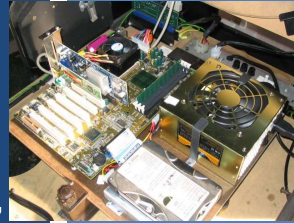# TCSS 422: OPERATING SYSTEMS

**Proportional Share Schedulers,
Linux Completely Fair Scheduler,
Introduction to Concurrency, Locks API,
Introduction to Locks**

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

October 17, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington  Tacoma

---

# FEEDBACK FROM 10/15

- **Multi-level Feedback Queue with I/O**
  - **Wikipedia explanation:** https://en.wikipedia.org/wiki/Multilevel_feedback_queue
  - **Each priority queue processes jobs in FIFO manner**
  - **Jobs always inserted at tail of FIFO queues**
  - **Scheduler selects first job in the highest priority queue to run**
  - **Only things that can happen to a job:**
    - **ANY JOB**: if finished executing is removed from queue
    - **I/O JOB**: Job goes from RUNNING→BLOCKED and is removed from the scheduler until it is READY and will be reinserted
    - **BATCH JOB**: Uses full quantum, is added to tail of next lower queue
  - **No job is run from a lower queue if higher queue is not empty**
- **KEY POINT (implicit in the textbook):
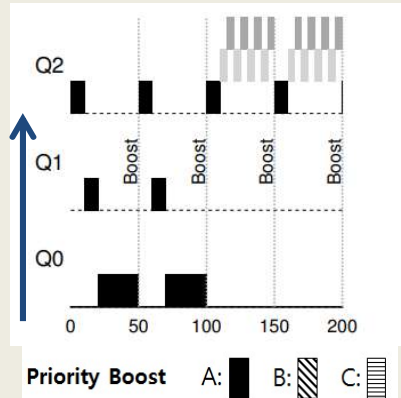  Starvation occurs because high priority queue is never empty**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington  - Tacoma | L7.2 |
|---|---|---|

# FEEDBACK - 2

- High priority queue must be empty "for a little while" for the scheduler to look at the lower queue for a job to run
- A single I/O job must go from RUNNING→BLOCKED and not use the full quantum
- At least two I/O jobs are required to cause starvation of lower jobs

- V1.0 of textbook corrects figure:
  - **_*A is on the left after the boost*_**
  - Priority boost: Prevents starvation

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma | L7.3 |
|---|---|---|

# FEEDBACK - 3

- Scheduling graph for MLFQ Scheduler
  - Using letters for jobs instead of blocks, where a letter is one timer unit (e.g. seconds or milliseconds) can be easier to debug

- How does the conversion from tickets to priority work with the Stride Scheduler?
  - Stride scheduler and lottery scheduler:
  - Jobs with highest number of tickets receive highest priority
  - Stride scheduler calculates a **stride value** that is inverse to the total number of tickets
  - Calculating **stride** requires knowing total # of system tickets

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma | L7.4 |
|---|---|---|

# OBJECTIVES

- C Tutorial (Sunday 10/21)
- Program 1 – MASH Shell (Friday 10/26)

- **CPU Scheduling cont'd:**
- Chapter 9 – Proportional Share Schedulers
- Linux - Completely Fair Scheduler (CFS)

- **Multi-threaded Programming**
- Chapter 26 – Concurrency Introduction
- Chapter 27 – Linux Thread API
- Chapter 28 – Introduction to Locks

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.5 |
|---|---|---|

# CHAPTER 9 - PROPORTIONAL SHARE SCHEDULER

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.6 |
|---|---|---|

## When used to make only a small number of job scheduling decisions which scheduling metric does the lottery scheduler perform poorly on?

**W**

Average Response Time

Average Turnaround Time

Fairness

Average Execution Time

Average Job Start Time

# PROPORTIONAL SHARE SCHEDULERS

- **How does the Lottery scheduler determine which job to run next?**

- **What problem does the job selection method cause for the Lottery scheduler?**

- **What is fundamentally different about how the stride scheduler performs job selection?**

- **Why does the different design of the stride scheduler solve the job selection problem of the lottery scheduler?**

## 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 <u>inverse in proportion</u> 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 (higher priority)

## Priority

C stride = 40

A stride = 100

B stride = 200

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.11 |

# STRIDE SCHEDULER EXAMPLE - 2

- Three-way tie: 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: two-way tie

| 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 | ... |

**Tickets**
C = 250
A = 100
B = 50

Initial job selection is random. All @ 0

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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.12 |

## 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):
- <u>Normal classes</u>: SCHED_OTHER (TS), SCHED_IDLE, SCHED_BATCH
  - TS = Time Sharing
- <u>Real-time classes</u>: SCHED_FIFO (FF), SCHED_RR (RR)

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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.13 |
|---|---|---|

## 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
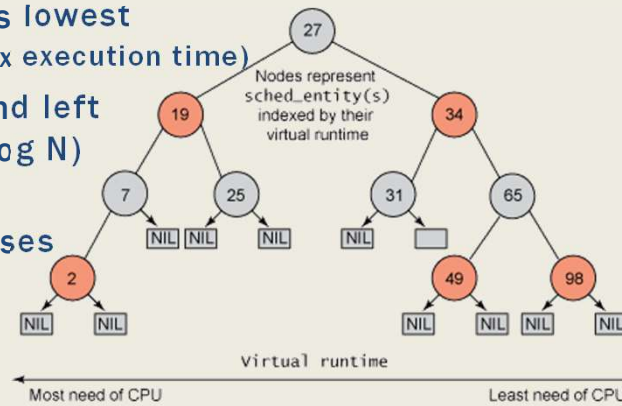    - No fixed time quantum (e.g. 10 ms)

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.14 |
|---|---|---|

## 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



| | October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.15 |

## 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!

| | October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.16 |

## 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.17 |

## 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.18 |

## 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.19 |
|---|---|---|

## 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington  - Tacoma | L7.20 |
|---|---|---|

# CHAPTER 26 - CONCURRENCY: AN INTRODUCTION

October 17, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L7.21

# OBJECTIVES

- **Introduction to threads**

- **Race condition**

- **Critical section**

- **Thread API**

October 17, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L7.22

# THREADS



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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.23 |
|---|---|---|

# THREADS - 2

- **Enables a single process (program) to have multiple "workers"**

- **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**

- **Code segment, memory, and heap are shared**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.24 |
|---|---|---|

# PROCESS AND THREAD METADATA

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



| Thread identification | Process identification |
| Thread state | Process status |
| CPU information: | Process state: |
|     Program counter |     Process status word |
|     Register contents |     Register contents |
| Thread priority |     Main memory |
| Pointer to process that created this thread |     Resources |
| Pointers to all other threads created by this thread |     Process priority |
| | Accounting |

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.25 |

---

# SHARED ADDRESS SPACE

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



**The code segment**: where instructions live

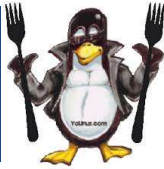**The heap segment**: contains malloc'd data dynamic data structures (it grows downward)

(it grows upward)

**The stack segment**: contains local variables arguments to routines, return values, etc.

A Single-Threaded Address Space

Two threaded Address Space

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.26 |

# PROCESSES VS. THREADS

- What's the difference between forks (processes) 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

| code | data | files |
| code | data | files |

single-threaded process

multithreaded process

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.27 |

---

# 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;
}
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.28 |

## 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' | | |

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.29 |
|---|---|---|

## 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' | | |

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.30 |
|---|---|---|

## POSSIBLE ORDERINGS OF EVENTS - 3

| Int main() | Thread 1 | Thread 2 |
|---|---|---|
| Starts running | | |
| Prints 'main: begin' | | |
| Creates Thread 1 | | |
| Creates Thread 2 | | |
| | | |
| | | |
| Waits for T: | | |
| | Runs | |
| | Prints 'A' | |
| | Returns | |
| Waits for T2 | | *Immediately returns* |
| Prints 'main: end' | | |

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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.31 |
|---|---|---|

## COUNTER EXAMPLE

- **Pthread create example (pthread_create.c)**
- **A + B : ordering**

- **Counter example (pthread.c)**
- **Counter: incrementing global variable by two threads**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.32 |
|---|---|---|

# 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
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.33 |
|---|---|---|

# 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_**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.34 |
|---|---|---|

# 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 (pthread_lock.c)**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.35 |
|---|---|---|

# CHAPTER 27 - LINUX THREAD API

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.36 |
|---|---|---|

# 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*)**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.37 |
|---|---|---|

# 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);
        …
}
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.38 |
|---|---|---|

# PASSING A SINGLE VALUE

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

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

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

 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  }
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.39 |
|---|---|---|

---

# 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 ??*

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.40 |
|---|---|---|

```
struct myarg {
  int a;                 What will this code do?
  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;        Data on thread stack
  output.a = 1;
  output.b = 2;                      $ ./pthread_struct
  return (void *) &output;           a=10 b=20
}                                    Segmentation fault (core dumped)

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

```
struct myarg {
  int a;               How about this code?
  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;
}                                    $ ./pthread_struct
                                     a=10 b=20
int main (int argc, char * argv[])   returned 1 2
{
  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;
}
```

## 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);
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.43 |
|---|---|---|

## 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;
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.44 |
|---|---|---|

# 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;
}
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.45 |
|---|---|---|

# 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.46 |
|---|---|---|

# 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.47 |
|---|---|---|

---

# 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.48 |
|---|---|---|

# 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.49 |
|---|---|---|

# 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()

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.50 |
|---|---|---|

# 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**

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

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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.51 |

# 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)

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.52 |

# 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.53 |
|---|---|---|

# 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.54 |
|---|---|---|

# CHAPTER 28 – LOCKS

# 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);
```

# 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.57 |
|---|---|---|

# 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.

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.58 |
|---|---|---|

## 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.59 |

## 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 (n
  node->
  node->subheading = str2;
  node->desc = str3;
  node->end = *e;
  node = node->next;
  i++
}
e = e – i;
pthread_mutex_unlock(&lock);
```

**Example of coarse-grained parallelism**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.60 |

## 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 . . .
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.61 |
| --- | --- | --- |

## 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.62 |
| --- | --- | --- |

# 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.63 |
|---|---|---|

# 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...**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.64 |
|---|---|---|

# 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  }
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.65 |
|---|---|---|

# 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.66 |
|---|---|---|

# 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…**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.67 |
|---|---|---|

# 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.68 |
|---|---|---|

# 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  }
```

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.69 |
|---|---|---|

# 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.70 |
|---|---|---|

# 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.71 |
|---|---|---|

---

# 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;
6
```

- **Spin loc**

```
1
2
3                ; // spin
4    }
```

**1-core VM:
Count is correct, no deadlock**

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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.72 |
|---|---|---|

# 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.73 |
|---|---|---|

# 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

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.74 |
|---|---|---|

## 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**

| October 17, 2018 | TCSS422: Operating Systems [Fall 2018]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.75 |
|---|---|---|

# QUESTIONS