

# Lecture 09: Transitioning from C to C++ Threads

- Introverts Revisited, in C++
  - Rather than deal with **pthread**s as a platform-specific extension of C, I'd rather use a thread package that's officially integrated into the language itself.
    - As of 2011, C++ provides [support for threading](#) and many synchronization directives.
    - Because C++ provides better alternatives for generic programming than C does, we avoid the **void \*** tomfoolery required when using **pthread**s .
  - Presented below is the object-oriented C++ equivalent of the **introverts** example we've already seen once before. The full program is online [right here](#).

```
static void recharge() {
    cout << oslock << "I recharge by spending time alone." << endl << osunlock;
}

static const size_t kNumIntroverts = 6;
int main(int argc, char *argv[]) {
    cout << "Let's hear from " << kNumIntroverts << " introverts." << endl
    thread introverts[kNumIntroverts]; // declare array of empty thread handles
    for (thread& introvert: introverts)
        introvert = thread(recharge); // move anonymous threads into empty handles
    for (thread& introvert: introverts)
        introvert.join();
    cout << "Everyone's recharged!" << endl;
    return 0;
}
```

# Lecture 09: Transitioning from C to C++ Threads

- We declare an array of empty **thread** handles as we did in the equivalent C version.
- We install the **recharge** function into temporary threads that are then moved (via the **thread's operator=(thread&& other)**) into a previously empty **thread** handle.
  - This is a relatively new form of **operator=** that fully transplants the contents of the **thread** on the right into the **thread** on the left, leaving the **thread** on the right fully gutted, as if it were zero-arg constructed. Restated, the left and right thread objects are effectively swapped.
  - This is an important distinction, because a traditional **operator=** would produce a second working copy of the same **thread**, and we don't want that.
- The **join** method mirrors the **pthread\_join** function we've already discussed.
- The prototype of the thread routine—in this case, **recharge**—can be anything (although the return type is always ignored, so it should generally be **void**).
- **operator<<**, unlike **printf**, isn't thread-safe.
  - I've constructed custom stream manipulators called **oslock** and **osunlock** that can be used to acquire and release exclusive access to an **ostream**.
  - These manipulators—which we can use by **#include**-ing "**ostreamlock.h**"—can be used to ensure at most one thread has permission to write into a stream at any one time.

# Lecture 09: Transitioning from C to C++ Threads

- Thread routines can accept any number of arguments using variable argument lists. (Variable argument lists—the C++ equivalent of the ellipsis in C—are supported via a recently added feature called [variadic templates](#).)
- Here's a [slightly more involved example](#), where **greet** threads are configured to say hello a variable number of times.

```
static void greet(size_t id) {
    for (size_t i = 0; i < id; i++) {
        cout << oslock << "Greeter #" << id << " says 'Hello!'" << endl << osunlock;
        struct timespec ts = {
            0, random() % 1000000000
        };
        nanosleep(&ts, NULL);
    }
    cout << oslock << "Greeter #" << id << " has issued all of his hellos, "
        << "so he goes home!" << endl << osunlock;
}

static const size_t kNumGreeters = 6;
int main(int argc, char *argv[]) {
    cout << "Welcome to Greetland!" << endl;
    thread greeters[kNumGreeters];
    for (size_t i = 0; i < kNumGreeters; i++) greeters[i] = thread(greet, i + 1);
    for (thread& greeter: greeters) greeter.join();
    cout << "Everyone's all greeted out!" << endl;
    return 0;
}
```

# Lecture 09: Transitioning from C to C++ Threads

- Multiple threads are often spawned to subdivide and collectively solve a larger problem. Full program is [right here](#).
- Consider the scenario where 10 ticket agents answer telephones—as they might have before the internet came along—at United Airlines to jointly sell 250 airline tickets.
  - Each ticket agent answers the telephone, and each telephone call always leads to the sale of precisely one ticket.
  - Rather than requiring each ticket agent sell 10% of the tickets, we'll account for the possibility that some ticket sales are more time consuming than others, some ticket agents need more time in between calls, etc. Instead, we'll require that all ticket agents keep answering calls and selling tickets until all have been sold.
- Here's our first stab at a **main** function.

```
int main(int argc, const char *argv[]) {
    thread agents[10];
    size_t remainingTickets = 250;
    for (size_t i = 0; i < 10; i++)
        agents[i] = thread(ticketAgent, 101 + i, ref(remainingTickets));
    for (thread& agent: agents) agent.join();
    cout << "End of Business Day!" << endl;
    return 0;
}
```

# Lecture 09: Transitioning from C to C++ Threads

- As with most multithreaded programs, the main thread elects to spawn child threads to subdivide and collaboratively solve the full problem at hand.
  - In this case, the **main** function declares the master copy of the remaining ticket count—aptly named **remainingTickets**—and initializes it to 250.
  - The main thread then spawns ten child threads to run some **ticketAgent** thread routine, yet to be fully defined. Each agent is assigned a unique id number between 101 and 110, inclusive, and a reference to **remainingTickets** is shared with each thread.
  - As is typical, the main thread blocks until all child threads have finished before exiting. Otherwise, the entire process might be torn down even though some child threads haven't finished.

```
int main(int argc, const char *argv[]) {
    thread agents[10];
    size_t remainingTickets = 250;
    for (size_t i = 0; i < 10; i++)
        agents[i] = thread(ticketAgent, 101 + i, ref(remainingTickets));
    for (thread& agent: agents) agent.join();
    cout << "End of Business Day!" << endl;
    return 0;
}
```

# Lecture 09: Transitioning from C to C++ Threads

- The `ticketAgent` thread routine accepts an id number (used for logging purposes) and a reference to the `remainingTickets`.
- It continually polls `remainingTickets` to see if any tickets remain, and if so, proceeds to answer the phone, sell a ticket, and publish a little note about the ticket sale to `cout`.
- `handleCall`, `shouldTakeBreak`, and `takeBreak` are all in place to introduce short, random delays and guarantee that each test run is different than prior ones. Full program is (still) [right here](#).

```
static void ticketAgent(size_t id, size_t& remainingTickets) {
    while (remainingTickets > 0) {
        handleCall(); // sleep for a small amount of time to emulate conversation time.
        remainingTickets--;
        cout << oslock << "Agent #" << id << " sold a ticket! (" << remainingTickets
            << " more to be sold)." << endl << osunlock;
        if (shouldTakeBreak()) // flip a biased coin
            takeBreak();        // if comes up heads, sleep for a random time to take a break
    }
    cout << oslock << "Agent #" << id << " notices all tickets are sold, and goes home!"
        << endl << osunlock;
}
```

# Lecture 09: Transitioning from C to C++ Threads

- Presented below right is the abbreviated output of a **confused-ticket-agents** run.
- In its current state, the program suffers from a serious race condition.
- Why? Because **remainingTickets > 0** test and **remainingTickets--** aren't guaranteed to execute within the same time slice.
- If a thread evaluates **remainingTickets > 0** to be **true** and commits to selling a ticket, the ticket might not be there by the time it executes the decrement. That's because the thread may be swapped off the CPU after the decision to sell but before the sale, and during the dead time, other threads—perhaps the nine others—all might get the CPU and do precisely the same thing.
- The solution? Ensure the decision to sell and the sale itself are executed without competition.

```
poohbear@myth61:$ ./confused-ticket-agents
Agent #110 sold a ticket! (249 more to be sold).
Agent #104 sold a ticket! (248 more to be sold).
Agent #106 sold a ticket! (247 more to be sold).
// some 245 lines omitted for brevity
Agent #107 sold a ticket! (1 more to be sold).
Agent #103 sold a ticket! (0 more to be sold).
Agent #105 notices all tickets are sold, and goes home!
Agent #104 notices all tickets are sold, and goes home!
Agent #108 sold a ticket! (4294967295 more to be sold).
Agent #106 sold a ticket! (4294967294 more to be sold).
Agent #102 sold a ticket! (4294967293 more to be sold).
Agent #101 sold a ticket! (4294967292 more to be sold).
// carries on for a very, very, very long time
```

# Lecture 09: Transitioning from C to C++ Threads

- Before we solve this problem, we should really understand why **remainingTickets--** itself isn't even thread-safe.
  - C++ statements aren't inherently atomic. Virtually all C++ statements—even ones as simple as **remainingTickets--**—compile to multiple assembly code instructions.
  - Assembly code instructions are atomic, but C++ statements are not.
  - **g++** on the myths compiles **remainingTickets--** to five assembly code instructions, as with:

```
0x0000000000401a9b <+36>:   mov     -0x20(%rbp), %rax
0x0000000000401a9f <+40>:   mov     (%rax), %eax
0x0000000000401aa1 <+42>:   lea    -0x1(%rax), %edx
0x0000000000401aa4 <+45>:   mov     -0x20(%rbp), %rax
0x0000000000401aa8 <+49>:   mov     %edx, (%rax)
```

- The first two lines drill through the **ticketsRemaining** reference to load a copy of the **ticketsRemaining** held in **main**. The third line decrements that copy, and the last two write the decremented copy back to the **ticketsRemaining** variable held in **main**.

# Lecture 09: Transitioning from C to C++ Threads

- Before we solve this problem, we should really understand why `remainingTickets--` itself isn't even thread-safe.
  - If a thread executes the first three (or two, or four) of these instructions and gets swapped off the CPU, it will eventually get the processor back, but by that time the local copy will be outdated if one or more other threads have since gotten processor time and executed some or all of the same five assembly code instructions below.
  - The problem: ALU operations are executed on copies of shared data

```
0x0000000000401a9b <+36>:    mov     -0x20(%rbp), %rax
0x0000000000401a9f <+40>:    mov     (%rax), %eax
0x0000000000401aa1 <+42>:    lea    -0x1(%rax), %edx
0x0000000000401aa4 <+45>:    mov     -0x20(%rbp), %rax
0x0000000000401aa8 <+49>:    mov     %edx, (%rax)
```

- A specific pathological example: `main`'s copy of `remainingTickets` is 250, and a thread manages to execute the first three of these five instructions before before swapped off the CPU. `%edx` stores a 249, but the thread has yet to push that 249 back to `main`'s `remainingTickets`. It will when it gets the processor again, but it might not get the processor until the other threads have cooperatively sold one more tickets (or perhaps even all of them).

# Lecture 09: Transitioning from C to C++ Threads

- We need to guarantee that the code that tests for remaining tickets, sells a ticket, and everything in between are executed as part of one large transaction, without interference from other threads. Restated, we must guarantee that at no other threads are permitted to even **examine** the value of **ticketsRemaining** if another thread is staged to modify it.
- One solution: provide a directive that allows a thread to ask that it not be swapped off the CPU while it's within a block of code that should be executed transactionally.
  - That, however, is not an option, and shouldn't be.
  - That would grant too much power to threads, which could abuse the option and block other threads from running for an indeterminate amount of time.
- The other option is to rely on a concurrency directive that can be used to prevent more than one thread from being anywhere in the same critical region at one time. That concurrency directive is the **mutex**, and in C++ it looks like this:

```
class mutex {  
public:  
    mutex();           // constructs the mutex to be in an unlocked state  
    void lock();      // acquires the lock on the mutex, blocking until it's unlocked  
    void unlock();   // releases the lock and wakes up another threads trying to lock it  
};
```

# Lecture 09: Transitioning from C to C++ Threads

- The name **mutex** is just a contraction of the words *mutual* and *exclusion*. It's so named because its primary use is to mark the boundaries of a critical region—that is, a stretch of code where at most one thread is permitted to be at any one moment.
  - Restated, a thread executing code within a critical region enjoys exclusive access.
- The constructor initializes the **mutex** to be in an unlocked state.
- The **lock** method will *eventually* acquire a lock on the **mutex**.
  - If the **mutex** is in an unlocked state, **lock** will lock it and return immediately.
  - If the **mutex** is in a locked state (presumably because another thread called **lock** but has yet to **unlock**), **lock** will pull the calling thread off the CPU and render it ineligible for processor time until notified the lock on the **mutex** was released.
- The **unlock** method will release the lock on a mutex. The only thread qualified to release the lock on the **mutex** is the one that holds the lock.

```
class mutex {  
public:  
    mutex();           // constructs the mutex to be in an unlocked state  
    void lock();      // acquires the lock on the mutex, blocking until it's unlocked  
    void unlock();    // releases the lock and wakes up another threads trying to lock it  
};
```

# Lecture 09: Transitioning from C to C++ Threads

- We can declare a single **mutex** aside the declaration of **remainingTickets** in the **main** function, and we can use that **mutex** to mark the boundaries of the critical region.
- This requires the **mutex** also be shared by reference with the **ticketAgent** thread routine so that all child threads compete to acquire the same lock.
- The new **ticketAgent** thread routine looks like this:

```
static void ticketAgent(size_t id, size_t& remainingTickets, mutex& ticketsLock) {
    while (true) {
        ticketsLock.lock();
        if (remainingTickets == 0) break;
        handleCall();
        remainingTickets--;
        cout << oslock << "Agent #" << id << " sold a ticket! (" << remainingTickets
            << " more to be sold)." << endl << osunlock;
        ticketsLock.unlock();
        if (shouldTakeBreak())
            takeBreak();
    }
    ticketsLock.unlock();
    cout << oslock << "Agent #" << id << " notices all tickets are sold, and goes home!"
        << endl << osunlock;
}
```

# Lecture 09: Transitioning from C to C++ Threads

- The **main** function might look like the one I present below (and [right here](#)).
  - **main** declares a single **mutex** called **ticketLock**. I can call it anything I want to, but I want to be clear that it's in place to guard any potential surgery on a variable with **tickets** in its name. That's why I choose the name I do.
  - Because the **ticketAgent** accepts a reference to a **mutex**, the **thread** constructor needs to pass a reference to a **mutex** via an addition argument that wasn't present in the first version. That's what **ref(ticketLock)** is all about.
  - Fundamentally, the **main** function maintains the state information needed to enable all child threads to collaboratively work in parallel without introducing any race conditions.

```
int main(int argc, const char *argv[]) {
    size_t remainingTickets = 250;
    mutex ticketsLock;
    thread agents[10];
    for (size_t i = 0; i < 10; i++)
        agents[i] = thread(ticketAgent, 101 + i, ref(remainingTickets), ref(ticketsLock));
    for (thread& agent: agents) agent.join();
    cout << "End of Business Day!" << endl;
    return 0;
}
```