# TCSS 422: OPERATING SYSTEMS

**Proportional Share Scheduling,
Linux Completely Fair Scheduler,
Introduction to Concurrency**

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

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington · Tacoma

---

## FEEDBACK FROM 1/23

- **What are batch jobs?**
- "Batch jobs" originates from the legacy concept of "batch processing" in computer systems
- Batch processing involves scripted running of one or more programs that run sequentially with no human interaction
- Examples include general data processing, system "housekeeping" tasks, report generation
- Tasks may be high-volume and repetitive
- Batch jobs are long-running tasks where most of the execution time requires long interrupted access to run code on the CPU

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.2

---

## FEEDBACK - 2

- **How does MLFQ priority switching work again?**
- **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).
  - Address *gaming the scheduler* through job accounting to track to execution time
- **Rule 5:** After some time period S, move all the jobs in the system to the topmost queue.
  - *Priority boost*

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.3

---

## FEEDBACK - 3

- **Is priority for processes scheduled using the MLFQ scheduler determined solely based on use of a time quantum (for each queue)?**
- For the classic MLFQ described in Ch. 8: **YES**
- For actual implementations of MLFQ, priority (which queue a job is in) could be influenced by the job's nice value

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.4

---

## FEEDBACK:
### *EXPLAIN THE SECOND EXAMPLE AGAIN*

- 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 (let's say *Job A*) gets at least 5% of the CPU?
- Key is: "**guarantee**" and "**starving**" → assume worst case scenario
- "**Single long-running**" → implies "BATCH" job
- WORST CASE: some combination of n short jobs consumes **all** remaining time of the 10ms quantum without relinquishing the CPU
  - 2 jobs=5ms ea; 3 jobs=3.33ms ea;... *does it matter how many jobs?*
  - **The quantum is gone!** n jobs ALWAYS uses full time quantum (10 ms)
  - Batch job A starts, runs for full quantum of 10ms
  - If 10ms is 5% of the CPU, when must the priority boost be ???
  - **ANSWER** → *Priority boost should occur every 200ms*

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.5

---

## FEEDBACK - 5

- **I'm confused about how to do a scheduling graph.**
- From the in class example, at T=3 C disappears
  - Where does it go?
- Then there are two A's
  - When do letters repeat?

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.6

---

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
```

---

## OBJECTIVES

- Quiz 2: Chapter 7 Schedulers
- Assignment 0 / Linux Tutorial
- C Tutorial
- Assignment 1 Posted

- **CPU Scheduling:**
- Chapter 9 – Proportional Share Scheduler
  Linux Completely Fair Scheduler (CFS)
- **Parallel programming with P-threads:**
- Chapter 26 – Intro to concurrency
- Chapter 27 – Linux Thread API
- Chapter 28 – Intro to locks

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.8 |
|---|---|---|

---

# CHAPTER 9 - PROPORTIONAL SHARE SCHEDULER

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.9 |
|---|---|---|

---

## 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

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.10 |
|---|---|---|

---

## 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

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.11 |
|---|---|---|

---

## LOTTERY SCHEDULER IMPLEMENTATION



```
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...
```

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.12 |
|---|---|---|

---

## 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)

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.13

## 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.

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.14

## 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 A B A B A

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

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.15

## COIN FLIPPING

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



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

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.16

## 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**.

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.17

## 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...

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.18

## STRIDE SCHEDULER

- Addresses statistical probability issues with lottery scheduling

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

January 28, 2019 | TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma | L6.19

## 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)

January 28, 2019 | TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma | L6.20

## 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)

- **KEY:** When the counter reaches a job's "PASS" value, the scheduler _passes_ on to the next job…

January 28, 2019 | TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma | L6.21

## STRIDE SCHEDULER - EXAMPLE

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

  Priority
  C stride = 40
  A stride = 100
  B stride = 200

January 28, 2019 | TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma | L6.22

## 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…

January 28, 2019 | TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma | L6.23

## 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 …

January 28, 2019 | TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma | L6.24

---

## 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) (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 | ... |

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.25

---

## LINUX: COMPLETELY FAIR SCHEDULER (CFS)

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

- 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)

- Show scheduling class and priority:
- `ps –elfc`
- `ps ax -o pid,ni,cls,pri,cmd`

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.26

---

## COMPLETELY FAIR SCHEDULER - 2

- 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

- Scheduling classes each have a runqueue
  - Groups process of same priority
  - Process priority groups use different sets of runqueues for priorities
  - Scheduler picks task with lowest accumulative runtime to run
  - Time quantum varies based on how many jobs in shared runqueue
    - Time quantum is proportional to system CPU load in the runqueue
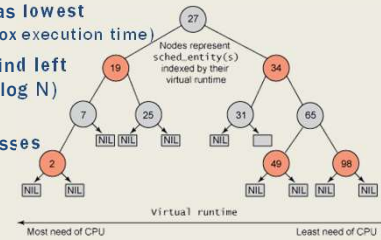    - No fixed time quantum (e.g. 10 ms)

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.27

---

## COMPLETELY FAIR SCHEDULER – 3

- 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 is ~O(log N) for N nodes
- Completed processes removed



January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.28

---

## COMPLETELY FAIR SCHEDULER - 4

- CFS tracks virtual 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, but weighted based on priority

  - Perfect scheduler → achieve equal `vruntime` for all processes of same priority

- Key takeaway:
  ***Identifying the next job to schedule is really fast!***

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.29

---

## CFS: JOB PRIORITY

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

- 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

January 28, 2019 · TCSS422: Operating Systems [Winter 2019] · School of Engineering and Technology, University of Washington - Tacoma · L6.30
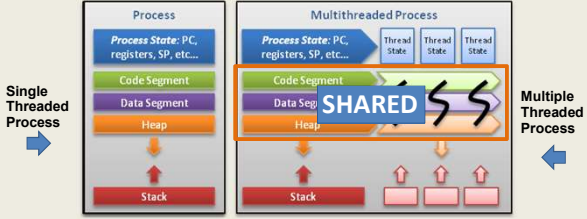
---

## CFS: TIME QUANTUM

- Scheduling quantum is calculated at runtime based on targeted latency and total number of running processes

- Will vary between:
- `cat /proc/sys/kernel/sched_min_granularity_ns` (3 ms – minimum quantum)
- `cat /proc/sys/kernel/sched_latency_ns` (24 ms – target quantum)

- Target quantum (latency):
  - Interval during which task should run at least once
  - Automatically increases as number of jobs increase

## CFS: TIME QUANTUM - 2

- How do we map a nice value to an actual CPU time quantum (timeslice) (ms)?  What is the best mapping?

- O(1) scheduler (< 2.6.23)
  - tried to map nice value to timeslice (fixed allotment)

- Linux completely fair scheduler
  - Nice value suggests priority to assign runqueue for job
  - Time proportion varies based on # of jobs in runqueue
  - With fewer jobs in runqueue, time proportion is larger

## COMPLETELY FAIR SCHEDULER - 5

- 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

- Introduction to threads

- Race condition

- Critical section

- Thread API

## THREADS



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

## 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?_**

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.37 |

## PROCESS AND THREAD METADATA

- **Thread Control Block vs. Process Control Block**

```
Thread identification
Thread state
CPU information:
        Program counter
        Register contents
Thread priority
Pointer to process that created this thread
Pointers to all other threads created by this thread
```

```
Process identification
Process status
Process state:
        Process status word
        Register contents
        Main memory
        Resources
        Process priority
Accounting
```

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.38 |

## SHARED ADDRESS SPACE

- **Every thread has it's own stack / PC**

A Single-Threaded Address Space:
- 0KB: Program Code — **The code segment:** where instructions live
- 1KB
- 2KB: Heap — **The heap segment:** contains malloc'd data dynamic data structures (it grows downward)
- (free)
- (it grows upward) **The stack segment:** contains local variables arguments to routines, return values, etc.
- 15KB
- 16KB: Stack (1)

Two threaded Address Space:
- 0KB: Program Code
- 1KB
- 2KB: Heap
- (free)
- Stack (2)
- (free)
- 15KB
- 16KB: Stack (1)

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.39 |

## THREAD CREATION EXAMPLE

```c
#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;
}
```

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.40 |

## 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' | | |

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.41 |

## 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' | | |

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.42 |

## 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?**

January 28, 2019     TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma     L6.43

---

## 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?*

January 28, 2019     TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma     L6.44

---

## 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



January 28, 2019     TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma     L6.45

---

## RACE CONDITION

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

```
                                              (after instruction)
   OS      Thread1        Thread2         PC    %eax   counter
           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
```

January 28, 2019     TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma     L6.46

---

## 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*

January 28, 2019     TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma     L6.47

---

## 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

January 28, 2019     TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma     L6.48

---

## CHAPTER 27 - LINUX THREAD API

January 28, 2019

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

L6.49

---

## THREAD CREATION

- **pthread_create**

```
#include <pthread.h>

int
pthread_create(      pthread_t*        thread,
                const pthread_attr_t* attr,
                      void*           (*start_routine)(void*),
                      void*            arg);
```

- **thread**: thread struct
- **attr**: stack size, scheduling priority… (*optional*)
- **start_routine**: function pointer to thread routine
- **arg**: argument to pass to thread routine (*optional*)

January 28, 2019

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

L6.50

---

## PTHREAD_CREATE – PASS ANY DATA

```
#include <pthread.h>

typedef struct __myarg_t {
        int a;
        int b;
} myarg_t;

void *mythread(void *arg) {
        myarg_t *m = (myarg_t *) arg;
        printf("%d %d\n", m->a, m->b);
        return NULL;
}

int main(int argc, char *argv[]) {
        pthread_t p;
        int rc;

        myarg_t args;
        args.a = 10;
        args.b = 20;
        rc = pthread_create(&p, NULL, mythread, &args);
        …
}
```

January 28, 2019

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

L6.51

---

## PASSING A SINGLE VALUE

**Using this approach on your Ubuntu VM,
How large (in bytes) can the primitive data type be?**

```
         printf("%d\n", m);
```

**How large (in bytes) can the primitive data type
be on a 32-bit operating system?**

```
9       int rc, m;
10      pthread_create(&p, NULL, mythread, (void *) 100);
11      pthread_join(p, (void **) &m);
12      printf("returned %d\n", m);
13      return 0;
14  }
```

January 28, 2019

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

L6.52

---

## WAITING FOR THREADS TO FINISH

```
int pthread_join(pthread_t thread, void **value_ptr);
```

- **thread:**    which thread?

- **value_ptr:**  pointer to return value
                type is dynamic / agnostic

- Returned values *must* be on the heap
- Thread stacks destroyed upon thread termination (join)
- Pointers to thread stack memory addresses are invalid
  - May appear as gibberish or lead to crash (seg fault)
- Not all threads join – *What would be Examples ??*

January 28, 2019

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

L6.53

---

```
struct myarg {
  int a;
  int b;
};

void *worker(void *arg)
{
  struct myarg *input = (struct myarg *) arg;
  printf("a=%d b=%d\n",input->a, input->b);
  struct myarg output;
  output.a = 1;
  output.b = 2;
  return (void *) &output;
}

int main (int argc, char * argv[])
{
  pthread_t p1;
  struct myarg args;
  struct myarg *ret_args;
  args.a = 10;
  args.b = 20;
  pthread_
  pthread_
  printf("
  return 0;
}
```

**What will this code do?**

Data on thread stack

```
$ ./pthread_struct
a=10 b=20
Segmentation fault (core dumped)
```

**How can this code be fixed?**

January 28, 2019

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

L6.54

## How about this code?

```
struct myarg {
  int a;
  int b;
};

void *worker(void *arg)
{
  struct myarg *input = (struct myarg *) arg;
  printf("a=%d b=%d\n",input->a, input->b);
  input->a = 1;
  input->b = 2;
  return (void *) &input;
}

int main (int argc, char * argv[])
{
  pthread_t p1;
  struct myarg args;
  struct myarg *ret_args;
  args.a = 10;
  args.b = 20;
  pthread_create(&p1, NULL, worker, &args);
  pthread_join(p1, (void *)&ret_args);
  printf("returned %d %d\n", ret_args->a, ret_args->b);
  return 0;
}
```

```
$ ./pthread_struct
a=10 b=20
returned 1 2
```

January 28, 2019    TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma    L6.55

---

## ADDING CASTS

- Casting
- Suppresses compiler warnings when passing "typed" data where (void) or (void *) is called for

- Example: uncasted capture in pthread_join
```
pthread_int.c: In function 'main':
pthread_int.c:34:20: warning: passing argument 2 of 'pthread_join'
from incompatible pointer type [-Wincompatible-pointer-types]
   pthread_join(p1, &p1val);
```

- Example: uncasted return
```
In file included from pthread_int.c:3:0:
/usr/include/pthread.h:250:12: note: expected 'void **' but argument
is of type 'int **'
 extern int pthread_join (pthread_t __th, void **__thread_return);
```

January 28, 2019    TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma    L6.56

---

## ADDING CASTS - 2

- pthread_join
```
int * p1val;
int * p2val;
pthread_join(p1, (void *)&p1val);
pthread_join(p2, (void *)&p2val);
```

- return from thread function
```
int * counterval = malloc(sizeof(int));
*counterval = counter;
return (void *) counterval;
```

January 28, 2019    TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma    L6.57

---

## LOCKS

- pthread_mutex_t data type
- /usr/include/bits/pthread_types.h

```
// Global Address Space
static volatile int counter = 0;
pthread_mutex_t lock;

void *worker(void *arg)
{
  int i;
  for (i=0;i<10000000;i++)  {
    int rc = pthread_mutex_lock(&lock);
    assert(rc==0);
    counter = counter + 1;
    pthread_mutex_unlock(&lock);
  }
  return NULL;
}
```

January 28, 2019    TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma    L6.58

---

## LOCKS - 2

- Ensure critical sections are executed atomically-*as a unit*
  - Provides implementation of "*Mutual Exclusion*"

- API
```
int pthread_mutex_lock(pthread_mutex_t *mutex);
int pthread_mutex_unlock(pthread_mutex_t *mutex);
```

- Example w/o initialization & error checking
```
pthread_mutex_t lock;
pthread_mutex_lock(&lock);
x = x + 1; // or whatever your critical section is
pthread_mutex_unlock(&lock);
```
  - Blocks forever until lock can be obtained
  - Enters critical section once lock is obtained
  - Releases lock

January 28, 2019    TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma    L6.59

---

## LOCK INITIALIZATION

- Assigning the constant
```
pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
```

- API call:
```
int rc = pthread_mutex_init(&lock, NULL);
assert(rc == 0); // always check success!
```

- Initializes mutex with attributes specified by 2nd argument
- If NULL, then default attributes are used
- Upon initialization, the mutex is initialized and unlocked

January 28, 2019    TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma    L6.60

## LOCKS - 3

- **Error checking wrapper**

```
// Use this to keep your code clean but check for failures
// Only use if exiting program is OK upon failure
void Pthread_mutex_lock(pthread_mutex_t *mutex) {
    int rc = pthread_mutex_lock(mutex);
    assert(rc == 0);
}
```

- **What if lock can't be obtained?**

```
int pthread_mutex_trylock(pthread_mutex_t *mutex);
int pthread_mutex_timelock(pthread_mutex_t *mutex,
                           struct timespec *abs_timeout);
```

- trylock – returns immediately (fails) if lock is unavailable
- timelock – tries to obtain a lock for a specified duration

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.61 |

---

## CONDITIONS AND SIGNALS

- **Condition variables support "signaling" between threads**

```
int pthread_cond_wait(pthread_cond_t *cond,
                      pthread_mutex_t *mutex);
int pthread_cond_signal(pthread_cond_t *cond);
```

- `pthread_cont_t` datatype

- `pthread_cond_wait()`
  - Puts thread to "sleep" (waits)   (THREAD is BLOCKED)
  - Threads added to FIFO queue, lock is released
  - Waits _(listens)_ for a "signal"   (NON-BUSY WAITING, no polling)
  - When signal occurs, interrupt fires, wakes up first thread, (THREAD is RUNNING), lock is provided to thread

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.62 |

---

## CONDITIONS AND SIGNALS - 2

```
int pthread_cond_signal(pthread_cond_t * cond);
int pthread_cond_broadcast(pthread_cond_t * cond);
```

- `pthread_cond_signal()`
  - Called to send a "signal" to wake-up first thread in FIFO "wait" queue
  - The goal is to unblock a thread to respond to the signal
- `pthread_cond_broadcast()`
  - Unblocks _all_ threads in FIFO "wait" queue, currently blocked on the specified condition variable
  - Broadcast is used when all threads should wake-up for the signal
- **Which thread is unblocked first?**
  - Determined by OS scheduler (based on priority)
  - Thread(s) awoken based on placement order in FIFO wait queue
  - When awoken threads acquire lock as in `pthread_mutex_lock()`

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.63 |

---

## CONDITIONS AND SIGNALS - 3

- **Wait example:**

```
pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
pthread_cond_t cond = PTHREAD_COND_INITIALIZER;

pthread_mutex_lock(&lock);
while (initialized == 0)
    pthread_cond_wait(&cond, &lock);
// Perform work that requires lock
a = a + b;
pthread_mutex_unlock(&lock);
```

- wait puts thread to sleep, releases lock
- when awoken, lock reacquired (but then released by this code)
- When initialized, another thread signals

State variable set, Enables other thread(s) to proceed above.

```
pthread_mutex_lock(&lock);
initialized = 1;
pthread_cond_signal(&init);
pthread_mutex_unlock(&lock);
```

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.64 |

---

## CONDITION AND SIGNALS - 4

```
pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
pthread_cond_t cond = PTHREAD_COND_INITIALIZER;

pthread_mutex_lock(&lock);
while (initialized == 0)
    pthread_cond_wait(&cond, &lock);
// Perform work that requires lock
a = a + b;
pthread_mutex_unlock(&lock);
```

- **Why do we wait inside a while loop?**

- **The while ensures upon awakening the condition is rechecked**
  - A signal is raised, but the pre-conditions required to proceed may have not been met.  **MUST CHECK STATE VARIABLE**
  - Without checking the state variable the thread may proceed to execute when it should not.  (e.g. too early)

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.65 |

---

## PTHREADS LIBRARY

- **Compilation**
  - gcc –pthread pthread.c –o pthread
  - Requires explicitly linking the library with compiler flag
  - Use makefile to provide compiler arguments

- **List of pthread manpages**
  - man –k pthread

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019] School of Engineering and Technology, University of Washington - Tacoma | L6.66 |

## SAMPLE MAKEFILE

```
CC=gcc
CFLAGS=-pthread -I. -Wall

binaries=pthread pthread_int pthread_lock_cond pthread_struct

all: $(binaries)

pthread_mult: pthread.c pthread_int.c
    $(CC) $(CFLAGS) $^ -o $@

clean:
    $(RM) -f $(binaries) *.o
```

- **Example builds multiple single file programs**
  - **All target**
- **pthread_mult**
  - **Example if multiple source files should produce a single executable**
- **clean target**

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.67

---

# CHAPTER 28 – LOCKS

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.68

---

## LOCKS

- **Ensure critical section(s) are executed atomically-*as a unit***
  - **Only one thread is allowed to execute a critical section at any given time**
  - **Ensures the code snippets are "mutually exclusive"**

- **Protect a global counter:**

  ```
  balance = balance + 1;
  ```

- **A "critical section":**

  ```
  1   lock_t mutex; // some globally-allocated lock 'mutex'
  2   …
  3   lock(&mutex);
  4   balance = balance + 1;
  5   unlock(&mutex);
  ```

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.69

---

## LOCKS - 2

- **Lock variables are called "MUTEX"**
  - **Short for mutual exclusion (that's what they guarantee)**

- **Lock variables store the state of the lock**

- **States**
  - **Locked** (acquired or held)
  - **Unlocked** (available or free)

- **Only 1 thread can hold a lock**

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.70

---

## LOCKS - 3

- **pthread_mutex_lock(&lock)**
  - **Try to acquire lock**
  - **If lock is free, calling thread will acquire the lock**
  - **Thread with lock enters critical section**
    - **Thread "owns" the lock**

- **No other thread can acquire the lock before the owner releases it.**

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.71
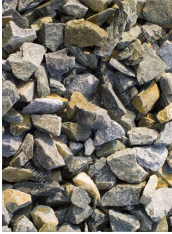
---

## LOCKS - 4

- **Program can have many mutex (lock) variables to "serialize" many critical sections**

- **Locks are also used to protect data structures**
  - **Prevent multiple threads from changing the same data simultaneously**
  - **Programmer can make sections of code "granular"**
    - **Fine grained – means just one grain of sand at a time through an hour glass**
  - **Similar to relational database transactions**
    - **DB transactions prevent multiple users from modifying a table, row, field**

January 28, 2019 — TCSS422: Operating Systems [Winter 2019]
School of Engineering and Technology, University of Washington - Tacoma — L6.72

## FINE GRAINED?

- Is this code a good example of *"fine grained parallelism"*?

```
pthread_mutex_lock(&lock);
a = b++;
b = a * c;
*d = a + b +c;
FILE * fp = fopen ("file.txt", "r");
fscanf(fp, "%s %s %s %d", str1, str2, str3, &e);
ListNode *node = mylist->head;
Int i=0
while (node) {
  node->title = str1;
  node->subheading = str2;
  node->desc = str3;
  node->end = *e;
  node = node->next;
  i++
}
e = e - i;
pthread_mutex_unlock(&lock);
```

## FINE GRAINED PARALLELISM

```
pthread_mutex_lock(&lock_a);
pthread_mutex_lock(&lock_b);
a = b++;
pthread_mutex_unlock(&lock_b);
pthread_mutex_unlock(&lock_a);

pthread_mutex_lock(&lock_b);
b = a * c;
pthread_mutex_unlock(&lock_b);

pthread_mutex_lock(&lock_d);
*d = a + b +c;
pthread_mutex_unlock(&lock_d);

FILE * fp = fopen ("file.txt", "r");
pthread_mutex_lock(&lock_e);
fscanf(fp, "%s %s %s %d", str1, str2, str3, &e);
pthread_mutex_unlock(&lock_e);

ListNode *node = mylist->head;
int i=0 . . .
```

## EVALUATING LOCK IMPLEMENTATIONS

- **Correctness**
  - **Does the lock work?**
  - **Are critical sections mutually exclusive?**
    (atomic-*as a unit*?)

- **Fairness**
  - **Are threads competing for a lock have a fair chance of acquiring it?**

- **Overhead**

## BUILDING LOCKS

- **Locks require hardware support**
  - **To minimize overhead, ensure fairness and correctness**

  - **Special "atomic-*as a unit*" instructions to support lock implementation**

  - **Atomic-*as a unit* exchange instruction**
    - XCHG

  - **Compare and exchange instruction**
    - CMPXCHG
    - CMPXCHG8B
    - CMPXCHG16B

## HISTORICAL IMPLEMENTATION

- **To implement mutual exclusion**
  - **Disable interrupts upon entering critical sections**

```
1   void lock() {
2       DisableInterrupts();
3   }
4   void unlock() {
5       EnableInterrupts();
6   }
```

- **Any thread could disable system-wide interrupt**
  - **What if lock is never released?**

- **On a multiprocessor processor each CPU has its own interrupts**
  - **Do we disable interrupts for all cores simultaneously?**

- **While interrupts are disabled, they could be lost**
  - **If not queued…**

## SPIN LOCK IMPLEMENTATION

- **Operate without atomic-*as a unit* assembly instructions**
- **"Do-it-yourself" Locks**
- **Is this lock implementation: Correct? Fair? Performant?**

```
1   typedef struct __lock_t { int flag; } lock_t;
2
3   void init(lock_t *mutex) {
4       // 0 → lock is available, 1 → held
5       mutex->flag = 0;
6   }
7
8   void lock(lock_t *mutex) {
9       while (mutex->flag == 1)  // TEST the flag
10          ;  // spin-wait (do nothing)
11      mutex->flag = 1;  // now SET it !
12  }
13
14  void unlock(lock_t *mutex) {
15      mutex->flag = 0;
16  }
```

## DIY: CORRECT?

- Correctness requires luck... (e.g. *DIY lock is incorrect*)

| Thread1 | Thread2 |
|---|---|
| call `lock()`<br>while (flag == 1)<br>interrupt: switch to Thread 2 | |
| | call `lock()`<br>while (flag == 1)<br>flag = 1;<br>interrupt: switch to Thread 1 |
| flag = 1; // set flag to 1 (too!) | |

- Here both threads have "acquired" the lock simultaneously

## DIY: PERFORMANT?

```
void lock(lock_t *mutex)
{
    while (mutex->flag == 1);   // while lock is unavailable, wait...
    mutex->flag = 1;
}
```

- What is wrong with while(<cond>); ?

- Spin-waiting wastes time actively waiting for another thread
- while (1); will "peg" a CPU core at 100%
  - Continuously loops, and evaluates mutex->flag value...
  - Generates heat...

## TEST-AND-SET INSTRUCTION

- C implementation: not atomic
  - Adds a simple check to basic spin lock
  - One a single core CPU system with preemptive scheduler:
  - Try this...

```
1   int TestAndSet(int *ptr, int new) {
2       int old = *ptr;   // fetch old value at ptr
3       *ptr = new;       // store 'new' into ptr
4       return old;       // return the old value
5   }
```

- lock() method checks that TestAndSet doesn't return 1
- Comparison is in the caller
- Single core systems are becoming scarce
- Try on a one-core VM

## DIY: TEST-AND-SET - 2

- Requires a preemptive scheduler on single CPU core system
- Lock is never released without a context switch
- 1-core VM: occasionally will deadlock, doesn't miscount

```
1    typedef struct __lock_t {
2        int flag;
3    } lock_t;
4
5    void init(lock_t *lock) {
6        // 0 indicates that lock is available,
7        // 1 that it is held
8        lock->flag = 0;
9    }
10
11   void lock(lock_t *lock) {
12       while (TestAndSet(&lock->flag, 1) == 1)
13           ;          // spin-wait
14   }
15
16   void unlock(lock_t *lock) {
17       lock->flag = 0;
18   }
```

## SPIN LOCK EVALUATION

- Correctness:
  - Spin locks guarantee: critical sections won't be executed simultaneously by (2) threads

- Fairness:
  - No fairness guarantee. Once a thread has a lock, nothing forces it to relinquish it...

- Performance:
  - Spin locks perform "busy waiting"
  - Spin locks are best for short periods of waiting
  - Performance is slow when multiple threads share a CPU
    - Especially for long periods

## COMPARE AND SWAP

- Checks that the lock variable has the expected value FIRST, before changing its value
  - If so, make assignment
  - Return value at location

- Adds a comparison to TestAndSet

- Useful for wait-free synchronization
  - Supports implementation of shared data structures which can be updated atomically (*as a unit*) using the HW support CompareAndSwap instruction
  - Shared data structure updates become "wait-free"
  - Upcoming in Chapter 32

## COMPARE AND SWAP

- Compare and Swap

```
1   int CompareAndSwap(int *ptr, int expected, int new) {
2       int actual = *ptr;
3       if (actual == expected)
4               *ptr = new;
5       return actual;
```

- Spin loc

1-core VM:
Count is correct, no deadlock

```
3               ; // spin
4   }
```

- X86 provides "cmpxchgl" compare-and-exchange instruction
  - cmpxchg8b
  - cmpxchg16b

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.85 |

---

## TWO MORE "LOCK BUILDING" CPU INSTRUCTIONS

- Cooperative instructions used together to support synchronization on RISC systems
- No support on x86 processors
  - Supported by RISC: Alpha, PowerPC, ARM
- Load-linked (LL)
  - Loads value into register
  - Same as typical load
  - Used as a mechanism to track competition
- Store-conditional (SC)
  - Performs "mutually exclusive" store
  - Allows only one thread to store value

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.86 |

---

## LL/SC LOCK

```
1   int LoadLinked(int *ptr) {
2       return *ptr;
3   }
4
5   int StoreConditional(int *ptr, int value) {
6       if (no one has updated *ptr since the LoadLinked to this address) {
7               *ptr = value;
8               return 1; // success!
9       } else {
10              return 0; // failed to update
11      }
12  }
```

- LL instruction loads pointer value (ptr)
- SC only stores if the load link pointer has not changed
- Requires HW support
  - C code is psuedo code

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.87 |

---

## LL/SC LOCK - 2

```
1   void lock(lock_t *lock) {
2       while (1) {
3           while (LoadLinked(&lock->flag) == 1)
4               ; // spin until it's zero
5           if (StoreConditional(&lock->flag, 1) == 1)
6               return; // if set-it-to-1 was a success: all done
7                       otherwise: try it all over again
8       }
9   }
10
11  void unlock(lock_t *lock) {
12      lock->flag = 0;
13  }
```

- Two instruction lock

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.88 |

---

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

| January 28, 2019 | TCSS422: Operating Systems [Winter 2019]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.90 |

## 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 | |

*Without limits on running programs, the OS wouldn't be in control of anything and would "just be a library"*

| OS | Program |
|---|---|
| argv | |
| 5. Clear registers | 7. Run main() |
| 6. Execute call main() | 8. Execute return from main() |
| 9. Free memory of process | |
| 10. Remove from process list | |

| | | |
|---|---|---|
| January 28, 2019 | TCSS422: Operating Systems [Winter 2019]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.91 |

## 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 28, 2019 | TCSS422: Operating Systems [Winter 2019]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.92 |

## 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 28, 2019 | TCSS422: Operating Systems [Winter 2019]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.93 |

## CONTEXT SWITCHING OVERHEAD



| | | |
|---|---|---|
| January 28, 2019 | TCSS422: Operating Systems [Winter 2019]<br>School of Engineering and Technology, University of Washington - Tacoma | L6.94 |