SCHED_FIFO (FF), SCHED_RR (RR)

- How to show scheduling class and priority:
- `#class`
  `ps –elfc`

- `#priority (nice value)`
  `ps ax -o pid,ni,cls,pri,cmd`

| April 16, 2020 | TCSS422: Operating Systems [Spring 2020]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.40 |
|---|---|---|

## COMPLETELY FAIR SCHEDULER - 3

- Linux ≥ 2.6.23: Completely Fair Scheduler (CFS)
- Linux < 2.6.23: O(1) scheduler

- Linux maintains simple counter (vruntime) to track how long each thread/process has run
- CFS picks process with lowest vruntime to run next

- CFS adjusts timeslice based on # of proc waiting for the CPU
- Kernel parameters that specify CFS behavior:
  ```
  $ sudo sysctl kernel.sched_latency_ns
  kernel.sched_latency_ns = 24000000
  $ sudo sysctl kernel.sched_min_granularity_ns
  kernel.sched_min_granularity_ns = 3000000
  $ sudo sysctl kernel.sched_wakeup_granularity_ns
  kernel.sched_wakeup_granularity_ns = 4000000
  ```

| April 16, 2020 | TCSS422: Operating Systems [Spring 2020]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.41 |
|---|---|---|

## COMPLETELY FAIR SCHEDULER - 4

- `Sched_min_granularity_ns` (3ms)
  - Time slice for a process: busy system (w/ full runqueue)
  - If system has idle capacity, time slice exceed the min as long as difference in `vruntime` between running process and process with lowest `vruntime` is less than `sched_wakeup_granularity_ns` (4ms)
- Scheduling time period is: total cycle time for iterating through a set of processes where each is allowed to run (like round robin)
- Example:
  `sched_latency_ns` (24ms)
  if (proc in runqueue < `sched_latency_ns/sched_min_granularity`)
  or
  `sched_min_granularity` * number of processes in runqueue

Ref: https://www.systutorials.com/sched_min_granularity_ns-sched_latency_ns-cfs-affect-timeslice-processes/

| April 16, 2020 | TCSS422: Operating Systems [Spring 2020]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.42 |
|---|---|---|

---

## CFS TRADEOFF

- **HIGH**     sched_min_granularity_ns (timeslice)
    sched_latency_ns
    sched_wakeup_granularity_ns

  reduced context switching → less overhead
  poor near-term fairness

- **LOW**     sched_min_granularity_ns (timeslice)
    sched_latency_ns
    sched_wakreup_granularity_ns

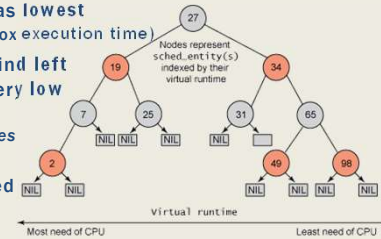  increased context switching → more overhead
  better near-term fairness

---

## COMPLETELY FAIR SCHEDULER - 5

- Runqueues are stored using a linux red-black tree
  - Self balancing binary tree - nodes indexed by **vruntime**
- Leftmost node has lowest **vruntime** (approx execution time)
- Walking tree to find left most node has very low big O complexity:
  *~O(log N) for N nodes*
- Completed processes removed

---

## CFS: JOB PRIORITY

- Time slice: Linux *"Nice value"*
  - Nice predates the CFS scheduler
  - Top shows nice values
  - Process command (nice & priority):
    `ps ax –o pid,ni,cmd,%cpu, pri`

```
static const int prio_to_weight[40] = {
 /*  -20 */ 88761, 71755, 56483, 46273, 36291,
 /*  -15 */ 29154, 23254, 18705, 14949, 11916,
 /*  -10 */  9548,  7620,  6100,  4904,  3906,
 /*   -5 */  3121,  2501,  1991,  1586,  1277,
 /*    0 */  1024,   820,   655,   526,   423,
 /*    5 */   335,   272,   215,   172,   137,
 /*   10 */   110,    87,    70,    56,    45,
 /*   15 */    36,    29,    23,    18,    15,
};
```

- Nice Values: from -20 to 19
  - Lower is *higher* priority, default is 0
  - Vruntime is a weighted time measurement
  - Priority weights the calculation of vruntime within a runqueue to give high priority jobs a boost.
    - Influences job's position in rb-tree

---

## COMPLETELY FAIR SCHEDULER - 6

- CFS tracks cumulative job run time in **vruntime** variable
- The task on a given runqueue with the lowest **vruntime** is scheduled next
- **struct sched_entity** contains **vruntime** parameter
  - Describes process execution time in nanoseconds
  - Value is not pure runtime, is weighted based on job priority
  - Perfect scheduler → achieve equal **vruntime** for all processes of same priority
- Sleeping jobs: upon return reset vruntime to lowest value in system
  - Jobs with frequent short sleep *SUFFER !!*
- Key takeaway:
  *Identifying the next job to schedule is really fast!*

---

## COMPLETELY FAIR SCHEDULER - 7

- More information:

- Man page: "man sched" : Describes Linux scheduling API
- http://manpages.ubuntu.com/manpages/bionic/man7/sched.7.html

- https://www.kernel.org/doc/Documentation/scheduler/sched-design-CFS.txt
- https://en.wikipedia.org/wiki/Completely_Fair_Scheduler

- See paper: The Linux Scheduler – a Decade of Wasted Cores
- http://www.ece.ubc.ca/~sasha/papers/eurosys16-final29.pdf

---

# CHAPTER 26 - CONCURRENCY: AN INTRODUCTION

---

## OBJECTIVES – 4/16

- Questions from 4/14
- C Tutorial
- Active Reading Quiz – Ch. 7
- Assignment 0
- Chapter 8: Multi-level Feedback Queue
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction

April 16, 2020 | TCSS422: Operating Systems [Spring 2020] School of Engineering and Technology, University of Washington - Tacoma | L6.49

## OBJECTIVES

- Introduction to threads

- Race condition

- Critical section

- Thread API

April 16, 2020 | TCSS422: Operating Systems [Spring 2020] School of Engineering and Technology, University of Washington - Tacoma | L6.50

## THREADS



©Alfred Park, http://randu.org/tutorials/threads

April 16, 2020 | TCSS422: Operating Systems [Spring 2020] School of Engineering and Technology, University of Washington - Tacoma | L6.51

## THREADS - 2

- Enables a single process (program) to have multiple "workers"
  - This is parallel programming…

- Supports independent path(s) of execution within a program with shared memory …

- Each thread has its own Thread Control Block (TCB)
  - PC, registers, SP, and stack

- Threads share code segment, memory, and heap are shared

- **What is an embarrassingly parallel program?**

April 16, 2020 | TCSS422: Operating Systems [Spring 2020] School of Engineering and Technology, University of Washington - Tacoma | L6.52

## PROCESS AND THREAD METADATA

- Thread Control Block vs. Process Control Block



April 16, 2020 | TCSS422: Operating Systems [Spring 2020] School of Engineering and Technology, University of Washington - Tacoma | L6.53

## SHARED ADDRESS SPACE

- Every thread has it's own stack / PC



April 16, 2020 | TCSS422: Operating Systems [Spring 2020] School of Engineering and Technology, University of Washington - Tacoma | L6.54

## THREAD CREATION EXAMPLE

```
#include <stdio.h>
#include <assert.h>
#include <pthread.h>

void *mythread(void *arg) {
    printf("%s\n", (char *) arg);
    return NULL;
}

int
main(int argc, char *argv[]) {
    pthread_t p1, p2;
    int rc;
    printf("main: begin\n");
    rc = pthread_create(&p1, NULL, mythread, "A"); assert(rc == 0);
    rc = pthread_create(&p2, NULL, mythread, "B"); assert(rc == 0);
    // join waits for the threads to finish
    rc = pthread_join(p1, NULL); assert(rc == 0);
    rc = pthread_join(p2, NULL); assert(rc == 0);
    printf("main: end\n");
    return 0;
}
```

April 16, 2020 | TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma | L6.55

## POSSIBLE ORDERINGS OF EVENTS

| Int main() | Thread 1 | Thread 2 |
|---|---|---|
| Starts running | | |
| Prints 'main: begin' | | |
| Creates Thread 1 | | |
| Creates Thread 2 | | |
| Waits for T1 | | |
| | Runs | |
| | Prints 'A' | |
| | Returns | |
| Waits for T2 | | |
| | | Runs |
| | | Prints 'B' |
| | | Returns |
| Prints 'main: end' | | |

April 16, 2020 | TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma | L6.56

## POSSIBLE ORDERINGS OF EVENTS - 2

| Int main() | Thread 1 | Thread 2 |
|---|---|---|
| Starts running | | |
| Prints 'main: begin' | | |
| Creates Thread 1 | | |
| | Runs | |
| | Prints 'A' | |
| | Returns | |
| Creates Thread 2 | | |
| | | Runs |
| | | Prints 'B' |
| | | Returns |
| Waits for T1 | Returns immediately | |
| Waits for T2 | | Returns immediately |
| Prints 'main: end' | | |

April 16, 2020 | TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma | L6.57

## POSSIBLE ORDERINGS OF EVENTS - 3

| Int main() | Thread 1 | Thread 2 |
|---|---|---|
| Starts running | | |
| Prints 'main: begin' | | |
| Creates Thread 1 | | |
| Creates Thread 2 | | |
| Waits for T1 | | |
| | Runs | |
| | Prints 'A' | |
| | Returns | |
| Waits for T2 | | Immediately returns |
| Prints 'main: end' | | |

**What if execution order of events in the program matters?**

April 16, 2020 | TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma | L6.58

## COUNTER EXAMPLE

- Counter example

- A + B : ordering
- Counter: incrementing global variable by two threads

- *Is the counter example embarrassingly parallel?*

- *What does the parallel counter program require?*

April 16, 2020 | TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma | L6.59

## PROCESSES VS. THREADS

- What's the difference between forks and threads?
  - Forks: duplicate a process
  - Think of *CLONING* - There will be two identical processes at the end
  - Threads: no duplicate of code/heap, lightweight execution threads



April 16, 2020 | TCSS422: Operating Systems [Spring 2020]
School of Engineering and Technology, University of Washington - Tacoma | L6.60

## RACE CONDITION

- What is happening with our counter?
  - When counter=50, consider code: counter = counter + 1
  - If synchronized, counter will = 52

| OS | Thread1 | Thread2 | PC | %eax | counter |
|---|---|---|---|---|---|
| | | (after instruction) | | | |
| | before critical section | | 100 | 0 | 50 |
| | mov 0x8049a1c, %eax | | 105 | 50 | 50 |
| | add $0x1, %eax | | 108 | 51 | 50 |
| interrupt | | | | | |
| save T1's state | | | | | |
| restore T2's state | | | 100 | 0 | 50 |
| | | mov 0x8049a1c, %eax | 105 | 50 | 50 |
| | | add $0x1, %eax | 108 | 51 | 50 |
| | | mov %eax, 0x8049a1c | 113 | 51 | 51 |
| interrupt | | | | | |
| save T2's state | | | | | |
| restore T1's state | | | 108 | 51 | 50 |
| | mov %eax, 0x8049a1c | | 113 | 51 | 51 |

## CRITICAL SECTION

- Code that accesses a shared variable must not be *concurrently* executed by more than one thread

- Multiple *active* threads inside a *critical section* produce a *race condition*.

- *Atomic execution* (*all code executed as a unit*) must be ensured in *critical* sections
  - These sections must be *mutually exclusive*

## LOCKS

- To demonstrate how critical section(s) can be executed "atomically-*as a unit*" Chapter 27 & beyond introduce locks

```
1    lock_t mutex;
2    . . .
3    lock(&mutex);
4    balance = balance + 1;          Critical section
5    unlock(&mutex);
```

- Counter example revisited

# QUESTIONS

# WILL RETURN IN A FEW MINUTES