# semaphores / reader/writer

# Changelog

Changes made in this version not seen in first lecture:

1 October 2019: fix mixup of 'result' and 'value' in semaphore exercise return

3 October 2019: correct reader-priority rwlock code to include readers $== 0$ check before signaling in ReadUnlock

# last time

monitors = mutex + condition variable

mutex protects shared data
    important: locked mutex = whether thread should wait wont' change

condition variable (CV): abstracts queue of waiting threads

CV wait: unlock a mutex + start waiting on queue
    done simultaneously so thread doesn't miss its signal to wake up
    *spurious wakeups* — need to double-check condition

CV broadcast: remove all threads from CV queue, have them reacquire lock

CV signal: remove one threads from CV queue, have it reacquire lock
    no guarantee that it reacquire lock first (except rare Hoare-style monitors)

# monitor exercise (1)

suppose we want producer/consumer, but...

but change to ConsumeTwo() which returns a pair of values
    and don't want two calls to ConsumeTwo() to wait...
    with each getting one item

what should we change below?

```
pthread_mutex_t lock;                      Consume() {
pthread_cond_t data_ready;                   pthread_mutex_lock(&lock);
UnboundedQueue buffer;                       while (buffer.empty()) {
                                               pthread_cond_wait(&data_ready, &lock)
Produce(item) {                              }
  pthread_mutex_lock(&lock);                 item = buffer.dequeue();
  buffer.enqueue(item);                      pthread_mutex_unlock(&lock);
  pthread_cond_signal(&data_ready);          return item;
  pthread_mutex_unlock(&lock);             }
}
```

# monitor exercise: solution (1)

(one of many possible solutions)
Assuming ConsumeTwo **replaces** Consume:

```
Produce() {
  pthread_mutex_lock(&lock);
  buffer.enqueue(item);
  if (buffer.size() > 1) { pthread_cond_signal(&data_ready); }
  pthread_mutex_unlock(&lock);
}
ConsumeTwo() {
    pthread_mutex_lock(&lock);
    while (buffer.size() < 2) { pthread_cond_wait(&data_ready, &lock); }
    item1 = buffer.dequeue(); item2 = buffer.dequeue();
    pthread_mutex_unlock(&lock);
    return Combine(item1, item2);
}
```

# monitor exercise: solution 2

(one of many possible solutions)

Assuming ConsumeTwo is **in addition to** Consume (using two CVs):

```
Produce() {
  pthread_mutex_lock(&lock);
  buffer.enqueue(item);
  pthread_cond_signal(&one_ready);
  if (buffer.size() > 1) { pthread_cond_signal(&two_ready); }
  pthread_mutex_unlock(&lock);
}
Consume() {
  pthread_mutex_lock(&lock);
  while (buffer.size() < 1) { pthread_cond_wait(&one_ready, &lock); }
  item = buffer.dequeue();
  pthread_mutex_unlock(&lock);
  return item;
}
ConsumeTwo() {
  pthread_mutex_lock(&lock);
  while (buffer.size() < 2) { pthread_cond_wait(&two_ready, &lock); }
  item1 = buffer.dequeue(); item2 = buffer.dequeue();
  pthread_mutex_unlock(&lock);
  return Combine(item1, item2);
}
```

# monitor exercise: slow solution

(one of many possible solutions)

Assuming ConsumeTwo is **in addition to** Consume (using one CV):

```
Produce() {
  pthread_mutex_lock(&lock);
  buffer.enqueue(item);
  // broadcast and not signal, b/c we might wakeup only ConsumeTwo() otherwise
  pthread_cond_broadcast(&data_ready);
  pthread_mutex_unlock(&lock);
}
Consume() {
  pthread_mutex_lock(&lock);
  while (buffer.size() < 1) { pthread_cond_wait(&data_ready, &lock); }
  item = buffer.dequeue();
  pthread_mutex_unlock(&lock);
  return item;
}
ConsumeTwo() {
  pthread_mutex_lock(&lock);
  while (buffer.size() < 2) { pthread_cond_wait(&data_ready, &lock); }
  item1 = buffer.dequeue(); item2 = buffer.dequeue();
  pthread_mutex_unlock(&lock);
  return Combine(item1, item2);
}
```

# monitor exercise (2)

suppose we want to implement a one-use barrier

what goes in the blanks?
```
struct BarrierInfo {
    pthread_mutex_t lock;
    int total_threads;  // initially total # of threads
    int number_reached; // initially 0
    _____
};

void BarrierWait(BarrierInfo *barrier) {
    pthread_mutex_lock(&barrier->lock);
    ++number_reached;

    _____
    _____
    _____
    pthread_mutex_unlock(&barrier->lock);
}
```

# mutex/cond var init/destroy

```
pthread_mutex_t mutex;
pthread_cond_t cv;
pthread_mutex_init(&mutex, NULL);
pthread_cond_init(&cv, NULL);
// --OR--
pthread_mutex_t mutex = PTHREAD_MUTEX_INITIALIZER;
pthread_cond_t cv = PTHREAD_COND_INITIALIZER;

// and when done:
...
pthread_cond_destroy(&cv);
pthread_mutex_destroy(&mutex);
```

# generalizing locks: semaphores

semaphore has a non-negative integer **value** and two operations:

**P()** or **down** or **wait**:
wait for semaphore to become positive ($> 0$),
then decerement by 1

**V()** or **up** or **signal** or **post**:
increment semaphore by 1 (waking up thread if needed)

P, V from Dutch: *proberen* (test), *verhogen* (increment)

# semaphores are kinda integers

semaphore like an integer, but...

cannot read/write directly
>     down/up operaion only way to access (typically)
>     exception: initialization

never negative — wait instead
>     down operation wants to make negative? thread waits

# reserving books

suppose tracking copies of library book…

```
Semaphore free_copies = Semaphore(3);
void ReserveBook() {
    // wait for copy to be free
    free_copies.down();
    ... // ... then take reserved copy
}

void ReturnBook() {
    ... // return reserved copy
    free_copies.up();
    // ... then wakeup waiting thread
}
```

# counting resources: reserving books

suppose tracking copies of same library book
non-negative integer count = # how many books used?
up = give back book; down = take book

| Copy 1 |
|--------|
| Copy 2 |
| Copy 3 |

free copies  3

# counting resources: reserving books

suppose tracking copies of same library book
non-negative integer count = # how many books used?
up = give back book; down = take book



taken out

| Copy 1 |
| Copy 2 |
| Copy 3 |

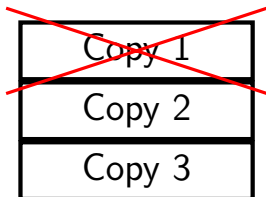free copies

after calling down to reserve

# counting resources: reserving books

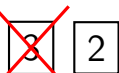suppose tracking copies of same library book
non-negative integer count = # how many books used?
up = give back book; down = take book

taken out



free copies  2

after calling down to reserve

# counting resources: reserving books

suppose tracking copies of same library book
non-negative integer count = # how many books used?
up = give back book; down = take book



taken out · Copy 1
taken out · Copy 2
taken out · Copy 3

free copies  $\boxed{0}$

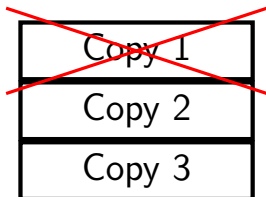after calling down three times
to reserve all copies

# counting resources: reserving books

suppose tracking copies of same library book
non-negative integer count = # how many books used?
up = give back book; down = take book



taken out    Copy 1
taken out    Copy 2
taken out    Copy 3

free copies  0

**reserve book**
call *down* again
start waiting…

# counting resources: reserving books

suppose tracking copies of same library book
non-negative integer count = # how many books used?
up = give back book; down = take book



taken out — Copy 1
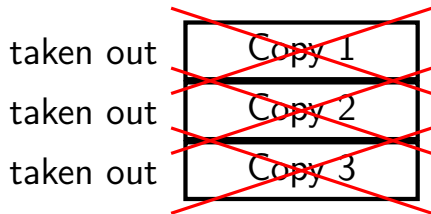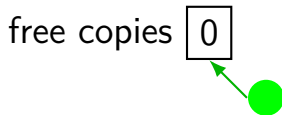taken out — Copy 2
taken out — Copy 3

free copies | 0 |

**return book**

call *up*
release waiter

**reserve book**

call *down*
~~waiting~~
done waiting

# implementing mutexes with semaphores

```
struct Mutex {
    Semaphore s; /* with inital value 1 */
    /* value = 1 --> mutex if free */
    /* value = 0 --> mutex is busy */
}

MutexLock(Mutex *m) {
    m->s.down();
}

MutexUnlock(Mutex *m) {
    m->s.up();
}
```

# implementing join with semaphores

```
struct Thread {
    ...
    Semaphore finish_semaphore; /* with initial value 0 */
    /* value = 0: either thread not finished OR already joined */
    /* value = 1: thread finished AND not joined */
};
thread_join(Thread *t) {
    t->finish_semaphore->down();
}

/* assume called when thread finishes */
thread_exit(Thread *t) {
    t->finish_semaphore->up();
    /* tricky part: deallocating struct Thread safely? */
}
```

# POSIX semaphores

```c
#include <semaphore.h>
...
sem_t my_semaphore;
int process_shared = /* 1 if sharing between processes */;
sem_init(&my_semaphore, process_shared, initial_value);
...
sem_wait(&my_semaphore);  /* down */
sem_post(&my_semaphore);  /* up */
...
sem_destroy(&my_semaphore);
```

## semaphore exercise

```
int value;  sem_t empty, ready;

void PutValue(int argument) {
    sem_wait(&empty);
    value = argument;
    sem_post(&ready);
}

int GetValue() {
    int result;
    _____
    result = value;
    _____
    return result;
}
```

A: sem_post(&empty) / sem_wait(&ready)
B: sem_wait(&ready) / sem_post(&empty)
C: sem_post(&ready) / sem_wait(&empty)
D: sem_post(&ready) / sem_post(&empty)
E: sem_wait(&empty) / sem_post(&ready)
F: something else

GetValue() waits for PutValue() to happen, then reutrns value, allows next PutValue() to happen. What goes in blanks?

# semaphore exercise [solution]

```
int value;
sem_t empty, ready;
void PutValue(int argument) {
    sem_wait(&empty);
    value = argument;
    sem_post(&ready);
}
int GetValue() {
    int result;
    sem_wait(&ready);
    result = value;
    sem_post(&empty);
    return result;
}
```

# semaphore intuition

What do you need to wait for?
> critical section to be finished
> queue to be non-empty
> array to have space for new items

what can you count that will be 0 when you need to wait?
> # of threads that can start critical section now
> # of threads that can join another thread without waiting
> # of items in queue
> # of empty spaces in array

use up/down operations to maintain count

# producer/consumer constraints

consumer waits for producer(s) if buffer is empty

producer waits for consumer(s) if buffer is full

any thread waits while a thread is manipulating the buffer

# producer/consumer constraints

consumer waits for producer(s) if buffer is empty

producer waits for consumer(s) if buffer is full

any thread waits while a thread is manipulating the buffer

one semaphore per constraint:

```
sem_t full_slots;    // consumer waits if empty
sem_t empty_slots;   // producer waits if full
sem_t mutex;         // either waits if anyone changing buffer
FixedSizedQueue buffer;
```

# producer/consumer pseudocode

```
sem_init(&full_slots, ..., 0 /* # buffer slots initially used */);
sem_init(&empty_slots, ..., BUFFER_CAPACITY);
sem_init(&mutex, ..., 1 /* # thread that can use buffer at once */);
buffer.set_size(BUFFER_CAPACITY);
...

Produce(item) {
    sem_wait(&empty_slots);  // wait until free slot, reserve it
    sem_wait(&mutex);
    buffer.enqueue(item);
    sem_post(&mutex);
    sem_post(&full_slots);  // tell consumers there is more data
}

Consume() {
    sem_wait(&full_slots);  // wait until queued item, reserve it
    sem_wait(&mutex);
    item = buffer.dequeue();
    sem_post(&mutex);
    sem_post(&empty_slots);  // let producer reuse item slot
    return item;
}
```

# producer/consumer pseudocode

```
sem_init(&full_slots, ..., 0 /* # buffer slots initially used */);
sem_init(&empty_slots, ..., BUFFER_CAPACITY);
sem_init(&mutex, ..., 1 /* # thread that can use buffer at once */);
buffer.set_size(BUFFER_CAPACITY);
...

Produce(item) {
    sem_wait(&empty_slots);  // wait until free slot, reserve it
    sem_wait(&mutex);
    buffer.enqueue(item);
    sem_post(&mutex);
    sem_post(&full_slots);  // tell consumers there is more data
}

Consume() {
    sem_wait(&full_slots);  // wait until queued item, reserve it
    sem_wait(&mutex);
    item = buffer.dequeue();
    sem_post(&mutex);
    sem_post(&empty_slots);  // let producer reuse item slot
    return item;
}
```

# producer/consumer pseudocode

```
sem_init(&full_slots, ..., 0 /* # buffer slots initially used */);
sem_init(&empty_slots, ..., BUFFER_CAPACITY);
sem_init(&mutex, ..., 1 /* # thread that can use buffer at once */);
buffer.set_size(BUFFER_CAPACITY);
...

Produce(item) {
    sem_wait(&empty_slots);  // wait until free slot, reserve it
    sem_wait(&mutex);
    buffer.enqueue(item);
    sem_post(&mutex);
    sem_post(&full_slots);   // tell consumers there is more data
}

Consume() {
    sem_wait(&full_slots);   // wait until queued item, reserve it
    sem_wait(&mutex);
    item = buffer.dequeue();
    sem_post(&mutex);
    sem_post(&empty_slots);  // let producer reuse item slot
    return item;
}
```

# producer/consumer pseudocode

```
sem_init(&full_slots, ..., 0 /* # buffer slots initially used */);
sem_init(&empty_slots, ..., BUFFER_CAPACITY);
sem_init(&mutex, ..., 1 /* # thread that can use buffer at once */);
buffer.set_size(BUFFER_CAPACITY);
...

Produce(item) {
    sem_wait(&empty_slots);   // wait until free slot, reserve it
    sem_wait(&mutex);
    buffer.enqueue(item);
    sem_post(&mutex);
    sem_post(&full_slots);                          re data
}

Consume() {
    sem_wait(&full_slots);   // wait until queued item, reserve it
    sem_wait(&mutex);
    item = buffer.dequeue();
    sem_post(&mutex);
    sem_post(&empty_slots);  // let producer reuse item slot
    return item;
}
```

Can we do
  sem_wait(&mutex);
  sem_wait(&empty_slots);
instead?

# producer/consumer pseudocode

```
sem_init(&full_slots, ..., 0 /* # buffer slots initially used */);
sem_init(&empty_slots, ..., BUFFER_CAPACITY);
sem_init(&mutex, ..., 1 /* # thread that can use buffer at once */);
buffer.set_size(BUFFER_CAPACITY);
...

Produce(item) {
    sem_wait(&empty_slots);   // wait until free slot, reserve it
    sem_wait(&mutex);
    buffer.enqueue(item);
    sem_post(&mutex);
    sem_post(&full_slots);                          re data
}

Consume() {
    sem_wait(&full_slots);
    sem_wait(&mutex);
    item = buffer.dequeue()
    sem_post(&mutex);
    sem_post(&empty_slots);
    return item;
}
```

Can we do
    sem_wait(&mutex);
    sem_wait(&empty_slots);
instead?

No. Consumer waits on `sem_wait(&mutex)`
so can't `sem_post(&empty_slots)`
(result: producer waits forever
problem called *deadlock*)

# producer/consumer: cannot reorder mutex/empty

```
ProducerReordered() {
  // BROKEN: WRONG ORDER
  sem_wait(&mutex);
  sem_wait(&empty_slots);

  ...

  sem_post(&mutex);
```

```
Consumer() {
  sem_wait(&full_slots);

  // can't finish until
  // Producer's sem_post(&mutex):
  sem_wait(&mutex);

  ...

  // so this is not reached
  sem_post(&full_slots);
```

# producer/consumer pseudocode

```
sem_init(&full_slots, ..., 0 /* # buffer slots initially used */);
sem_init(&empty_slots, ..., BUFFER_CAPACITY);
sem_init(&mutex, ..., 1 /* # thread that can use buffer at once */);
buffer.set_size(BUFFER_CAPACITY);
...

Produce(item) {
    sem_wait(&empty_slots);  // wait until free slot, reserve it
    sem_wait(&mutex);
    buffer.enqueue(item);
    sem_post(&mutex);
    sem_post(&full_slots                              s more data
}

Consume() {
    sem_wait(&full_slots                              , reserve it
    sem_wait(&mutex);
    item = buffer.dequeu
    sem_post(&mutex);
    sem_post(&empty_slots);  // let producer reuse item slot
    return item;
}
```

Can we do
   sem_post(&full_slots);
   sem_post(&mutex);
instead?
Yes — post never waits

# producer/consumer summary

producer: wait (down) empty_slots, post (up) full_slots

consumer: wait (down) full_slots, post (up) empty_slots

two producers or consumers?
    still works!

# binary semaphores

*binary semaphores* — semaphores that are <span style="color:red">only zero or one</span>

as powerful as normal semaphores
> exercise: simulate counting semaphores with binary semaphores (more than one) and an integer

# counting semaphores with binary semaphores

```
// assuming initialValue > 0
BinarySemaphore mutex(1);
int value = initialValue ;
BinarySemaphore gate(1 /* if initialValue >= 1 */);
    /* gate = # threads that can Down() now */
```

```
void Down() {                          void Up() {
  gate.Down();                           mutex.Down();
  // wait, if needed                     value += 1;
  mutex.Down();                          if (value == 1) {
  value -= 1;                              gate.Up();
  if (value > 0) {                         // because down should finish now
    gate.Up();                             // but could not before
    // because next down should finish   }
    // now (but not marked to before)    mutex.Up();
  }                                    }
  mutex.Up();
}
```

# gate intuition/pattern

gate is open (value $= 1$): Down() can proceed

gate is closed (Value $= 0$): Down() waits

# gate intuition/pattern

gate is open (value $= 1$): Down() can proceed

gate is closed (Value $= 0$): Down() waits

common pattern with semaphores:

allow threads one-by-one past 'gate'
> keep gate open forever? thread passing gate allows next in

# Anderson-Dahlin and semaphores

Anderson/Dahlin complains about semaphores

"Our view is that programming with locks and condition variables is superior to programming with semaphores."

argument 1: clearer to have separate constructs for

waiting for condition to be come true, and
allowing only one thread to manipulate a thing at a time

arugment 2: tricky to verify thread calls up exactly once for every down

alternatives allow one to be sloppier (in a sense)

# building semaphore with monitors

```
pthread_mutex_t lock;
```

lock to protect shared state

# building semaphore with monitors

```
pthread_mutex_t lock;
unsigned int count;
```

lock to protect shared state
    shared state: semaphore tracks a count

# building semaphore with monitors

```
pthread_mutex_t lock;
unsigned int count;
/* condition, broadcast when becomes count > 0 */
pthread_cond_t count_is_positive_cv;
```

lock to protect shared state
    shared state: semaphore tracks a count

add cond var for each reason we wait
    semaphore: wait for count to become positive (for down)

# building semaphore with monitors

```
pthread_mutex_t lock;
unsigned int count;
/* condition, broadcast when becomes count > 0 */
pthread_cond_t count_is_positive_cv;
void down() {
    pthread_mutex_lock(&lock);
    while (!(count > 0)) {
        pthread_cond_wait(
            &count_is_positive_cv,
            &lock);
    }
    count -= 1;
    pthread_mutex_unlock(&lock);
}
```

lock to protect shared state
    shared state: semaphore tracks a count

add cond var for each reason we wait
    semaphore: wait for count to become positive (for down)

wait using condvar; broadcast/signal when condition changes

# building semaphore with monitors

```
pthread_mutex_t lock;
unsigned int count;
/* condition, broadcast when becomes count > 0 */
pthread_cond_t count_is_positive_cv;
void down() {
    pthread_mutex_lock(&lock);
    while (!(count > 0)) {
        pthread_cond_wait(
            &count_is_positive_cv,
            &lock);
    }
    count -= 1;
    pthread_mutex_unlock(&lock);
}
```

```
void up() {
    pthread_mutex_lock(&lock);
    count += 1;
    /* count must now be
       positive, and at most
       one thread can go per
       call to Up() */
    pthread_cond_signal(
        &count_is_positive_cv
    );
    pthread_mutex_unlock(&lock);
}
```

lock to protect shared state
    shared state: semaphore tracks a count

add cond var for each reason we wait
    semaphore: wait for count to become positive (for down)

wait using condvar; broadcast/signal when condition changes

# building semaphore with monitors (version B)

```
pthread_mutex_t lock;
unsigned int count;
/* condition, broadcast when becomes count > 0 */
pthread_cond_t count_is_positive_cv;
void down() {
    pthread_mutex_lock(&lock);
    while (!(count > 0)) {
        pthread_cond_wait(
            &count_is_positive_cv,
            &lock);
    }
    count -= 1;
    pthread_mutex_unlock(&lock);
}
```

```
void up() {
    pthread_mutex_lock(&lock);
    count += 1;
    /* condition *just* became true */
    if (count == 1) {
        pthread_cond_broadcast(
            &count_is_positive_cv
        );
    }
    pthread_mutex_unlock(&lock);
}
```

before: signal every time

can check if condition just became true instead?

# building semaphore with monitors (version B)

```
pthread_mutex_t lock;
unsigned int count;
/* condition, broadcast when becomes count > 0 */
pthread_cond_t count_is_positive_cv;
void down() {                              void up() {
    pthread_mutex_lock(&lock);                 pthread_mutex_lock(&lock);
    while (!(count > 0)) {                     count += 1;
        pthread_cond_wait(                     /* condition *just* became true */
            &count_is_positive_cv,             if (count == 1) {
            &lock);                                pthread_cond_broadcast(
    }                                                  &count_is_positive_cv
    count -= 1;                                     );
    pthread_mutex_unlock(&lock);               }
}                                              pthread_mutex_unlock(&lock);
                                           }
```

before: signal every time

can check if condition just became true instead?

but do we really need to broadcast?

# exercise: why broadcast?

```
pthread_mutex_t lock;
unsigned int count;
/* condition, broadcast when becomes count > 0 */
pthread_cond_t count_is_positive_cv;
void down() {                            void up() {
    pthread_mutex_lock(&lock);               pthread_mutex_lock(&lock);
    while (!(count > 0)) {                    count += 1;
        pthread_cond_wait(                   if (count == 1) { /* became > 0 */
            &count_is_positive_cv,               pthread_cond_broadcast(
            &lock);                                  &count_is_positive_cv
    }                                            );
    count -= 1;                              }
    pthread_mutex_unlock(&lock);             pthread_mutex_unlock(&lock);
}                                        }
```

exercise: why can't this be pthread_cond_signal?

hint: think of two threads calling down + two calling up?

brute force: only so many orders they can get the lock in

# broadcast problem

| Thread 1 | Thread 2 | Thread 3 | Thread 4 |
|---|---|---|---|
| Down() | | | |
| lock | | | |
| count == 0? yes | | | |
| unlock/wait | | | |
| | Down() | | |
| | lock | | |
| | count == 0? yes | | |
| | unlock/wait | | |
| | | Up() | |
| | | lock | |
| | | count += 1 (now 1) | Up() |
| stop waiting on CV | | signal | wait for lock |
| wait for lock | | unlock | wait for lock |
| wait for lock | | | lock |
| wait for lock | | | count += 1 (now 2) |
| wait for lock | | | count != 1: don't signal |
| lock | | | unlock |
| count == 0? no | | | |
| count -= 1 (becomes 1) | | | |
| unlock | | | |
| | still waiting??? | | |

# broadcast problem

| Thread 1 | Thread 2 | Thread 3 | Thread 4 |
|---|---|---|---|
| Down() | | | |
| lock | | | |
| count == 0? yes | | | |
| unlock/wait | | | |
| | Down() | | |
| | lock | | |
| | count == 0? yes | | |
| | unlock/wait | | |
| | | Up() | |
| | | lock | |
| | | count += 1 (now 1) | Up() |
| stop waiting on CV | | signal | wait for lock |
| wait for lock | | unlock | wait for lock |
| wait for lock | | | lock |
| wait for lock | | | count += 1 (now 2) |
| wait for lock | | | count != 1: don't signal |
| lock | | | unlock |
| count == 0? no | | | |
| count -= 1 (becomes 1) | | | |
| unlock | | | |
| | still waiting??? | | |

# broadcast problem

| Thread 1 | Thread 2 | Thread 3 | Thread 4 |
|---|---|---|---|
| Down() | | | |
| lock | | | |
| count == 0? yes | | | |
| unlock/wait | | | |
| | Down() | | |
| | lock | | |
| | count == 0? yes | | |
| | unlock/wait | | |
| | | Up() | |
| | | lock | |
| | | count += 1 (now 1) | Up() |
| stop waiting on CV | | signal | wait for lock |
| wait for lock | | unlock | wait for lock |
| wait for lock | | | lock |
| wait for lock | | | count += 1 (now 2) |
| wait for lock | | | count != 1: don't signal |
| lock | | | unlock |
| count == 0? no | | | |
| count -= 1 (becomes 1) | | | |
| unlock | | | |
| | still waiting??? | | |

Mesa-style monitors
signalling doesn't
"hand off" lock

# semaphores with monitors: no condition

```
pthread_mutex_t lock;
unsigned int count;
/* condition, broadcast when becomes count > 0 */
pthread_cond_t count_is_positive_cv;
void down() {                              void up() {
    pthread_mutex_lock(&lock);                 pthread_mutex_lock(&lock);
    while (!(count > 0)) {                      count += 1;
        pthread_cond_wait(                      pthread_cond_signal(
            &count_is_positive_cv,                  &count_is_positive_cv
            &lock);                             );
    }                                           pthread_mutex_unlock(&lock);
    count -= 1;                             }
    pthread_mutex_unlock(&lock);
}
```

same as where we started…

# semaphores with monitors: alt w/ signal

```
pthread_mutex_t lock;
unsigned int count;
/* condition, broadcast when becomes count > 0 */
pthread_cond_t count_is_positive_cv;
void down() {                              void up() {
    pthread_mutex_lock(&lock);                 pthread_mutex_lock(&lock);
    while (!(count > 0)) {                      count += 1;
        pthread_cond_wait(                      if (count == 1) {
            &count_is_positive_cv,                  pthread_cond_signal(
            &lock);                                     &count_is_positive_cv
    }                                               );
    count -= 1;                                 }
    if (count > 0) {                            pthread_mutex_unlock(&lock);
        pthread_cond_signal(                }
            &count_is_positive_cv
        );
    }
    pthread_mutex_unlock(&lock);
}
```

# on signal/broadcast generally

whenever using signal need to ask
what if more than one thread is waiting?

be concerned about "skipping" cases where thread would wake up
    unfortunately, Mesa-style scheduling/spurious wakeups make this harder

# monitors with semaphores: locks

```
sem_t semaphore;  // initial value 1

Lock() {
    sem_wait(&semaphore);
}

Unlock() {
    sem_post(&semaphore);
}
```

# monitors with semaphores: cvs

condition variables are more challenging

start with only wait/signal:

```
sem_t threads_to_wakeup;  // initially 0
Wait(Lock lock) {
    lock.Unlock();
    sem_wait(&threads_to_wakeup);
    lock.Lock();
}
Signal() {
    sem_post(&threads_to_wakeup);
}
```

# monitors with semaphores: cvs

condition variables are more challenging

start with only wait/signal:

```
sem_t threads_to_wakeup;   // initially 0
Wait(Lock lock) {
    lock.Unlock();
    sem_wait(&threads_to_wakeup);
    lock.Lock();
}
Signal() {
    sem_post(&threads_to_wakeup);
}
```

annoying: signal wakes up non-waiting threads (in the far future)

# monitors with semaphores: cvs (better)

condition variables are more challenging

start with only wait/signal:

```
sem_t private_lock;  // initially 1
int num_waiters;
sem_t threads_to_wakeup;  // initially 0
Wait(Lock lock) {                          Signal() {
  sem_wait(&private_lock);                   sem_wait(&private_lock);
  ++num_waiters;                             if (num_waiters > 0) {
  sem_post(&private_lock);                     sem_post(&threads_to_wakeup);
  lock.Unlock();                               --num_waiters;
  sem_wait(&threads_to_wakeup);              }
  lock.Lock();                               sem_post(&private_lock);
}                                          }
```

# monitors with semaphores: broadcast

now allows broadcast:

```
sem_t private_lock;  // initially 1
int num_waiters;
sem_t threads_to_wakeup;  // initially 0
Wait(Lock lock) {
  sem_wait(&private_lock);
  ++num_waiters;
  sem_post(&private_lock);
  lock.Unlock();
  sem_wait(&threads_to_wakeup);
  lock.Lock();
}
```

```
Broadcast() {
  sem_wait(&private_lock);
  while (num_waiters > 0) {
    sem_post(&threads_to_wakeup);
    --num_waiters;
  }
  sem_post(&private_lock);
}
```

# monitors with semaphores: chosen order

if we want to make sure threads woken up <span style="color:red">in order</span>

```
ThreadSafeQueue<sem_t> waiters;
Wait(Lock lock) {
  sem_t private_semaphore;
  ... /* init semaphore
         with count 0 */
  waiters.Enqueue(&semaphore);          Signal() {
  lock.Unlock();                          sem_t *next = waiters.DequeueOrNull();
  sem_post(private_semaphore);            if (next != NULL) {
  lock.Lock();                              sem_post(next);
}                                         }
                                        }
```

# monitors with semaphores: chosen order

if we want to make sure threads woken up <span style="color:red">in order</span>

```
ThreadSafeQueue<sem_t> waiters;
Wait(Lock lock) {
  sem_t private_semaphore;
  ... /* init semaphore
         with count 0 */              Signal() {
  waiters.Enqueue(&semaphore);          sem_t *next = waiters.DequeueOrNull();
  lock.Unlock();                        if (next != NULL) {
  sem_post(private_semaphore);            sem_post(next);
  lock.Lock();                          }
}                                     }
```

(but now implement queue with semaphores…)

# reader/writer problem

some shared data

only one thread modifying (read+write) at a time

read-only access from multiple threads is safe

# reader/writer problem

some shared data

only one thread modifying (read+write) at a time

read-only access from multiple threads is safe

could use lock — but doesn't allow multiple readers

# reader/writer locks

abstraction: lock that distinguishes readers/writers

operations:
   read lock: wait until no writers
   read unlock: stop being registered as reader
   write lock: wait until no readers and no writers
   write unlock: stop being registered as writer

# reader/writer locks

abstraction: lock that distinguishes readers/writers

operations:
    read lock: wait until no writers
    read unlock: stop being registered as reader
    write lock: wait until no readers and no writers
    write unlock: stop being registered as writer

# pthread rwlocks

```
pthread_rwlock_t rwlock;
pthread_rwlock_init(&rwlock, NULL /* attributes */);
...
    pthread_rwlock_rdlock(&rwlock);
    ... /* read shared data */
    pthread_rwlock_unlock(&rwlock);

    pthread_rwlock_wrlock(&rwlock);
    ... /* read+write shared data */
    pthread_rwlock_unlock(&rwlock);

...
pthread_rwlock_destroy(&rwlock);
```

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
```

lock to protect shared state

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
```

state: number of active readers, writers

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
/* condition, signal when writers becomes 0 */
cond_t ok_to_read_cv;
/* condition, signal when readers + writers becomes 0 */
cond_t ok_to_write_cv;
```

conditions to wait for (no readers or writers, no writers)

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
/* condition, signal when writers becomes 0 */
cond_t ok_to_read_cv;
/* condition, signal when readers + writers becomes 0 */
cond_t ok_to_write_cv;
```

```
ReadLock() {
  mutex_lock(&lock);
  while (writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);
  }
  ++readers;
  mutex_unlock(&lock);
}
ReadUnlock() {
  mutex_lock(&lock);
  --readers;
  if (readers == 0) {
    cond_signal(&ok_to_write_cv);
  }
  mutex_unlock(&lock);
}
```

```
WriteLock() {
  mutex_lock(&lock);
  while (readers + writers != 0) {
    cond_wait(&ok_to_write_cv);
  }
  ++writers;
  mutex_unlock(&lock);
}
WriteUnlock() {
  mutex_lock(&lock);
  --writers;
  cond_signal(&ok_to_write_cv);
  cond_broadcast(&ok_to_read_cv);
  mutex_unlock(&lock);
}
```

broadcast — wakeup all readers when no writers

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
/* condition, signal when writers becomes 0 */
cond_t ok_to_read_cv;
/* condition, signal when readers + writers becomes 0 */
cond_t ok_to_write_cv;
ReadLock() {                          WriteLock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  while (writers != 0) {                while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);     cond_wait(&ok_to_write_cv);
  }                                     }
  ++readers;                            ++writers;
  mutex_unlock(&lock);                  mutex_unlock(&lock);
}                                     }
ReadUnlock() {                        WriteUnlock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  --readers;                            --writers;
  if (readers == 0) {                   cond_signal(&ok_to_write_cv);
    cond_signal(&ok_to_write_cv);       cond_broadcast(&ok_to_read_cv);
  }                                     mutex_unlock(&lock);
  mutex_unlock(&lock);                }
}
```

wakeup a single writer when no readers or writers

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
/* condition, signal when writers becomes 0 */
cond_t ok_to_read_cv;
/* condition, signal when readers + writers becomes 0 */
cond_t ok_to_write_cv;
ReadLock() {                          WriteLock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  while (writers != 0) {                while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);     cond_wait(&ok_to_write_cv);
  }                                     }
  ++readers;                            ++writers;
  mutex_unlock(&lock);                  mutex_unlock(&lock);
}                                     }
ReadUnlock() {                        WriteUnlock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  --readers;                            --writers;
  if (readers == 0) {                   cond_signal(&ok_to_write_cv);
    cond_signal(&ok_to_write_cv);       cond_broadcast(&ok_to_read_cv);
  }                                     mutex_unlock(&lock);
  mutex_unlock(&lock);                }
}
```

problem: wakeup readers first or writer first?

this solution: wake them all up and they fight! inefficient!

# reader/writer-priority

policy question: writers first or readers first?

    writers-first: no readers go when writer waiting

    readers-first: no writers go when reader waiting

previous implementation: whatever randomly happens

    writers signalled first, maybe gets lock first?

    …but non-determinstic in pthreads

can make explicit decision

# writer-priority (1)

```
mutex_t lock; cond_t ok_to_read_cv; cond_t ok_to_write_cv;
int readers = 0, writers = 0;
int waiting_writers = 0;
ReadLock() {                            WriteLock() {
  mutex_lock(&lock);                      mutex_lock(&lock);
  while (writers != 0                     ++waiting_writers;
         || waiting_writers != 0) {       while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);       cond_wait(&ok_to_write_cv, &lock);
  }                                       }
  ++readers;                              --waiting_writers;
  mutex_unlock(&lock);                    ++writers;
}                                         mutex_unlock(&lock);
                                        }
ReadUnlock() {
  mutex_lock(&lock);                    WriteUnlock() {
  --readers;                              mutex_lock(&lock);
  if (readers == 0) {                     --writers;
    cond_signal(&ok_to_write_cv);         if (waiting_writers != 0) {
  }                                         cond_signal(&ok_to_write_cv);
  mutex_unlock(&lock);                    } else {
}                                           cond_broadcast(&ok_to_read_cv);
                                          }
                                          mutex_unlock(&lock);
                                        }
```

# writer-priority (1)

```
mutex_t lock; cond_t ok_to_read_cv; cond_t ok_to_write_cv;
int readers = 0, writers = 0;
int waiting_writers = 0;
ReadLock() {                          WriteLock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  while (writers != 0                   ++waiting_writers;
         || waiting_writers != 0) {     while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);     cond_wait(&ok_to_write_cv, &lock);
  }                                     }
  ++readers;                            --waiting_writers;
  mutex_unlock(&lock);                  ++writers;
}                                       mutex_unlock(&lock);
                                      }
ReadUnlock() {
  mutex_lock(&lock);                  WriteUnlock() {
  --readers;                            mutex_lock(&lock);
  if (readers == 0) {                   --writers;
    cond_signal(&ok_to_write_cv);       if (waiting_writers != 0) {
  }                                       cond_signal(&ok_to_write_cv);
  mutex_unlock(&lock);                  } else {
}                                         cond_broadcast(&ok_to_read_cv);
                                        }
                                        mutex_unlock(&lock);
                                      }
```

# writer-priority (1)

```
mutex_t lock; cond_t ok_to_read_cv; cond_t ok_to_write_cv;
int readers = 0, writers = 0;
int waiting_writers = 0;
ReadLock() {                          WriteLock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  while (writers != 0                   ++waiting_writers;
         || waiting_writers != 0) {     while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);     cond_wait(&ok_to_write_cv, &lock);
  }                                     }
  ++readers;                            --waiting_writers;
  mutex_unlock(&lock);                  ++writers;
}                                       mutex_unlock(&lock);
                                      }
ReadUnlock() {
  mutex_lock(&lock);                  WriteUnlock() {
  --readers;                            mutex_lock(&lock);
  if (readers == 0) {                   --writers;
    cond_signal(&ok_to_write_cv);       if (waiting_writers != 0) {
  }                                       cond_signal(&ok_to_write_cv);
  mutex_unlock(&lock);                  } else {
}                                         cond_broadcast(&ok_to_read_cv);
                                        }
                                        mutex_unlock(&lock);
                                      }
```

# reader-priority (1)

```
...
int waiting_readers = 0;
ReadLock() {
  mutex_lock(&lock);
  ++waiting_readers;
  while (writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);
  }
  --waiting_readers;
  ++readers;
  mutex_unlock(&lock);
}

ReadUnlock() {
  ...
  if (waiting_readers == 0) {
    cond_signal(&ok_to_write_cv);
  }
}
```

```
WriteLock() {
  mutex_lock(&lock);
  while (waiting_readers +
         readers + writers != 0) {
    cond_wait(&ok_to_write_cv);
  }
  ++writers;
  mutex_unlock(&lock);
}
WriteUnlock() {
  mutex_lock(&lock);
  --writers;
  if (readers == 0 && waiting_readers == 0) {
    cond_signal(&ok_to_write_cv);
  } else {
    cond_broadcast(&ok_to_read_cv);
  }
  mutex_unlock(&lock);
}
```

# reader-priority (1)

```
...
int waiting_readers = 0;
ReadLock() {
  mutex_lock(&lock);
  ++waiting_readers;
  while (writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);
  }
  --waiting_readers;
  ++readers;
  mutex_unlock(&lock);
}

ReadUnlock() {
  ...
  if (waiting_readers == 0) {
    cond_signal(&ok_to_write_cv);
  }
}
```

```
WriteLock() {
  mutex_lock(&lock);
  while (waiting_readers +
         readers + writers != 0) {
    cond_wait(&ok_to_write_cv);
  }
  ++writers;
  mutex_unlock(&lock);
}
WriteUnlock() {
  mutex_lock(&lock);
  --writers;
  if (readers == 0 && waiting_readers == 0) {
    cond_signal(&ok_to_write_cv);
  } else {
    cond_broadcast(&ok_to_read_cv);
  }
  mutex_unlock(&lock);
}
```

# choosing orderings?

can use monitors to implement lots of lock policies

want $X$ to go first/last — add extra variables
(number of waiters, even lists of items, etc.)

need way to write condition "you can go now"
e.g. writer-priority: readers can go if no writer waiting