

# Ch 5: Processor Management

Bill Cheng

<http://merlot.usc.edu/cs402-s16>



Copyright © William C. Cheng

Operating Systems - CSCI 402

Operating Systems - CSCI 402

- Processor Management
- Threads *implementation*
- lock/mutex implementation on multiprocessors
- Interrupts
- Scheduling
- Linux/Windows Scheduler



Copyright © William C. Cheng

## Threads Implementation

- The ultimate goal of the OS is to support user-level applications
- we will discuss various strategies for supporting threads
- Where are operations on threads implemented?
  - in the kernel?
  - or in user-level library?
- Approaches
  - one-level model (threads are implemented in the kernel)
  - variable-weight processes
  - two-level model (threads are implemented in user library)
  - $N \times 1$
  - $M \times N$
  - scheduler activations model



Copyright © William C. Cheng

4

## 5.1 Threads Implementations

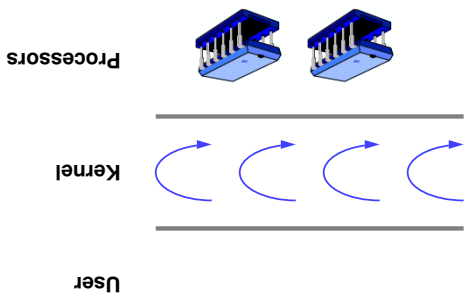
- *Strategies*
- A Simple Thread Implementation
- Multiple Processors



Copyright © William C. Cheng

Operating Systems - CSCI 402

## One-Level Model



Copyright © William C. Cheng

## One-Level Model

- The simplest and most direct approach is the *one-level model*
- all aspects of the *thread implementation* are in the kernel
- i.e., all thread routines (e.g., `pthread_mutex_lock`) called by user code are all system calls
- each *user thread* is mapped one-to-one to a *kernel thread*
- If a thread calls `pthread_create()`
  - it's a system call, so it traps into the kernel
  - the kernel creates a thread control block
  - associate it with the process control block
  - the kernel creates a kernel and a user stack for this thread
- What about `pthread_mutex_lock()`
  - why does it have to be done in the kernel?
  - it's not necessary to protect the threads from each other!
  - you definitely don't need the kernel to protect threads from each other



Copyright © William C. Cheng

6





Copyright © William C. Cheng

- ## Two-Level Model - One Kernel Thread
- Within a process, user threads are multiplexed not on the processor, but on a kernel-supported thread
  - the OS multiplexes kernel threads (or equivalently, processes) on the processor
  - kernel does *not* know about *the existence of user threads*
  - User thread creation
    - a stack and a thread control block is allocated
    - thread is put on a queue of runnable threads
    - wait for its turn to become the running thread
  - Synchronization implementation
    - relative straightforward
    - e.g., mutex (one queue per mutex)
    - if a thread must block, it simply queues itself on a wait queue and calls context-switch routine to pass control to the first thread on the runnable queue



Copyright © William C. Cheng

- ## Two-Level Model - One Kernel Thread
- This is called the *N-to-1 Model*
- Major advantage
    - fast, because no system calls for thread-related APIs
  - Major disadvantage
    - what if a thread makes a system call (for a non-thread-related API)?
    - it gets blocked in the kernel
    - no other user thread in the process can run



Copyright © William C. Cheng



Copyright © William C. Cheng

## Coping ...

```

size_t read(int fd, void *buf, size_t count)
{
    while (1) {
        if (ret = real_read(fd, buf, count)) == -1) {
            if (errno == EWOULDBLOCK) {
                sem_wait(&fileSemaphore[fd]);
                continue;
            }
            break;
        }
        return ret;
    }
}

```

Solution is to have a non-blocking read () called real\_read ()

- real\_read () either returns immediately with data in buf or returns immediately with an error code in errno
- EWOULDBLOCK means that a real read () would block, i.e., data is not ready to be read



Copyright © William C. Cheng

## Coping ...

```

size_t read(int fd, void *buf, size_t count)
{
    while (1) {
        if (ret = real_read(fd, buf, count)) == -1) {
            if (errno == EWOULDBLOCK) {
                sem_wait(&fileSemaphore[fd]);
                continue;
            }
            break;
        }
        return ret;
    }
}

```

- One semaphore for each open file
  - perhaps a signal handler will invoke sem\_post () to when data is ready to be read
- Major drawback
  - only works for some I/O objects - not a general solution



Copyright © William C. Cheng



Copyright © William C. Cheng

## Two-Level Model: Multiple Kernel Threads

This is called the *M-to-N model*

- Implementation is similar to the two-level model with a single kernel thread
- no system calls (for thread-related APIs)
- if we don't have enough kernel threads per user process, we end up having the same problem with the N-to-1 model



Copyright © William C. Cheng

## Deadlock

Ex: two threads are communicating using a pipe (this is essentially a kernel implementation of the producer-consumer problem)

- first user thread writes to a full pipe and get blocked in the kernel
- first thread just happened to use the last kernel thread
- 2nd thread wants to read the pipe to unblock the first thread, but cannot because no kernel thread left



Copyright © William C. Cheng

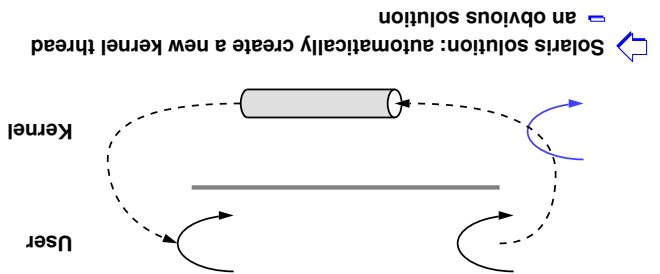
## Scheduler Activations Model Example

- ➡ Let's say a process starts up running a single thread
  - ➡ kernel scheduler assigns a processor to the process
  - ➡ if the thread blocks, the process gives up the processor to the kernel scheduler
- ➡ Suppose the user program creates a new thread and parallelism is desired
  - ➡ code in user-level library notifies the kernel that it needs two processors
  - ➡ when a processor becomes available, the kernel creates a new kernel context
  - ➡ the kernel places an upcall to the user-level library, effectively giving it the processor
  - ➡ the user-level library code assigns this processor to the new thread

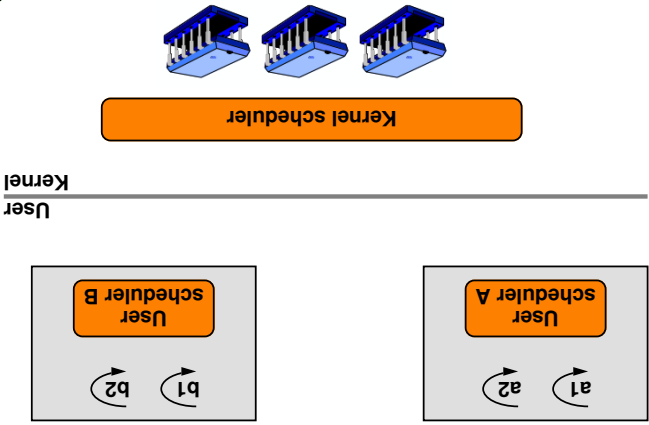
## Schedulier Activations Model

- The scheduler activations model is radically different from the other models
  - in other models, we think of the kernel as providing some kernel thread contexts
    - then multiplexing these contexts on processors using the kernel's scheduler
  - in scheduler activations model, we divvy up processors to processes, and processes determine which threads get to use these processors
    - the kernel should supply however many kernel contexts it finds necessary

## Deadlock

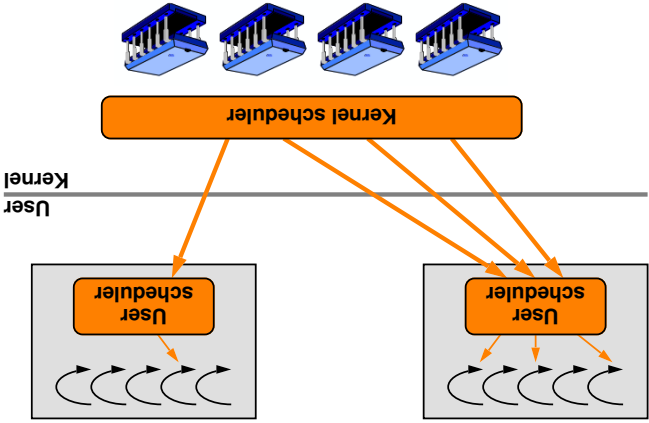


## Scheduler Activations Model Example



**Kernel scheduler does not schedule threads**

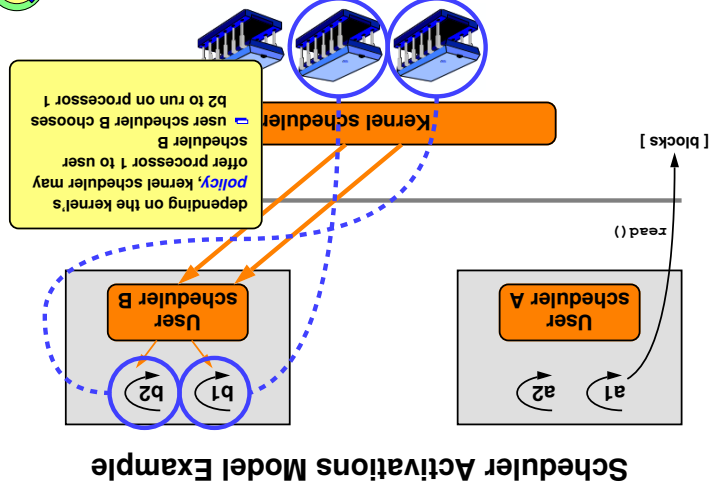
## Schedulier Activations Model Example



## Recap - Problems

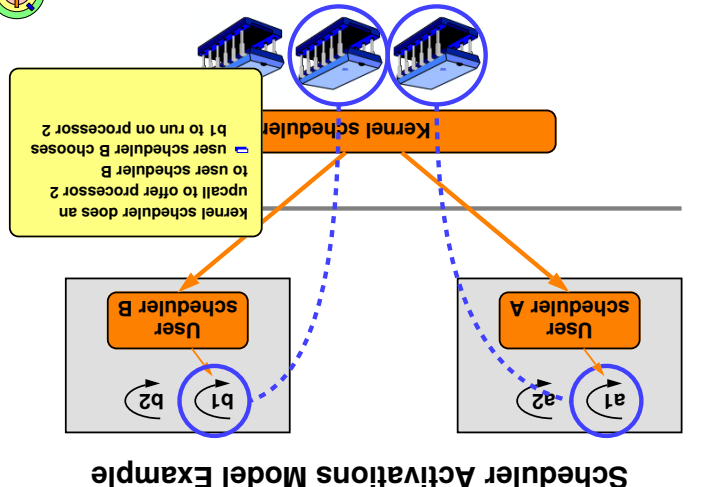
- ➡ Two-level model does not solve the I/O blocking problem
  - ➡ if there are  $N$  kernel threads and if  $N$  user threads are blocked in I/O
    - no other user threads can make progress
- ➡ Another problem: *Priority Inversion*
  - ➡ user-level thread schedulers are not aware of the kernel-level thread scheduler
    - it may know the number of kernel threads how can the user-level scheduler talk to the kernel-level scheduler?
    - people have tried this, but it's complicated
      - ➡ it's possible to have a higher priority user thread scheduled on a lower priority kernel thread and vice versa

Kernel scheduler can have various scheduling policies



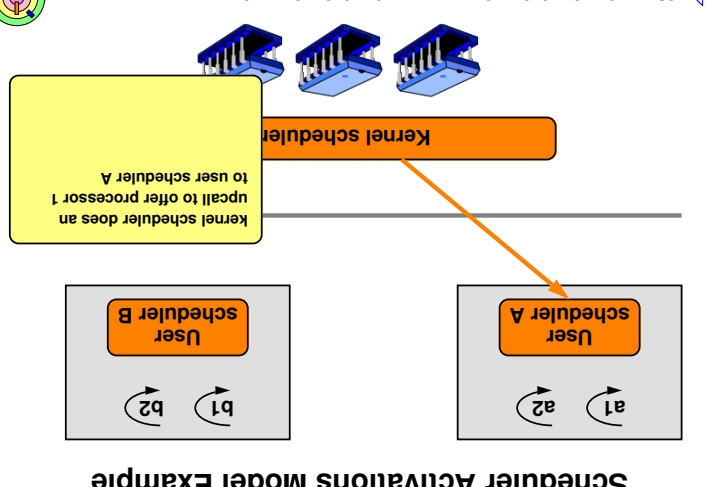
## Scheduler Activations Model Example

Kernel scheduler does not schedule threads



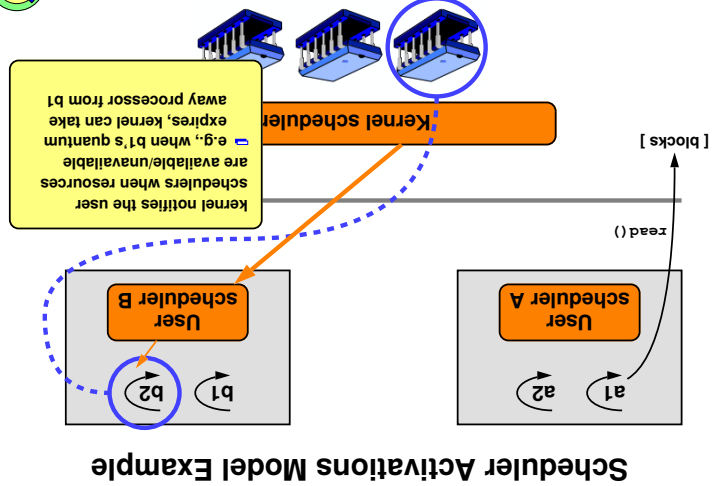
## Schedulier Activations Model Example

Kernel scheduler does not schedule threads



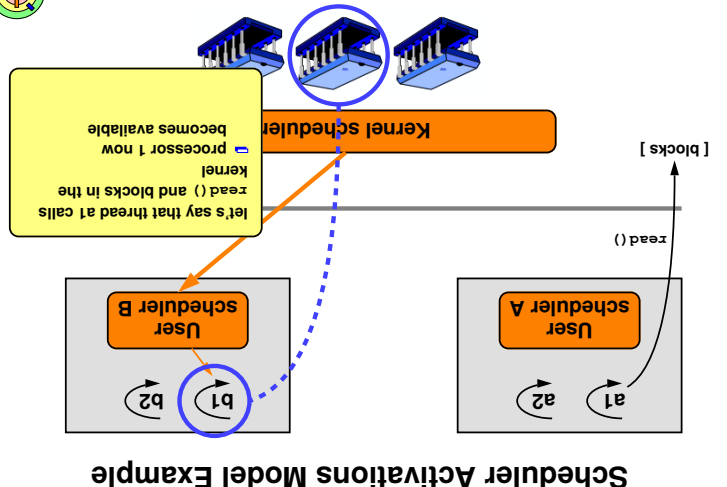
### Scheduier Activations Model Example

Kernel scheduler can have various scheduling policies



## Scheduler Activations Model Example

Kernel scheduler can have various scheduling policies



## Scheduler Activations Model Example

# 5.1 Threads Implementations

- Strategies
- *A Simple Thread Implementation*
- Multiple Processors



Copyright © William C. Cheng



- Threads implementation considerations
- data structures
- thread switching
- synchronization
- how to implement mutexes?
- spin locks
- sleep/blocking locks
- mutexes
- please keep in mind that a mutex *can be* implemented in the kernel *and* in the user space

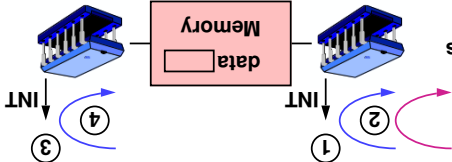
## A Simple Threads Implementation



Copyright © William C. Cheng



- The challenge with implementing mutexes is that you have to ensure that they perform correctly under different kinds of concurrency
- Asynchronous activities that may require concurrency control



Copyright © William C. Cheng



- 1) an *interrupt handler* running on the *same processor* that accesses the same data structure
- 2) *another thread* running on the *same processor* may *preempt* this thread and accesses the same data structure
- 3) an *interrupt handler* running on *another processor* might access the same data structure
- 4) *another thread* running on *another processor* might access the same data structure

## A Simple Threads Implementation

- This implementation is the basis for user-level threads package
- "thread" can mean kernel thread or user thread
- mutex does not need to be a kernel data structure
- *Straight-threads* implementation
- everything happens in thread contexts
- *no interrupt*
- therefore, *no preemption*
- *one processor*
- this is like your kernel 1 with DRIVERS=0 in Config.mk

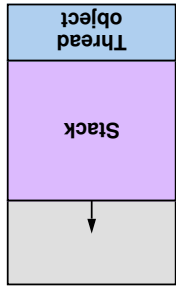
## A Simple Threads Implementation



Copyright © William C. Cheng



## Basic Representation



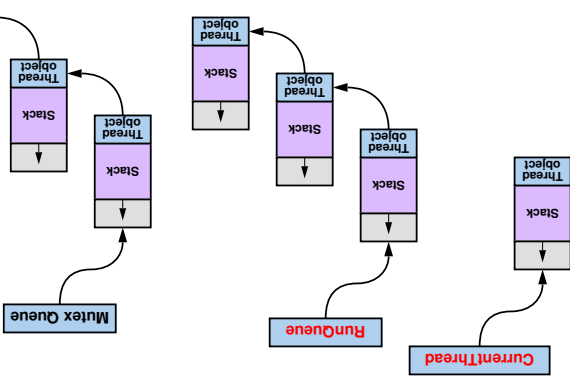
- We will depict a thread like this (to be more compact) although we know that a thread control block is separated from a thread's stack



Copyright © William C. Cheng



## A Collection of Threads



- Each thread must be in one of these data structures
- your kernel assignment looks like this
- *at any time*, you should know where your threads are



Copyright © William C. Cheng



41

Copyright © William C. Cheng

### Straight-threads - Synchronization

```

void mutex_lock(mutex_t *m) {
    if (m->locked) {
        enqueue(m->queue, CurrentThread);
        thread_switch();
    } else {
        m->locked = 1;
    }
}

void mutex_unlock(mutex_t *m) {
    if (queue_empty(m->queue)) {
        m->locked = 0;
    } else {
        dequeue(runqueue, dequeue(m->queue));
    }
}

```

Why is the code atomic?

- single process and no interrupts
- no way to preempt a thread's execution
- a thread holds on to the processor as long as it wants, until it relinquishes processor all by itself

Operating Systems - CSCI 402

39

Copyright © William C. Cheng

### Straight-threads - Synchronization

According to the textbook

```

void mutex_lock(mutex_t *m) {
    if (m->locked) {
        enqueue(m->queue, CurrentThread);
        thread_switch();
    } else {
        m->locked = 1;
    }
}

void mutex_unlock(mutex_t *m) {
    if (queue_empty(m->queue)) {
        m->locked = 0;
    } else {
        dequeue(runqueue, dequeue(m->queue));
    }
}

```

mutex\_unlock() does not seem to work because when it returns, the mutex can be locked and the new mutex holder is not holding the mutex

- after further analysis, it actually does work!

Operating Systems - CSCI 402

37

Copyright © William C. Cheng

### Context Pointer

Recall from Ch 3

if this thread is not currently running, "stack frame" corresponds to switch()

Operating Systems - CSCI 402

42

Copyright © William C. Cheng

### 5.1 Threads Implementations

- Strategies
- A Simple Thread Implementation
- Multiple Processors

Operating Systems - CSCI 402

40

Copyright © William C. Cheng

### Straight-threads - Synchronization

Why is the code atomic?

```

void mutex_lock(mutex_t *m) {
    if (m->locked) {
        enqueue(m->queue, CurrentThread);
        thread_switch();
    } else {
        m->locked = 1;
    }
}

void mutex_unlock(mutex_t *m) {
    if (queue_empty(m->queue)) {
        m->locked = 0;
    } else {
        dequeue(runqueue, dequeue(m->queue));
    }
}

```

Operating Systems - CSCI 402

38

Copyright © William C. Cheng

### Straight-threads - Thread Switch

Need a thread\_switch() function to yield the processor

```

void thread_switch() {
    thread_t NextThread, OldCurrent;
    NextThread = dequeue(Runqueue);
    OldCurrent = CurrentThread;
    CurrentThread = NextThread;
    swapcontext(&OldCurrent->context,
                &NextThread->context);
    // We're now in the new thread's context
}

```

- switch() in Ch 3 has a target thread argument
- swapcontext(old, new) **saves** the caller's context into the old context and **restores** from the new context
- note that the RunQueue may be empty, so this code is incomplete
- before you get here, the current thread is queued onto somewhere else already (e.g., a mutex queue)

Operating Systems - CSCI 402

Copyright © William C. Cheng

47

LOCK  
WR  
RD  
D[0..31]  
A[0..31] lock

○ e.g., assume mutex is *unlocked*, call CAS (lock, 0, 1)  
○ mutex is represented as a bit, 0 if unlocked, 1 if locked

```

int CAS(int *ptr, int old, int new) {
    int tmp = *ptr;
    if (tmp == old) // get the value of mutex
        *ptr = *tmp; // set it to new
    return tmp;
}

```

Compare and swap *machine instruction*

### Hardware Support

Operating Systems - CSCI 402

Copyright © William C. Cheng

45

LOCK  
WR  
RD  
D[0..31]  
A[0..31] lock

○ often implemented as a machine-level instruction  
○ must execute atomically  
◇ how?

```

int CAS(int *ptr, int old, int new) {
    int tmp = *ptr;
    if (tmp == old) // get the value of mutex
        *ptr = *tmp; // set it to new
    return tmp;
}

```

Compare and swap *machine instruction*

### Hardware Support

Operating Systems - CSCI 402

Copyright © William C. Cheng

43

```

void thread_switch() {
    thread_t NextThread, OldCurrent;
    NextThread = dequeue(RunQueue);
    OldCurrent = CurrentThread;
    CurrentThread = NextThread;
    swapcontext(&OldCurrent->context,
                &NextThread->context);
}
}
}
enqueue(RunQueue, CurrentThread);
thread_switch();
}
}

```

○ code is incomplete (because thread\_switch() is incomplete, the way it was presented here)  
○ this thread never blocks, so there is always something to run to avoid boundary condition  
= normal threads join the RunQueue when ready

longer sufficient  
= it's meant for uniprocessor  
Simple approach  
= run on each *processor* an *idle thread*

thread\_switch() is no

### Straight-threads - Multiple Processors

Operating Systems - CSCI 402

Copyright © William C. Cheng

46

LOCK  
WR  
RD  
D[0..31]  
A[0..31] lock

○ e.g., assume mutex is *unlocked*, call CAS (lock, 0, 1)  
○ mutex is represented as a bit, 0 if unlocked, 1 if locked

```

int CAS(int *ptr, int old, int new) {
    int tmp = *ptr;
    if (tmp == old) // get the value of mutex
        *ptr = *tmp; // set it to new
    return tmp;
}

```

Compare and swap *machine instruction*

### Hardware Support

Operating Systems - CSCI 402

Copyright © William C. Cheng

46

LOCK  
WR  
RD  
D[0..31]  
A[0..31] lock

○ e.g., assume mutex is *unlocked*, call CAS (lock, 0, 1)  
○ mutex is represented as a bit, 0 if unlocked, 1 if locked

```

int CAS(int *ptr, int old, int new) {
    int tmp = *ptr;
    if (tmp == old) // get the value of mutex
        *ptr = *tmp; // set it to new
    return tmp;
}

```

Compare and swap *machine instruction*

### Hardware Support

Operating Systems - CSCI 402

Copyright © William C. Cheng

44

Memory

When there are multiple processors, the difficulty lies in locking

```

if (i == -1) {
    m->locked = 1;
}

```

if both threads execute the above code concurrently, in different processors, both threads think they got the lock

### Straight-threads - Multiple Processors

Operating Systems - CSCI 402



Copyright © William C. Cheng

## Blocking Locks

- Spin locks are wasteful
- processor time wasted waiting for the lock to be released
- barely acceptable if locks are held only briefly
- A better approach is to have a blocking lock
- threads wait by having their execution suspended
- a thread much yield the processor and join a queue of waiting threads
- later on, get resumed explicitly

53

Copyright © William C. Cheng

## Hardware Support

Compare and swap *machine instruction*

```
int CAS(int *ptr, int old, int new) {
    int tmp = *ptr;
    if (tmp == old) // get the value of mutex
        *ptr = new; // if it equals to old
    return tmp; // set it to new
}
```

- often implemented as a machine-level instruction
- must execute atomically
- Spin lock
  - mutex is represented as a bit, 0 if unlocked, 1 if locked

51

Copyright © William C. Cheng

## Hardware Support

Compare and swap *machine instruction*

```
int CAS(int *ptr, int old, int new) {
    int tmp = *ptr;
    if (tmp == old) // get the value of mutex
        *ptr = new; // if it equals to old
    return tmp; // set it to new
}
```

→

→ e.g., assume mutex is **unlocked**, call CAS(&lock, 0, 1)

→ e.g., assume mutex is **locked**, call CAS(&lock, 0, 1)

49

Copyright © William C. Cheng

## Blocking Locks

```
void blocking_lock(mutex_t *m) {
    if (m->holder != 0) {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    } else {
        m->holder = CurrentThread;
    }
}

void blocking_unlock(mutex_t *m) {
    if (!queue_empty(m->wait_queue)) {
        m->holder = 0;
        dequeue(RunQueue, m->holder);
    }
}
```

→ This code only works on a uniprocessor

54

Copyright © William C. Cheng

## Spin Lock

Naive spin lock

```
void spin_lock(int *mutex) {
    while (*mutex == 1) {}
}

void spin_unlock(int *mutex) {
    *mutex = 0;
}
```

Better spin lock

```
void spin_lock(int *mutex) {
    while (*mutex == 1) {}
    if (*mutex == 0) {
        // the mutex was at least momentarily unlocked
        if (!CAS(mutex, 0, 1))
            break; // we have locked the mutex
        // some other thread beat us to it, try again
    }
}
```

→ textbook is wrong

52

Copyright © William C. Cheng

## Hardware Support

Compare and swap *machine instruction*

```
int CAS(int *ptr, int old, int new) {
    int tmp = *ptr;
    if (tmp == old) // get the value of mutex
        *ptr = new; // if it equals to old
    return tmp; // set it to new
}
```

→ e.g., assume mutex is **unlocked**, call CAS(&lock, 0, 1)

→ e.g., assume mutex is **locked**, call CAS(&lock, 0, 1)

50

Copyright © William C. Cheng

59

Thread 2 can move thread 1 to another processor! (Can it?)

```

void blocking_unlock(mutex_t *m)
{
    if (m->holder != 0)
    {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    }
    else
    {
        m->holder = CurrentThread;
        spin_unlock(m->spinlock);
    }
}

void blocking_lock(mutex_t *m)
{
    if (queue_empty(m->wait_queue))
    {
        m->holder = 0;
        spin_lock(m->spinlock);
    }
    else
    {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    }
}

enqueue(RunQueue, m->holder);
m->holder = dequeue(m->wait_queue);
spin_unlock(m->spinlock);
}

```

Working Blocking Locks (?)

Copyright © William C. Cheng

60

Can you do spin\_unlock() inside thread\_switch()?

```

void blocking_unlock(mutex_t *m)
{
    if (queue_empty(m->wait_queue))
    {
        m->holder = 0;
        spin_lock(m->spinlock);
    }
    else
    {
        enqueue(RunQueue, m->holder);
        spin_unlock(m->spinlock);
    }
}

void blocking_lock(mutex_t *m)
{
    if (m->holder != 0)
    {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    }
    else
    {
        m->holder = CurrentThread;
        spin_unlock(m->spinlock);
    }
}

```

Working Blocking Locks (?)

Copyright © William C. Cheng

57

Will deadlock because of thread\_switch()

```

void blocking_unlock(mutex_t *m)
{
    if (queue_empty(m->wait_queue))
    {
        m->holder = 0;
        spin_lock(m->spinlock);
    }
    else
    {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    }
}

void blocking_lock(mutex_t *m)
{
    if (m->holder != 0)
    {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    }
    else
    {
        m->holder = CurrentThread;
        spin_unlock(m->spinlock);
    }
}

```

Working Blocking Locks (?)

Copyright © William C. Cheng

58

Has a different problem

```

void blocking_unlock(mutex_t *m)
{
    if (queue_empty(m->wait_queue))
    {
        m->holder = 0;
        spin_lock(m->spinlock);
    }
    else
    {
        enqueue(RunQueue, m->holder);
        spin_unlock(m->spinlock);
    }
}

void blocking_lock(mutex_t *m)
{
    if (m->holder != 0)
    {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    }
    else
    {
        m->holder = CurrentThread;
        spin_unlock(m->spinlock);
    }
}

```

Working Blocking Locks (?)

Copyright © William C. Cheng

55

On a multiprocessor, it may not work

```

void blocking_unlock(mutex_t *m)
{
    if (queue_empty(m->wait_queue))
    {
        m->holder = 0;
        spin_lock(m->spinlock);
    }
    else
    {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    }
}

void blocking_lock(mutex_t *m)
{
    if (m->holder != 0)
    {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    }
    else
    {
        m->holder = CurrentThread;
        spin_unlock(m->spinlock);
    }
}

```

Blocking Locks

Copyright © William C. Cheng

56

On a multiprocessor, it may not work

```

void blocking_unlock(mutex_t *m)
{
    if (queue_empty(m->wait_queue))
    {
        m->holder = 0;
        spin_lock(m->spinlock);
    }
    else
    {
        enqueue(RunQueue, m->holder);
        spin_unlock(m->spinlock);
    }
}

void blocking_lock(mutex_t *m)
{
    if (m->holder != 0)
    {
        enqueue(m->wait_queue, CurrentThread);
        thread_switch();
    }
    else
    {
        m->holder = CurrentThread;
        spin_unlock(m->spinlock);
    }
}

```

Blocking Locks

Copyright © William C. Cheng

## Thread Synchronization Summary

- Spin locks
  - used if the duration of waiting is expected to be small
  - as in the case at the beginning of `blocking_lock()`
- Sleep (or blocking) locks
  - used if the duration of waiting is expected to be long
- Futexes
  - optimized version of blocking locks
- In your kernel assignment #1, you need to implement *kernel* threads
  - very different from *user* threads
  - keep in mind that the *wee* kernel is *non-preemptive*
  - the kernel is very powerful (and therefore, must be bug free)
  - in kernel assignment #3, you need to implement user threads/processes (well, still one thread per process)

65

Copyright © William C. Cheng

## Attempt 1

```

futex->val
= 0 means unlocked; otherwise, locked

void lock(futex_t *futex)
{
    unsigned int c;
    while ((c = atomic_inc(&futex->val)) != 0)
        futex_wait(futex, c+1);
}

void unlock(futex_t *futex)
{
    futex->val = 0;
    futex_wake(futex);
}

```

- Problem with `unlock()`
  - slow because `futex_wake()` is a system call
- Problem with `lock()`
  - threads run in lock steps in a multiprocessor environment!
  - `futex->val` may wrap-around

63

Copyright © William C. Cheng

## Futexes

- Futex**: fast user-space mutex
  - safe, efficient kernel conditional queueing in Linux
  - most of the time when you try to lock a mutex, it's unlocked; so just go ahead and lock it (no system call)
  - if it's locked (by another thread), then a system call is required for this thread to obtain the lock
  - contained in it is an unsigned integer state called `val` and a queue of waiting threads
- Two **system calls** are provided to support futexes
  - `futex_wait(futex_t *futex, int val)`
    - if (`futex->val == val`)
      - `sleep()`;
  - `futex_wake(futex_t *futex)`
    - wake up one thread from wait queue if there is any

61

Copyright © William C. Cheng

## Attempt 2

```

futex->val can only take on values of 0, 1, and 2
= 0 means unlocked
= 1 means locked but no waiting thread
= 2 means locked with the possibility of waiting threads

void lock(futex_t *futex)
{
    unsigned int c;
    if ((c = CAS(&futex->val, 0, 1)) != 0)
        do {
            if (c == 2 || (CAS(&futex->val, 1, 2) == 1))
                futex_wait(futex, 2);
        } while ((c = CAS(&futex->val, 0, 2)) != 0);
}

void unlock(futex_t *futex)
{
    if (atomic_dec(&futex->val) != 1)
        futex_wake(futex);
}

```

textbook is wrong

64

Copyright © William C. Cheng

## Ancillary Functions

- Add 1 to `*val`, return its original value
 

```

      unsigned int atomic_inc(unsigned int *val)
      {
          // performed atomically
          return ((*val)++); // textbook is wrong
      }
      
```
- Subtract 1 to `*val`, return its original value
 

```

      unsigned int atomic_dec(unsigned int *val)
      {
          // performed atomically
          return ((*val)--); // textbook is wrong
      }
      
```
- Just like `CAS()`, both functions return the *previous* lock value

62