

Lecture 10 — Use of Locks, Reentrancy

Jeff Zarnett & Patrick Lam

{jzarnett, patrick.lam}@uwaterloo.ca

Department of Electrical and Computer Engineering
University of Waterloo

October 29, 2019

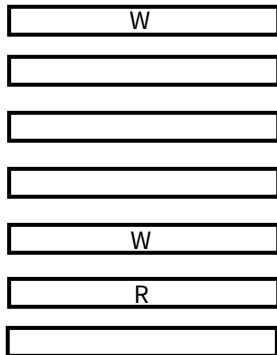
Part I

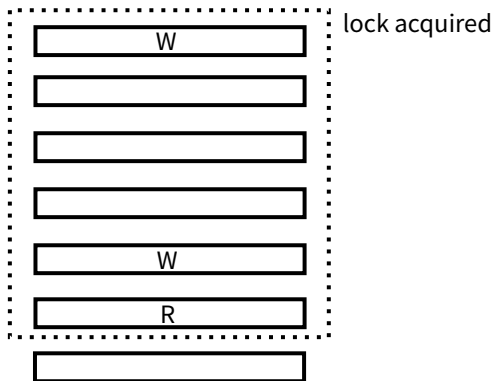
Use of Locks

In previous courses: learned about locking and how it works.

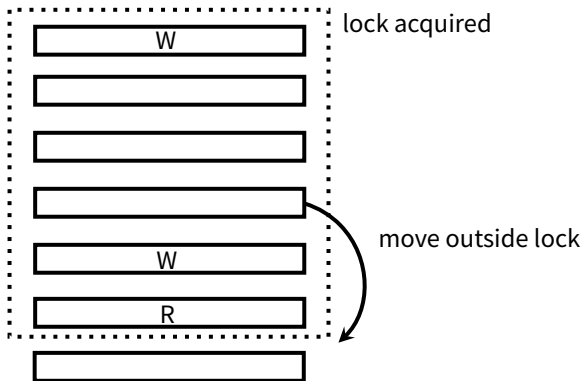
Before: was enough to avoid bad stuff.

Now: Still important, but not enough.





Critical sections should be as large as they need to be but no larger.



Sometimes control flow or other very short statements might get swept into the critical section.

Those should be the exception rather than the rule.

Remember the Producer-Consumer Problem?

Short code example from the producer-consumer problem:

```
sem_t spaces;  
sem_t items;  
int counter;  
int* buffer;  
int pindex = 0;  
int cindex = 0;  
int ctotat = 0;  
pthread_mutex_t prod_mutex;  
pthread_mutex_t con_mutex;  
  
void consume( int to_consume );
```

Remember the Producer-Consumer Problem?

Single-threaded consumer code:

```
void* consumer( void* arg ) {  
    while( ctotol < MAX_ITEMS_CONSUMED ) {  
        sem_wait( &items );  
        consume( buffer[cindex] );  
        buffer[cindex] = -1;  
        cindex = (cindex + 1) % BUFFER_SIZE;  
        ++ctotol;  
        sem_post( &spaces );  
    }  
}
```

Must add mutual exclusion to allow multiple simultaneous consumers.

Could allow exactly one consumer to run at a time. We can do better...

```
void* consumer( void* arg ) {  
    while( ctotal < MAX_ITEMS_CONSUMED ) {  
        pthread_mutex_lock( &con_mutex );  
        sem_wait( &items );  
        consume( buffer[cindex] );  
        buffer[cindex] = -1;  
        cindex = (cindex + 1) % BUFFER_SIZE;  
        ++ctotal;  
        sem_post( &spaces );  
        pthread_mutex_unlock( &con_mutex );  
    }  
}
```

Anything to Kick Out?



At first glance it is probably not very obvious but the consume function takes a regular integer, any old integer, not a pointer of some sort.

So we could, inside the critical section, read the value of the buffer at the current index into a temp variable.

That temp variable then can be given to the consume function at any time...outside of the critical section.

Everything else inside our lock and unlock statements seems to be shared data: operates on `cindex` or `ctotal`.

```
void* consumer( void* arg ) {  
    while( ctotal < MAX_ITEMS_CONSUMED ) {  
        pthread_mutex_lock( &con_mutex );  
        sem_wait( &items );  
        int temp = buffer[cindex]; // store for now  
        buffer[cindex] = -1;  
        cindex = (cindex + 1) % BUFFER_SIZE;  
        ++ctotal;  
        sem_post( &spaces );  
        pthread_mutex_unlock( &con_mutex );  
        consume( temp ); // consume outside mutex  
    }  
}
```

The while condition checks the value of `ctotal` and that is a read of shared data.

Now we maybe have a problem. How do we get that inside the critical section?

One idea we might have is to read the value of `ctotal` into a temporary variable and use that, but it might cause some headaches with the timing...

Not Forgetting The Loop Condition

Instead what I'd recommend is to make the loop a while true loop and then have a test of the value to determine when we should break out of the loop.

```
void* consumer( void* arg ) {  
    while( 1 ) {  
        pthread_mutex_lock( &con_mutex );  
        if ( ctotal == MAX_ITEMS_CONSUMED ) {  
            pthread_mutex_unlock( &con_mutex );  
            break;  
        }  
        sem_wait( &items );  
        int temp = buffer[cindex];  
        buffer[cindex] = -1;  
        cindex = (cindex + 1) % BUFFER_SIZE;  
        ++ctotal;  
        pthread_mutex_unlock( &con_mutex );  
        sem_post( &spaces );  
        consume( temp );  
    }  
    pthread_exit( NULL );  
}
```

Are we done?

“Do you think you got him?”



Keeping the critical section as small as possible is important because it speeds up performance (reduces the serial portion of your program).

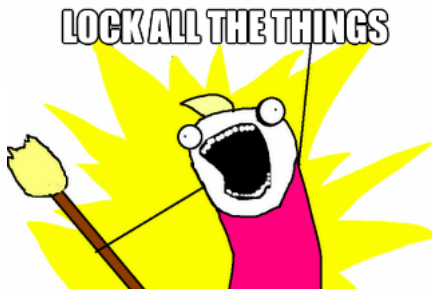
But that's not the only reason.

The lock is a resource, and contention for that resource is itself expensive.

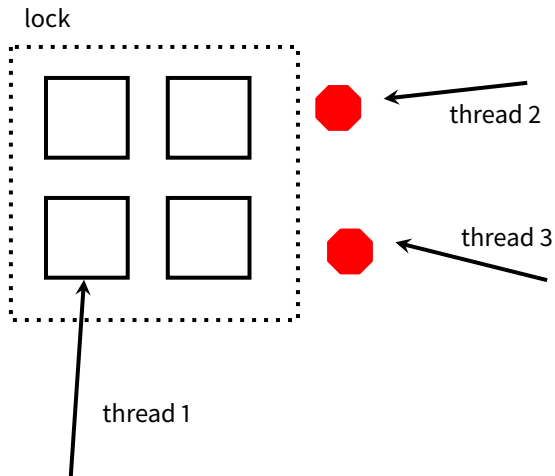
Alright, we already know that locks prevent data races.

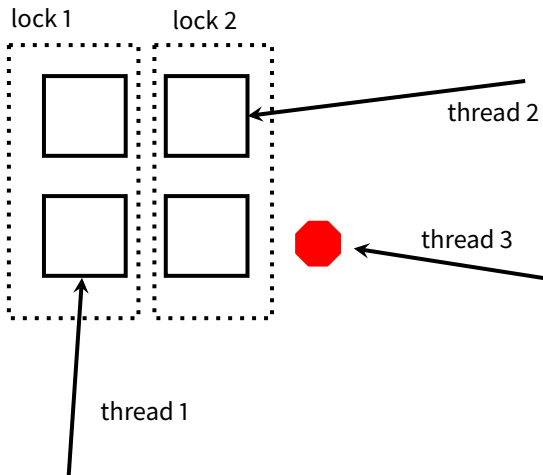
So we need to use them to prevent data races, but it's not as simple as it sounds.

We have choices about the granularity of locking, and it is a trade-off.



Coarse-Grained Lock





Coarse-Grained vs Fine-Grained Locking

Coarse-grained locking is easier to write and harder to mess up, but it can significantly reduce opportunities for parallelism.

Fine-grained locking requires more careful design, increases locking overhead and is more prone to bugs (deadlock etc).

Locks' extents constitute their granularity.

We'll discuss three major concerns when using locks:

- overhead;
- contention; and
- deadlocks.

We aren't even talking about under-locking (i.e., remaining race conditions).

Using a lock isn't free. You pay:

- allocated memory for the locks;
- initialization and destruction time; and
- acquisition and release time.

These costs scale with the number of locks that you have.

Most locking time is wasted waiting for the lock to become available.

We can fix this by:

- making the locking regions smaller (more granular); or
- making more locks for independent sections.

Finally, the more locks you have, the more you have to worry about deadlocks.



As you know, the key condition for a deadlock is waiting for a lock held by process X while holding a lock held by process X' . ($X = X'$ is allowed).

Okay, in a formal sense, the four conditions for deadlock are:

- 1 Mutual Exclusion**
- 2 Hold-and-Wait**
- 3 No Preemption**
- 4 Circular-Wait**

Consider, for instance, two processors trying to get two locks.

Thread 1

Get Lock 1
Get Lock 2
Release Lock 2
Release Lock 1

Thread 2

Get Lock 2
Get Lock 1
Release Lock 1
Release Lock 2

To avoid deadlocks, always be careful if your code **acquires a lock while holding one**.

You have two choices: (1) ensure consistent ordering in acquiring locks; or (2) use trylock.

As an example of consistent ordering:

```
void f1 () {  
    lock(&l1);  
    lock(&l2);  
    // protected code  
    unlock(&l2);  
    unlock(&l1);  
}
```

```
void f2 () {  
    lock(&l1);  
    lock(&l2);  
    // protected code  
    unlock(&l2);  
    unlock(&l1);  
}
```

Alternately, you can use trylock.

Recall that Pthreads' `trylock` returns 0 if it gets the lock.

But if it doesn't, your thread doesn't get blocked.

```
void f1 () {  
    lock(&l1);  
    while (trylock(&l2) != 0) {  
        unlock(&l1);  
        // wait  
        lock(&l1);  
    }  
    // protected code  
    unlock(&l2);  
    unlock(&l1);  
}
```

This prevents the hold and wait condition.

There Can Be Only One



Photo Credit: Craig Middleton¹

¹<https://www.flickr.com/photos/craigmiddleton/3579668458>

One way of avoiding problems due to locking is to use few locks (1?).

This is *coarse-grained locking*; it does have a couple of advantages:

- it is easier to implement;
- with one lock, there is no chance of deadlocking; and
- it has the lowest memory usage and setup time possible.

Your parallel program will quickly become sequential.

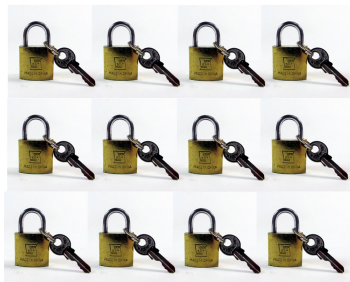
Python puts a lock around the whole interpreter (known as the *global interpreter lock*).

This is the main reason (most) scripting languages have poor parallel performance; Python's just an example.

Two major implications:

- The only performance benefit you'll see from threading is if one of the threads is waiting for I/O.
- But: any non-I/O-bound threaded program will be **slower** than the sequential version (plus, the interpreter will slow down your system).

On the other end of the spectrum is *fine-grained locking*.



The big advantage: it maximizes parallelization in your program.

However, it also comes with a number of disadvantages:

- if your program isn't very parallel, it'll be mostly wasted memory and setup time;
- plus, you're now prone to deadlocks; and
- fine-grained locking is generally more error-prone (be sure you grab the right lock!)

Databases may lock fields / records / tables. (fine-grained → coarse-grained).

You can also lock individual objects (but beware: sometimes you need transactional guarantees.)

Part II

Reentrancy



The trivial example of a non-reentrant C function:

```
int tmp;  
  
void swap( int x, int y ) {  
    tmp = y;  
    y = x;  
    x = tmp;  
}
```

Why is this non-reentrant?

How can we make it reentrant?

⇒ A function can be suspended in the middle and **re-entered** (called again) before the previous execution returns.

Does not always mean **thread-safe** (although it usually is).

- Recall: **thread-safe** is essentially “no data races”.

Moot point if the function only modifies local data, e.g. `sin()`.

Courtesy of Wikipedia (with modifications):

```
int t;

void swap(int *x, int *y) {
    t = *x;
    *x = *y;
    // hardware interrupt might invoke isr() here!
    *y = t;
}

void isr() {
    int x = 1, y = 2;
    swap(&x, &y);
}

...
int a = 3, b = 4;
...
    swap(&a, &b);
```

Reentrancy Example—Explained (a trace)

```
call swap(&a, &b);  
  t = *x;           // t = 3 (a)  
  *x = *y;          // a = 4 (b)  
  call isr();  
    x = 1; y = 2;  
    call swap(&x, &y)  
      t = *x;       // t = 1 (x)  
      *x = *y;      // x = 2 (y)  
      *y = t;       // y = 1  
    *y = t;         // b = 1
```

Final values:

a = 4, b = 1

Expected values:

a = 4, b = 3

```
int t;

void swap(int *x, int *y) {
    int s;

    @\alert{s = t}@\; // save global variable
    t = *x;
    *x = *y;
    // hardware interrupt might invoke isr() here!
    *y = t;
    @\alert{t = s}@\; // restore global variable
}

void isr() {
    int x = 1, y = 2;
    swap(&x, &y);
}

...
int a = 3, b = 4;
...
    swap(&a, &b);
```

Reentrancy Example, Fixed—Explained (a trace)

```
call swap(&a, &b);  
s = t;           // s = UNDEFINED  
t = *x;          // t = 3 (a)  
*x = *y;         // a = 4 (b)  
call isr();  
    x = 1; y = 2;  
    call swap(&x, &y)  
        s = t;           // s = 3  
        t = *x;          // t = 1 (x)  
        *x = *y;         // x = 2 (y)  
        *y = t;          // y = 1  
        t = s;           // t = 3  
    *y = t;         // b = 3  
    t = s;          // t = UNDEFINED
```

Final values:

a = 4, b = 3

Expected values:

a = 4, b = 3

Remember that in things like interrupt subroutines (ISRs) having the code be reentrant is very important.

Interrupts can get interrupted by higher priority interrupts and when that happens the ISR may simply be restarted, or we pause and resume.

Either way, if the code is not reentrant we will run into problems.

Let us also draw a distinction between thread safe code and reentrant code.

A thread safe operation is one that can be performed from more than one thread at the same time.

On the other hand, a reentrant operation can be invoked while the operation is already in progress, possibly from within the same thread.

Or it can be re-started without affecting the outcome.

If a function is not reentrant, it may not be possible to make it thread safe.

And furthermore, a reentrant function cannot call a non-reentrant one (and maintain its status as reentrant).

Thread Safe Non-Reentrant Example

```
int length = 0;
char *s = NULL;

// Note: Since strings end with a 0, if we want to
// add a 0, we encode it as "\0", and encode a
// backslash as "\\".

// WARNING! This code is buggy — do not use!
void AddToString(int ch)
{
    EnterCriticalSection(&someCriticalSection);
    // +1 for the character we're about to add
    // +1 for the null terminator
    char *newString = realloc(s, (length+1) * sizeof(char));
    if (newString) {
        if (ch == '\0' || ch == '\\') {
            AddToString('\\'); // escape prefix
        }
        newString[length++] = ch;
        newString[length] = '\0';
        s = newString;
    }
    LeaveCriticalSection(&someCriticalSection);
}
```

Is it thread safe? Sure—there is a critical section protected by the mutex `someCriticalSection`.

But is it re-entrant? Nope.

The internal call to `AddToString` causes a problem because the attempt to use `realloc` will use a pointer to `s`.

That is no longer valid because it got stomped by the earlier call to `realloc`.

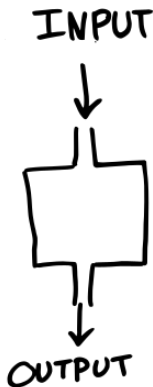
Interestingly, functional programming languages (NOT procedural like C) such as Scala and so on, lend themselves very nicely to being parallelized.

Why?

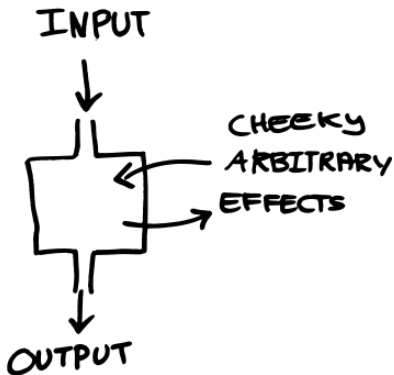
Because a purely functional program has no side effects and they are very easy to parallelize.

Any impure function has to indicate that in its function signature.

Functions



Procedures



<https://www.fpcomplete.com/blog/2017/04/pure-functional-programming>

Without understanding functional programming, you can't invent MapReduce, the algorithm that makes Google so massively scalable. The terms Map and Reduce come from Lisp and functional programming. MapReduce is, in retrospect, obvious to anyone who remembers from their 6.001-equivalent programming class that purely functional programs have no side effects and are thus trivially parallelizable.

— Joel Spolsky

Object oriented programming kind of gives us some bad habits in this regard.

We tend to make a lot of `void` methods.

In functional programming these don't really make sense, because if it's purely functional, then there are some inputs and some outputs.

If a function returns nothing, what does it do?

For the most part it can only have side effects which we would generally prefer to avoid if we can, if the goal is to parallelize things.

`algorithms` has been part of C++ since C++11.

C++17 (not well supported yet) introduces parallel and vectorized algorithms;
you specify an execution policy, compiler does it:
(`sequenced`, `parallel`, or `parallel_unsequenced`).

Some examples of algorithms:

`sort`, `reverse`, `is_heap`...

Other examples of algorithms:

`for_each_n`, `exclusive_scan`, `reduce`

If you know functional programming (e.g. Haskell), these are:

`map`, `scanl`, `foldl`/`foldl1`.

Side effects are sort of undesirable (and definitely not functional), but not necessarily bad.

Printing to console is unavoidably making use of a side effect, but it's what we want.

We don't want to call print reentrantly, because two threads trying to write at the same time to the console would result in jumbled output.

Or alternatively, restarting the print routine might result in some doubled characters on the screen.

The notion of purity is related to side effects.

A function is *pure* if it has no side effects and if its outputs depend solely on its inputs.

Pure functions should be implemented as thread-safe and reentrant.