

+ 21

# Back to Basics: Concurrency

MIKE SHAH



20  
21 | A graphic of three white mountain peaks with a yellow peak on top, positioned next to the year '2021'.  
October 24-29

---

Please do not redistribute slides without  
prior permission.



# Back to Basics: Concurrency

Mike Shah, Ph.D.

[@MichaelShah](https://twitter.com/MichaelShah) | [mshah.io](https://mshah.io) | [www.youtube.com/c/MikeShah](https://www.youtube.com/c/MikeShah)

4:45 pm MDT, Mon. October 25

60 minutes | Introductory Audience

# Abstract

## The abstract that you read and enticed you to join me is here!

You have spent your hard earned money on a multi-core machine. But what does that mean for you as a programmer or for the consumers of your software who also spent their hard earned money on a multi-core machine? Well the deal is, you only get an increase in performance, if you know how to take advantage of your hardware. Perhaps you have also heard something about the free lunch being over for programmers?

In this talk we provide a gentle introduction to concurrency with the modern C++ std::thread library. We will introduce topics with pragmatic code examples introducing the ideas of threads and locks, and showing how these programming primitives enable concurrency. In addition, we will show the common pitfalls to be aware of with concurrency: data races, starvation, and deadlock (the most extreme form of starvation!). But don't worry--I will show you how to fix these concurrency bugs!

The session will wrap up with discussion and examples of other concurrency primitives and how to use them to write concurrent programs for common patterns(different types of locks, conditional variables, promises/futures). Attendees will leave this session being able to use threads, locks, and start thinking about architecting multithreaded software. All materials and code samples will also be posted online.

# Code for the talk

---

Available here: <https://github.com/MikeShah/cppcon2021>

The screenshot shows a GitHub repository interface. At the top, there's a list of files:

- concurrency (thread example) - 12 seconds ago
- pointers (More pointer examples) - 4 hours ago
- README.md (Update README.md) - 10 hours ago

Below this is a preview of the README.md file content:

**README.md**

**cppcon2021**

---

Examples and materials for my talks during Cppcon 2021!

# Who Am I?

by Mike Shah

- Assistant Teaching Professor at Northeastern University in Boston, Massachusetts.
  - I teach courses in computer systems, computer graphics, and game engine development.
  - My research in program analysis is related to performance building static/dynamic analysis and software visualization tools.
- I do consulting and technical training on modern C++, Concurrency, OpenGL, and Vulkan projects
  - (Usually graphics or games related)
- I like teaching, guitar, running, weight training, and anything in computer science under the domain of computer graphics, visualization, concurrency, and parallelism.
- Contact information and more on: [www.mshah.io](http://www.mshah.io)



# What you are going to learn today (1/2)

---

- (So you know what to pay attention to)
  - **The foundations of Concurrency (From a beginners standpoint)**
- Pragmatically
  - You are someone who wants to learn about concurrency, and wants to learn about some of the fundamental building blocks to write for instance a multi-threaded program
  - We are going to focus on std::thread for creating concurrent programs.
    - We will review briefly other methods of writing concurrent programs (e.g., std::async)
  - I'll focus on the foundations, and my hope is this will help you also choose the right higher level primitives to write more scalable concurrent software.

Concurrency is a massive topic--so let's start from the beginning

# What you are going to learn today (2/2)

---

- (So you know what to pay attention to)
  - **The foundations of Concurrency (From a beginners standpoint)**
- Pragmatically
  - Talk contains two core components
    1. Motivation and the ‘why’ to concurrency
    2. Several toy examples to learn
  - You are going to learn about some of the fundamental concepts of concurrency
  - We are going to learn about some of the fundamental concepts of concurrency
  - I’ll focus on how to use them at the right higher level

I hope this is a talk that you can revisit in the future and do a deep dive with the resources at the end

Concurrency is a massive topic--so let's start from the beginning

---

# What is Concurrency?

And the possible motivation to work concurrently

con·cur·rence

/kən'kərəns/ ⓘ

*noun*

noun: **concurrency**

1. the fact of two or more events or circumstances happening or existing at the same time.

# Human Psychology

- The reality is that we have to multitask:
  - We balance work, family, to-do lists, 8 half-finished hobby C++ programming projects, etc.
- Our brains (i.e., human cpus) however tend to perform best when we can focus on one task at a time
  - (e.g., I am unable to code and listen to music with any lyrics...unfortunately)

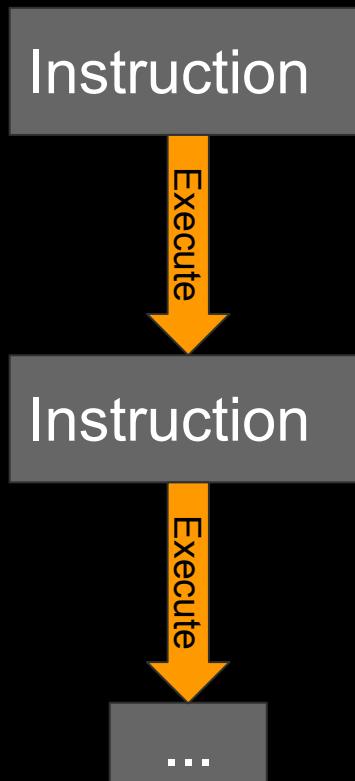
The screenshot shows the homepage of the Association for Psychological Science (APS). At the top is the APS logo, followed by a navigation bar with links for News, Research Topics, Conventions, Journals, and Observer Magazine. Below the navigation is a large, bold headline from LiveScience: "Why Humans Are Bad at Multitasking". Underneath the headline are the tags: COGNITIVE PROCESSES | COGNITIVE PSYCHOLOGY | MULTITASKING. The text "LiveScience:" is also present.

Thinking about all the *potential sequences* of events is difficult for humans:  
<https://www.psychologicalscience.org/news/why-humans-are-bad-at-multitasking.html>

# This somewhat reflects how we write software (1/3)

---

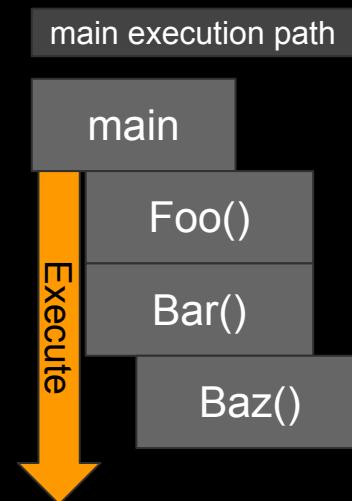
- Typically we have one main, sequential thread of execution
- One CPU core executes code sequentially
  - i.e. One instruction after the other.



# This somewhat reflects how we write software (2/3)

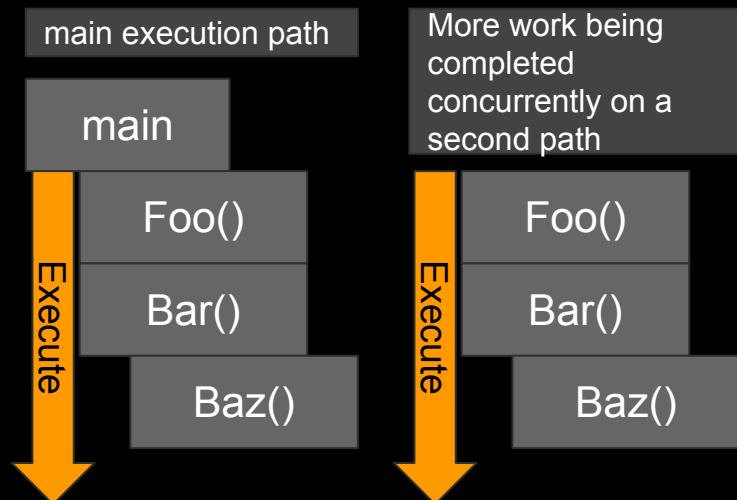
---

- Typically we have one main, sequential thread of execution
- One CPU core executes code sequentially
  - i.e. One instruction after the other.
  - We can abstract our visualization, and just show the call stack.
    - (One function calling the other, with the indentation indicating the call stack)



# This somewhat reflects how we write software (3/3)

- Typically we have one main, sequential thread of execution
- One CPU core executes code serially
  - i.e. One instruction after the other.
  - We can abstract our visualization, and just show the call stack.
    - (One function calling the other, with the indentation indicating the call stack)
  - But! What if we could have a second path in our code executing.

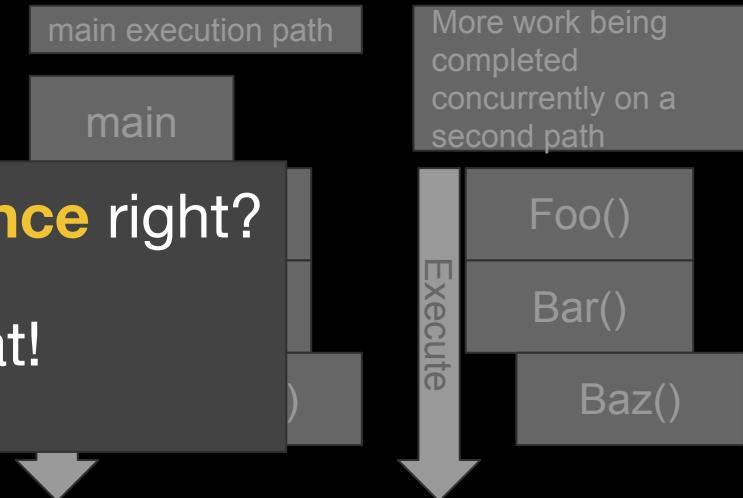


# This somewhat reflects how we write software (3/3)

- Typically we have one main, sequential thread of execution
- One CPU core executes
  - i.e. One instruction at a time
  - We can abstract away from the call stack.
    - (One function at a time)
  - indentation indicating the call stack)
  - But! What if we could have a second path in our code executing.

Twice the **performance** right?

Sounds great!

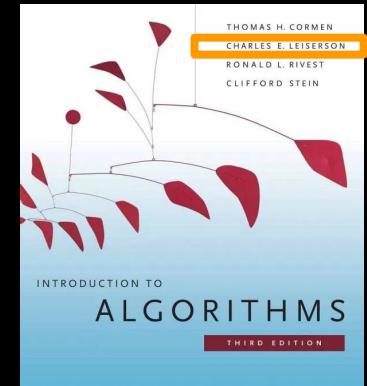


---

**Performance is the currency of computing.**

---

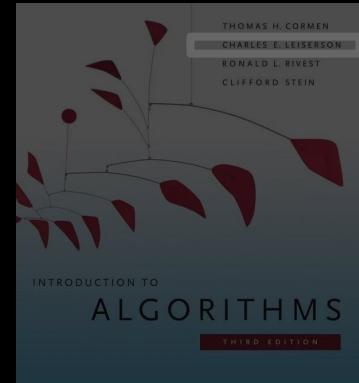
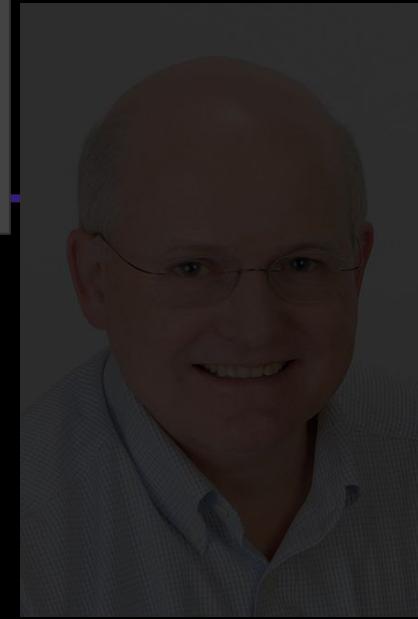
**“Performance is the currency of computing. You can often “buy” needed properties [of software] with performance”** - Charles Leiserson



So how can we get more performance? What are our tools?

Back to our original question--*what is concurrency* ?

**“Performance is the currency of computing.** You can often “buy” needed properties [of software] with performance” - Charles Leiserson



# Parallelism vs Concurrency (programming context) (1/4)

---

Concurrency is often used interchangeably with parallelism--so let's separate those two terms.

1. Concurrency Definition: Multiple things can happen at once, the order matters, and sometimes tasks have to wait on shared resources.
2. Parallelism Definition: Everything happens at once, instantaneously

# Parallelism vs Concurrency (programming context) (2/4)

---

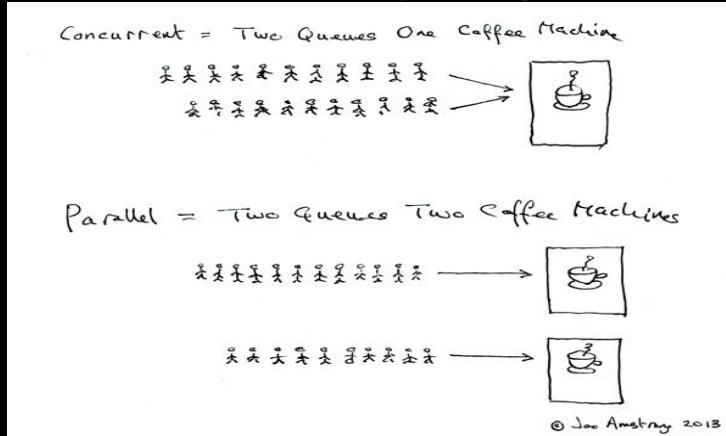
Concurrency is often used interchangeably with parallelism--so let's separate those two terms.

1. Concurrency Definition: Multiple things can happen at once, **the order matters, and sometimes tasks have to wait on shared resources.**
2. Parallelism Definition: Everything happens at once, instantaneously

# Parallelism vs Concurrency (programming context) (3/4)

Concurrency is often used interchangeably with parallelism--so let's separate those two terms.

1. Concurrency Definition: Multiple things can happen at once, **the order matters, and sometimes tasks have to wait on shared resources.**
2. Parallelism Definition: Everything happens at once, instantaneously

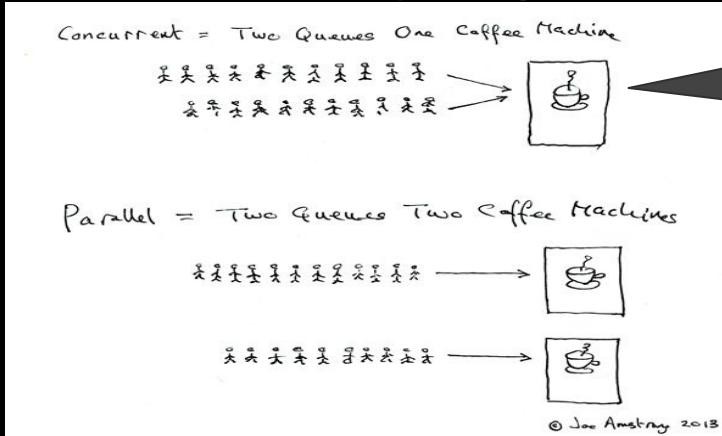


"The world is concurrent." - [Joe Armstrong](#)

# Parallelism vs Concurrency (programming context) (4/4)

Concurrency is often used interchangeably with parallelism--so let's separate those two terms.

1. Concurrency Definition: Multiple things can happen at once, **the order matters, and sometimes tasks have to wait on shared resources.**
2. Parallelism Definition: Everything happens at once, instantaneously



Concurrency and parallelism (implemented correctly) should yield better performance.

"The world is concurrent." - [Joe Armstrong](#)

# The Necessity of Concurrency\*

---

- In general, concurrency (like parallelism) is used because it is necessary for a system to function.
  - Real world concurrency examples
    - e.g. an orchestra, a subway transit system, cars at a traffic stop
  - Computer Science examples
    - e.g. A memory allocator, File I/O, Network requests (awaiting data)
- Again, concurrency is largely motivated by increased performance
  - *The potential* for more tasks to happen at once can thus increase performance (typically if we have multiple cores on our machine)

# Concurrency – as Heard in Music

---

- Ideally the instruments do not all play at once
- Ideally, at least some of the instruments are playing as well.
  - (as opposed to sitting silently)
- There is some notion of *synchronization* here.



# Good Concurrency = Good Conversation (1/3)

---



I don't know who these people are??

# Good Concurrency = Good Conversation (2/3)



Notice  
shared  
resources  
here as well

I don't know who these people are??

# Good Concurrency = Good Conversation (3/3)

Each person  
is *competing*  
for a  
resource



I don't know who these people are??

---

# Does our hardware support Concurrency?

Brief Architecture History

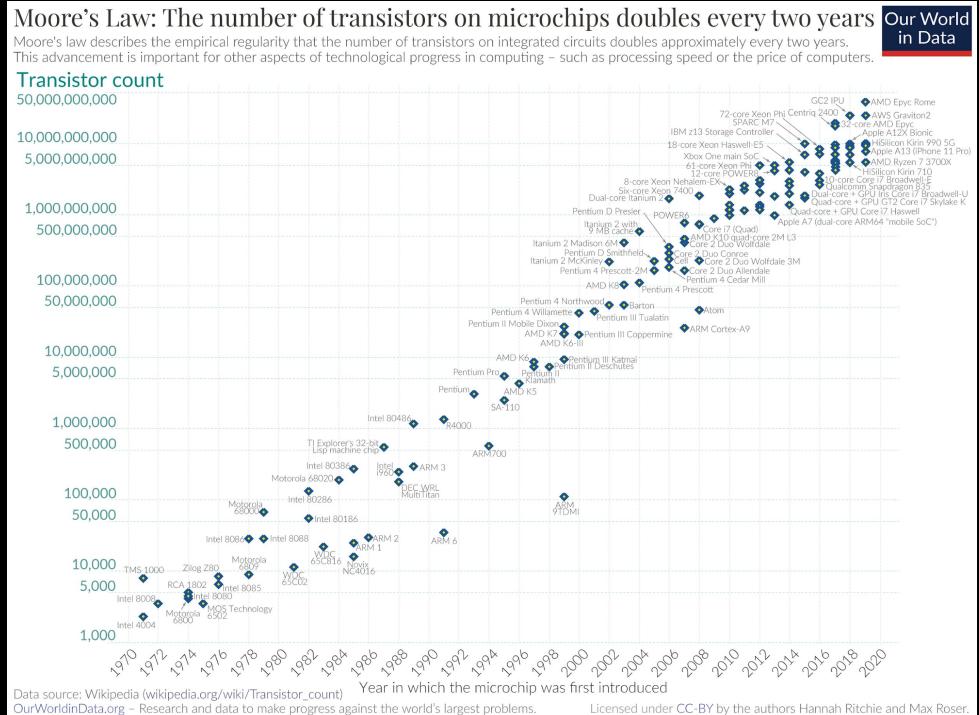
A long time ago in a galaxy far,  
far away....

*"The number of transistors incorporated in a chip will approximately double every 24 months."*

--Gordon Moore, Intel co-founder

# Moore's Law (1/2)

- Around 1965 Gordon Moore predicted the number of transistors would roughly double every 18-24 months on a single core
  - And largely this held true!

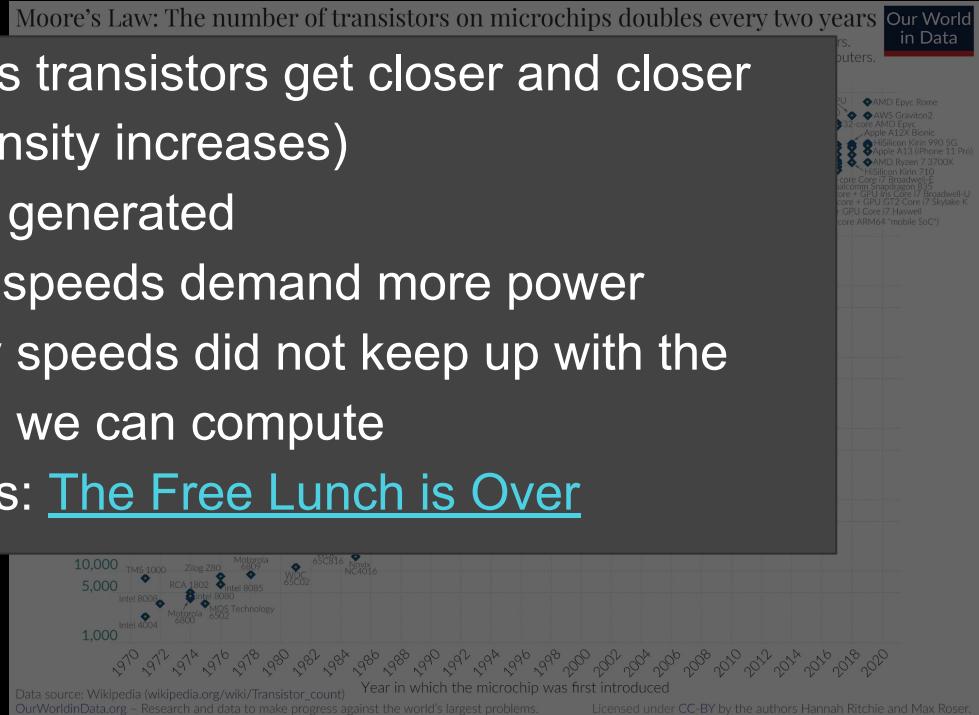


*"The number of transistors incorporated in a chip will approximately double every 24 months."*

--Gordon Moore, Intel co-founder

## Moore's Law (2/2)

- Around 1965 Gordon Moore predicted that the number of transistors on microchips would double every two years on a single chip
  - The problem is as transistors get closer and closer (i.e. transistor density increases)
    - More heat is generated
    - Faster clock speeds demand more power
    - And memory speeds did not keep up with the rate at which we can compute
  - And later...
- See Herb Sutter's: [The Free Lunch is Over](#)

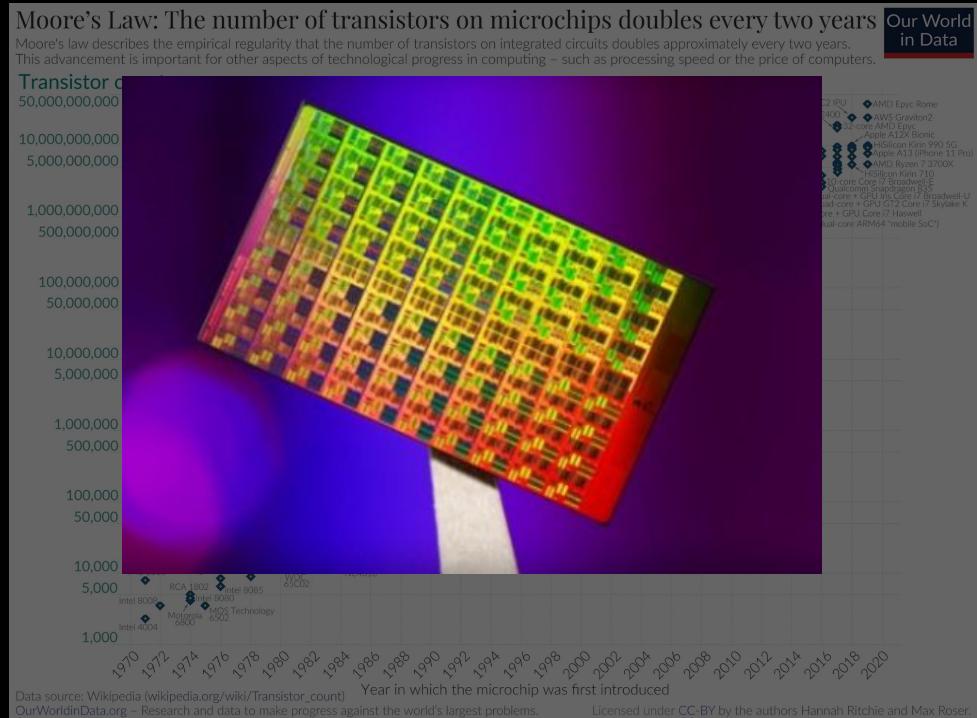


# Dennard Scaling (1/3)

*"The number of transistors incorporated in a chip will approximately double every 24 months."*

--Gordon Moore, Intel co-founder

- Physically (on the atomic scale) transistors are packed very tightly together
- Heat becomes a problem
- Energy consumption increases
  - (i.e. Dennard Scaling)



# Dennard Scaling (2/3)

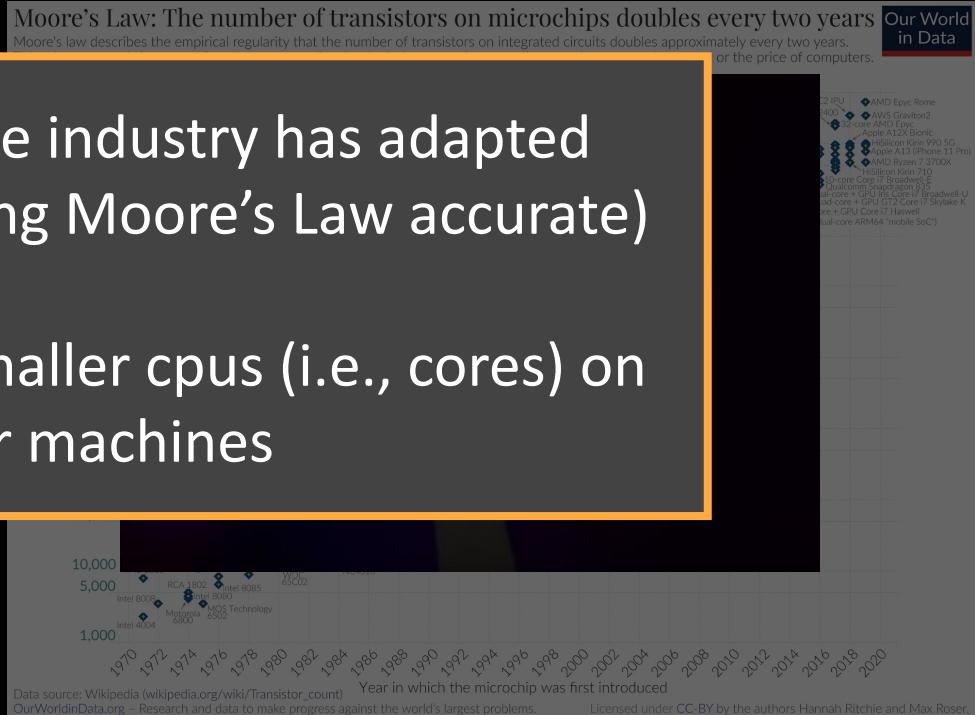
*"The number of transistors incorporated in a chip will approximately double every 24 months."*

--Gordon Moore, Intel co-founder

- Physically (as transistors are tightly together)
- Heat becomes a problem
- Energy consumption increases (i.e. Dennard scaling)

So the hardware industry has adapted  
(effectively keeping Moore's Law accurate)

We have more smaller cpus (i.e., cores) on our machines



# Dennard Scaling (3/3)

*"The number of transistors incorporated in a chip will approximately double every 24 months."*

--Gordon Moore, Intel co-founder

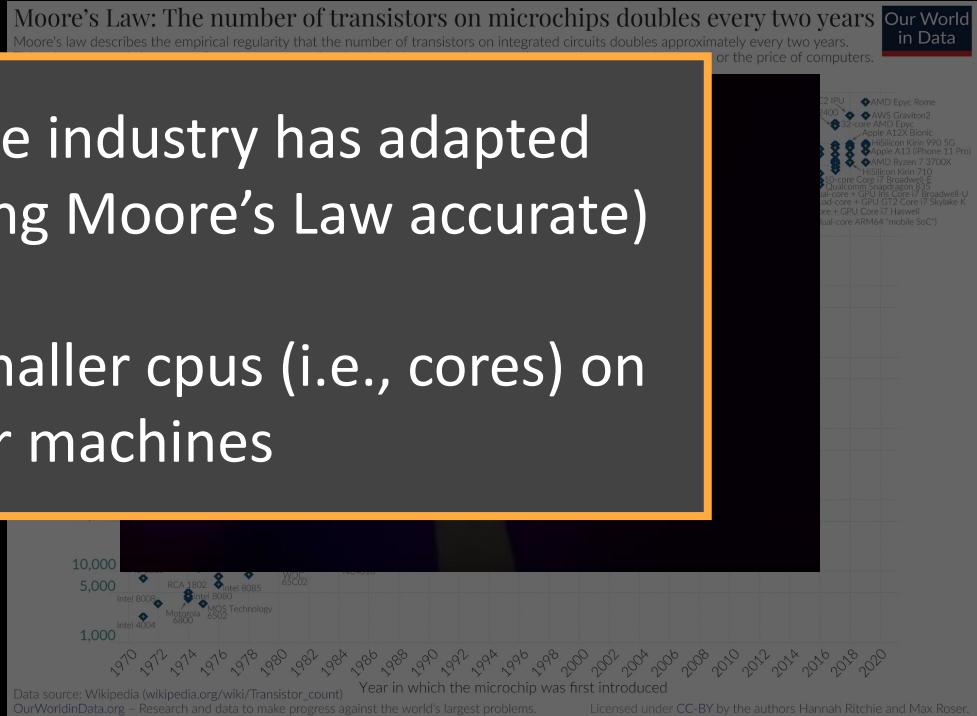
- Physically (as transistors are tightly together)
- Heat becomes a problem
- Energy consumption increases

So--does our hardware support concurrency?

Yes! Each core can be working on a separate task. (Note: A single core can also handle concurrency)

So the hardware industry has adapted (effectively keeping Moore's Law accurate)

We have more smaller cpus (i.e., cores) on our machines



# (Quick Examples) X-Box One

---

- Yes, the hardware has multiple CPUs
- Source: <http://www.videogames101.net/videogame.htm>

## Quick Specs

- CPU: AMD 1.75GHz 8 Core CPU
- GPU: AMD APU with 5GB dedicated VRAM
- Memory: 8GB DDR3 RAM
- Ports: 3 USB 3.0 ports, 1 Ethernet port, and 1 HDMI Output port
- Additional Storage: Upgradable internal hard drive storage and removable hard drive support.



# (Quick Examples) X-Box

- X-Box Series X
  - Also 8 cores

- Source:

<https://www.xbox.com/en-US/consoles/xbox-series-x#specs>



## TECH SPECS

### PROCESSOR

**CPU.** 8X Cores @ 3.8 GHz (3.66 GHz w/SMT) Custom Zen 2 CPU

**GPU.** 12 TFLOPS, 52 CUs @1.825 GHz Custom RDNA 2 GPU

**SOC Die Size.** 360.45 mm

**Process.** 7nm Enhanced

### MEMORY & STORAGE

**Memory.** 16GB GDDR6 w/320 bit-wide bus

**Memory Bandwidth.** 10GB @ 560 GB/s, 6GB @ 336 GB/s.

**Internal Storage.** 1TB Custom NVME SSD

**I/O Throughput.** 2.4 GB/s (Raw), 4.8 GB/s (Compressed, with custom hardware decompression block)

**Expandable Storage.** Support for 1TB Seagate Expansion Card for Xbox Series X|S matches internal storage exactly (sold separately). Support for USB 3.1 external HDD (sold separately).

### VIDEO CAPABILITIES

**Gaming Resolution.** True 4K

**High Dynamic Range.** Up to 8K HDR

**Optical Drive.** 4K UHD Blu-Ray

**Performance Target.** Up to 120 FPS

**HDMI Features.** Auto Low Latency Mode, HDMI Variable Refresh Rate, AMD FreeSync.

# (Quick Examples) PS5

- Yes, also utilizing multiple cores
- [https://en.wikipedia.org/wiki/PlayStation\\_5](https://en.wikipedia.org/wiki/PlayStation_5)



## CPU

Custom 8-core AMD Zen 2,  
variable frequency, up to 3.5 GHz

## Memory

16 GB GDDR6 SDRAM  
512 MB DDR4 RAM (for background tasks)<sup>[1]</sup>

## Storage

Custom 825 GB SSD

## Removable storage

Internal (user upgradeable) NVMe M.2 SSD, or external USB-based HDD

## Display

Video output formats  
HDMI: 720p, 1080i, 1080p, 4K UHD, 8K UHD

## Graphics

Custom AMD RDNA 2,  
36 CUs @ variable frequency up to 2.23 GHz

# More cores = The general trend in architecture.

---

- The idea of increasing cores to has existed for a long time looking back at early supercomputers
  - 1964 - Seymour Cray pioneers pipelining
  - 1972 - CDC 8600 Multiprocessor
  - 1975 - Cray-1 pioneered vector processing (i.e. SIMD)



---

# Now for code...(almost)

Let's see our first mechanism for concurrent programming

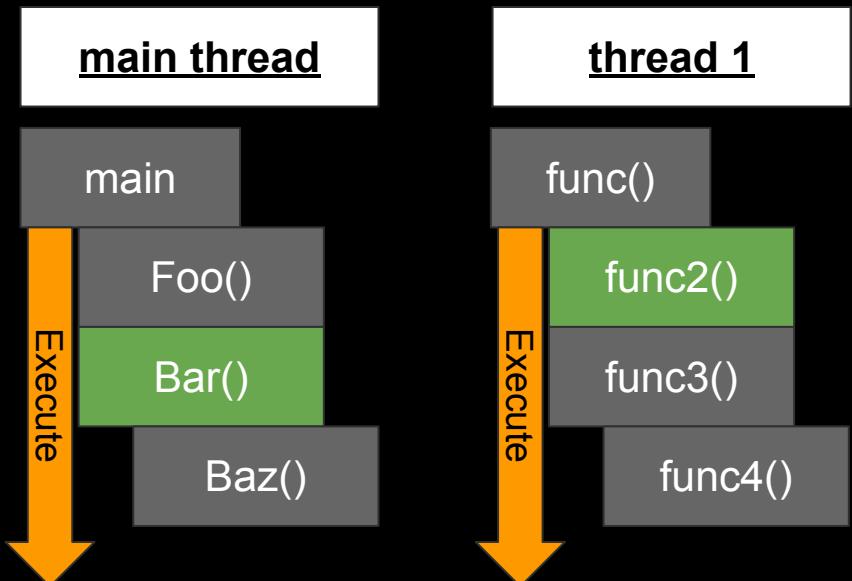


SHOW ME Concurrent Programming

Reference to American movie Jerry Maguire in which actor Tom Cruise yells "Show me the money" loudly into the phone to win a deal.

# Concurrency Mechanism - Thread

- One mechanism for achieving concurrency is a ‘thread’
  - A ‘thread’ allows us to execute two control flows at the same time
  - The ‘main thread’ is where our program starts
    - We may then have 1 or more additional threads:
      - executing a block of code
      - executing other functions
      - And overall--sharing the same code, and the same data
        - (all while our main thread also executes. )





# Threads

Often defined as a “lightweight process”

# High Level View of Thread

- 1 Process (i.e. your application) can have many threads:
  - Each thread shares the same code, data, and kernel context
  - A thread has its own thread id (TID)
  - A thread has its own logical control flow
  - A thread has its own stack for local variables

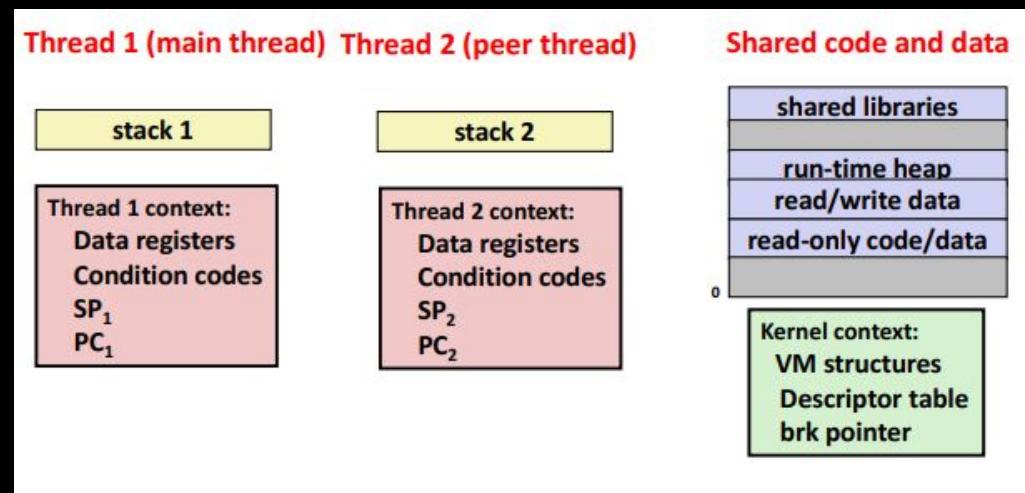


Figure from: Computer Systems a Programmer's Perspective 3rd Edition

# When to use threads

---

- Heavy Computations
  - Use threads to work on a heavy computation
    - The most common case is actually using threads on your GPU for graphics
    - GPUs have 100s or 1000s of threads that are good for massively parallel tasks.
      - (You could also use things like CUDA to take advantage of your graphics hardware)
  - You may need to use a series of threads to otherwise resolve complex computations on your CPU where decisions may need to be made.
- Using threads to separate work
  - Gives performance (Same as above)
  - But also simplifies the logic of your problem

---

# Threads in Modern C++

The `std::thread`

# Thread Libraries

---

- Before C++11/14/17/20, there existed threading libraries with different semantics
  - Libraries like “Boost”, Intel “Thread Building Blocks”, or “pthread” were used
    - Perhaps you have used pthread at least in C
    - (`std::thread` I believe is implemented with pthread most posix systems)
- Typically today *I would personally recommend* using the standard C++ threading library for portability reasons as the default choice.

#include <thread> | <http://en.cppreference.com/w/cpp/thread>

---

- With Modern C++
  - Here's what you get!
  - Let's start with 'thread'

C++ Thread support library

## Thread support library

C++ includes built-in support for threads, mutual exclusion, condition variables, and futures.

### Threads

Threads enable programs to execute across several processor cores.

Defined in header `<thread>`

<code>thread</code> (C++11)	manages a separate thread (class)
<code>jthread</code> (C++20)	<code>std::thread</code> with support for auto-joining and cancellation (class)

### Functions managing the current thread

Defined in namespace `this_thread`

<code>yield</code> (C++11)	suggests that the implementation reschedule execution of threads (function)
<code>get_id</code> (C++11)	returns the thread id of the current thread (function)
<code>sleep_for</code> (C++11)	stops the execution of the current thread for a specified time duration (function)
<code>sleep_until</code> (C++11)	stops the execution of the current thread until a specified time point (function)

# Thread Example - Launching a thread (1/2)

- `#include <thread>`
  - `std::thread`
- (Aside: For those familiar, this example is essentially going to do ‘fork-join’ parallelism)

```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"    << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     //           finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

## Threads

- `#include <thread>`
  - `std::thread`
- (Aside: For those familiar, this is essentially going to do ‘fork-join’ parallelism)

```
mike:concurrency$ g++ -std=c++17 thread2.cpp -o prog -lpthread
mike:concurrency$ ./prog
Hello from our thread!
Argument passed in:100
Hello from the main thread!
```

```
3 #include <iostream>
4 #include <thread> // Include the thread library
5
6 // Test function which we'll launch in a thread
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in: " << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     //           finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

Don't forget to link in the pthread library for posix users.

# Visual execution of “Hello Thread” (1/13)

```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"   << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     //           finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Visual execution of “Hello Thread” (2/13)

Main Thread

main() function where all C++ programs start.

We have 1 thread in our program (the main thread)

```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"   << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     //           finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Visual execution of “Hello Thread” (3/13)

Main Thread



We begin constructing  
std::thread  
myThread

```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"   << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     //           finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Visual execution of “Hello Thread” (4/13)

Main Thread



std::thread myThread

```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"   << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     //           finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Visual execution of “Hello Thread” (5/13)

Main Thread



```
std::thread myThread(&test, 100)
```

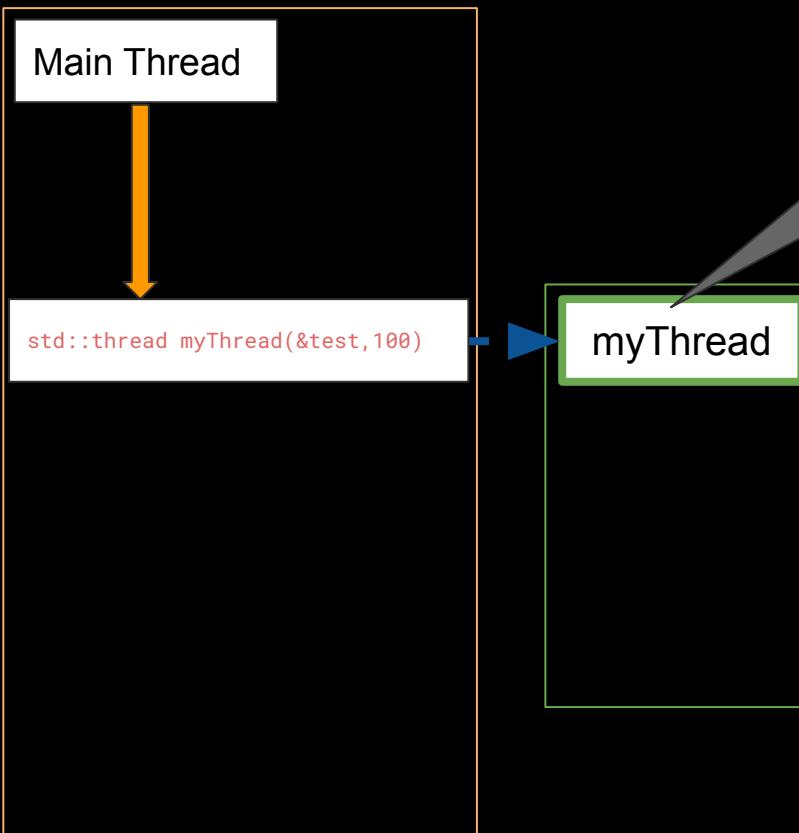
Our new thread will begin executing it's logical control flow from the ‘test’ function. *separately* from main()

The thread will start executing immediately on construction

(Remember, threads shares code and the heap)

```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"   << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100); // Join with the main thread, which is the same as
15     // saying "hey, main thread--wait until myThread
16     // finishes before executing further."
17     myThread.join();
18
19     // Continue executing the main thread
20     std::cout << "Hello from the main thread!" << std::endl;
21
22
23     return 0;
24 }
```

# Visual execution of “Hello from our thread!” (C/C++)



So now we have two  
“threads” executing

```
prog -lpthread
4 #include <thread> // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:" << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     // finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Visual execution

Both threads are executing concurrently!

(maybe on separate cores, or maybe on the same one)

Main Thread

```
std::thread myThread(&test, 100)
```



myThread

```
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"    << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread - wait until myThread
17     // finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

ad" (7/13)

# Visual exec

- We just happen to execute the next line in main thread
- myThread.join() tells this thread ('main') to wait on our other thread (tid) to finish.
  - We 'wait' in the main thread, because this is where we are calling join from

Main Thread

std::thread myThread(&test, 100)

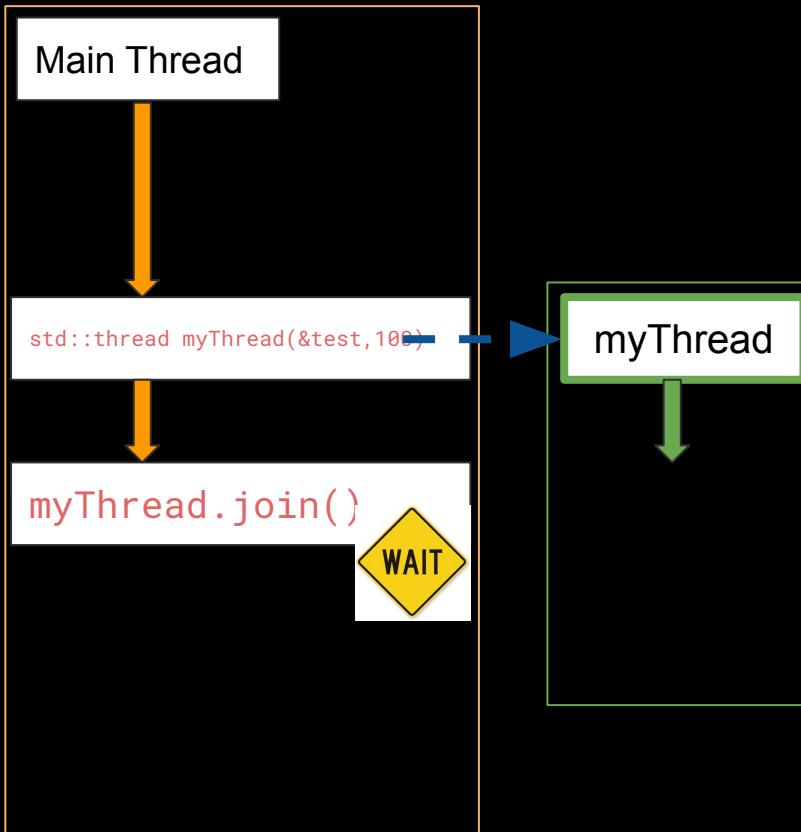
myThread.join()

myThread

```
g -lpthread
thread library

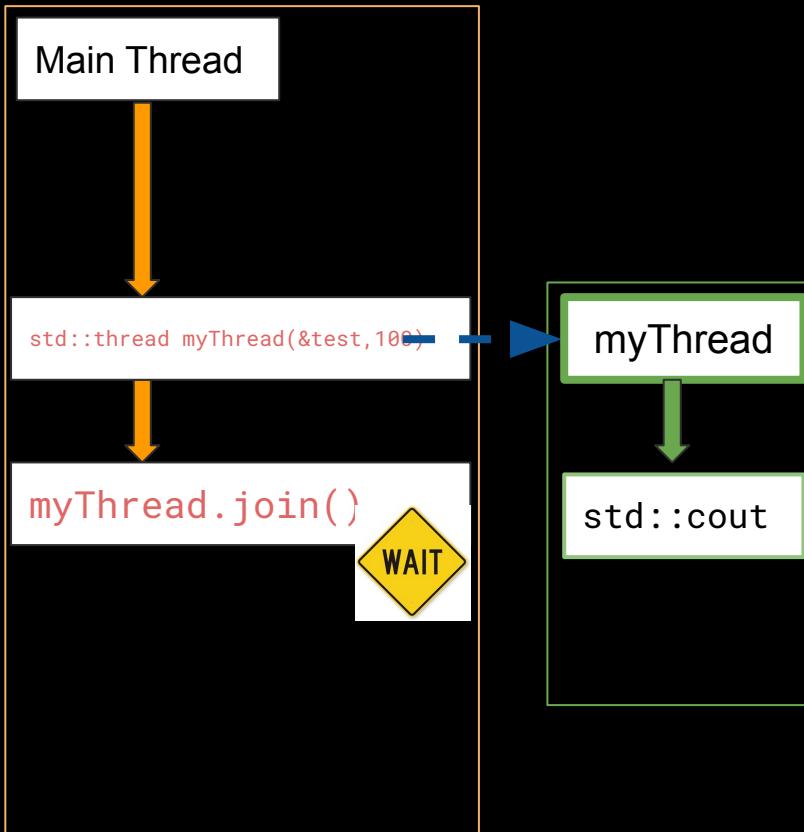
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"    << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     // finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Visual execution of “Hello Thread” (9/13)



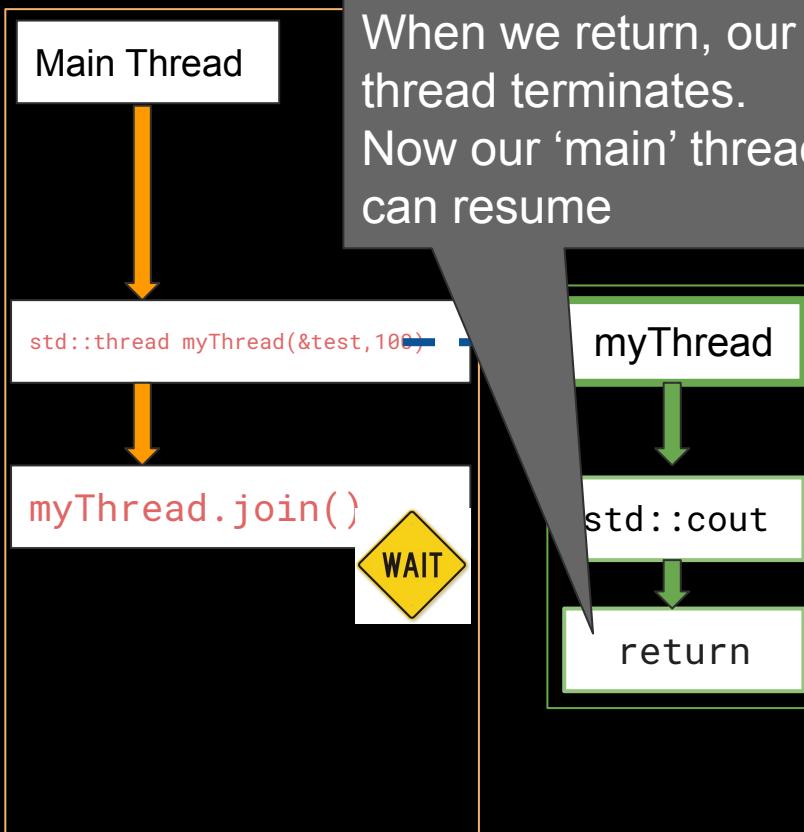
```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"    << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     // finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Visual execution of “Hello Thread” (10/13)



```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"    << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     // finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Visual execution of “Hello Thread” (11/13)



```
// @file thread1.cpp
// g++ -std=c++17 thread1.cpp -o prog -lpthread
#include <iostream>
#include <thread> // Include the thread library

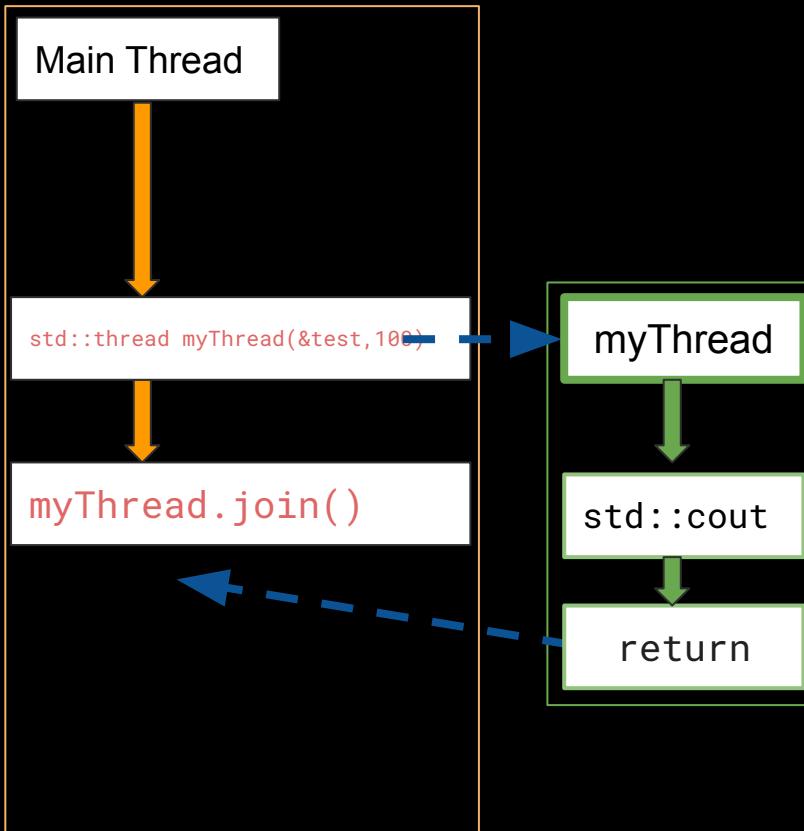
// Test function which we'll launch threads from
void test(int x) {
    std::cout << "Hello from our thread!" << std::endl;
    std::cout << "Argument passed in:" << x << std::endl;
}

int main() {
    // Create a new thread and pass one parameter
    std::thread myThread(&test, 100);
    // Join with the main thread, which is the same as
    // saying "hey, main thread--wait until myThread
    // finishes before executing further."
    myThread.join();

    // Continue executing the main thread
    std::cout << "Hello from the main thread!" << std::endl;

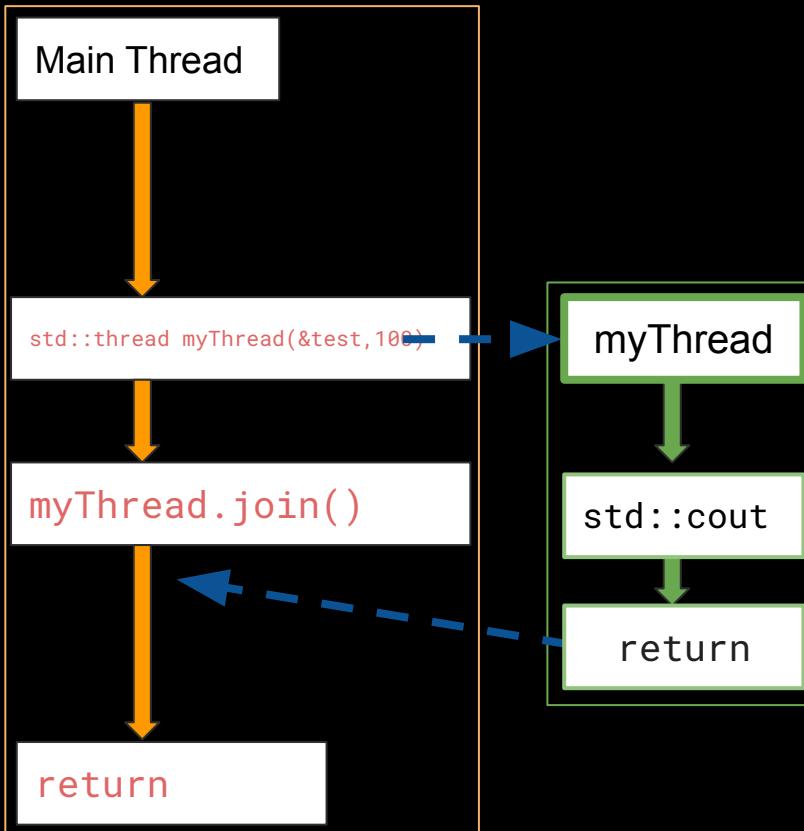
    return 0;
}
```

# Visual execution of “Hello Thread” (12/13)



```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"   << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     // finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Visual execution of “Hello Thread” (13/13)



```
1 // @file thread1.cpp
2 // g++ -std=c++17 thread1.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 // Test function which we'll launch threads from
7 void test(int x) {
8     std::cout << "Hello from our thread!" << std::endl;
9     std::cout << "Argument passed in:"    << x << std::endl;
10 }
11
12 int main() {
13     // Create a new thread and pass one parameter
14     std::thread myThread(&test, 100);
15     // Join with the main thread, which is the same as
16     // saying "hey, main thread--wait until myThread
17     // finishes before executing further."
18     myThread.join();
19
20     // Continue executing the main thread
21     std::cout << "Hello from the main thread!" << std::endl;
22
23     return 0;
24 }
```

# Same example as before -- but with a lambda!

- Same example as before, but instead of a function, I have a lambda with 1 parameter (and no return type)
  - std::thread takes a callable as the parameter--so lambdas, functions, etc. are all fine!

```
1 // @file thread2.cpp
2 // g++ -std=c++17 thread2.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5
6 int main() {
7
8     // This time create a lambda function
9     auto lambda = [](int x){
10         std::cout << "Hello from our thread!" << std::endl;
11         std::cout << "Argument passed in:" << x << std::endl;
12     };
13
14     // Create a new thread with our lambda this time
15     std::thread myThread(lambda,100);
16     // Join with the main thread, which is the same as
17     // saying "hey, main thread-wait until myThread
18     // finishes before executing further."
19     myThread.join();
20
21     // Continue executing the main thread
22     std::cout << "Hello from the main thread!" << std::endl;
23
24     return 0;
25 }
```

# Now how about if we wanted 10 threads

- Let's create a `std::vector<std::thread>`
  - Then we'll launch 10 threads from a loop
- It's important however, that we also join each of the threads!

```
1 // @file thread3.cpp
2 // g++ -std=c++17 thread3.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6
7 int main() {
8
9     // This time create a lambda function
10    auto lambda = [] (int x){
11        std::cout << "thread.get_id:" << std::this_thread::get_id() << std::endl;
12        std::cout << "Argument passed in:" << x << std::endl;
13    };
14
15    std::vector<std::thread> threads;
16    // Create a collection of threads
17    for (int i=0; i < 10; i++){
18        threads.push_back(std::thread(lambda, i));
19        threads[i].join();
20    }
21
22    // Continue executing the main thread
23    std::cout << "Hello from the main thread!" << std::endl;
24
25    return 0;
26 }
```

# Now how about if we wanted 10 threads (1/5)

- So here we create each of our threads and join them

```
1 // @file thread3.cpp
2 // g++ -std=c++17 thread3.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6
7 int main() {
8
9     // This time create a lambda function
10    auto lambda = [] (int x){
11        std::cout << "thread.get_id:" << std::this_thread::get_id() << std::endl;
12        std::cout << "Argument passed in:" << x << std::endl;
13    };
14
15    std::vector<std::thread> threads;
16    // Create a collection of threads
17    for (int i=0; i < 10; i++){
18        threads.push_back(std::thread(lambda, i));
19        threads[i].join();
20    }
21
22    // Continue executing the main thread
23    std::cout << "Hello from the main thread!" << std::endl;
24
25    return 0;
26 }
```

# Now how about if we wanted 10 threads (2/5)

- So here we create each of our threads and join them

```
mike:concurrency$ g++ -std=c++17 thread3.cpp -o prog -lpthread
mike:concurrency$ ./prog
thread.get_id:140658209871616
Argument passed in:0
thread.get_id:140658209871616
Argument passed in:1
thread.get_id:140658209871616
Argument passed in:2
thread.get_id:140658209871616
Argument passed in:3
thread.get_id:140658209871616
Argument passed in:4
thread.get_id:140658209871616
Argument passed in:5
thread.get_id:140658209871616
Argument passed in:6
thread.get_id:140658209871616
Argument passed in:7
thread.get_id:140658209871616
Argument passed in:8
thread.get_id:140658209871616
Argument passed in:9
Hello from the main thread!
```

```
1 // @file thread3.cpp
2 // g++ -std=c++17 thread3.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6
7 int main() {
8
9     // This time create a lambda function
10    auto lambda = [] (int x){
11        std::cout << "thread.get_id:" << std::this_thread::get_id() << std::endl;
12        std::cout << "Argument passed in:" << x << std::endl;
13    };
14
15    std::vector<std::thread> threads;
16    // Create a collection of threads
17    for (int i=0; i < 10; i++){
18        threads.push_back(std::thread(lambda, i));
19        threads[i].join();
20    }
21
22    // Continue executing the main thread
23    std::cout << "Hello from the main thread!" << std::endl;
24
25    return 0;
26 }
```

# Now how about if we

- So here we create each of our threads and join them

```
mike:concurrency$ g++ -std=c++17 thread3.cpp -lpthread  
mike:concurrency$ ./prog  
thread.get_id:140658209871616  
Argument passed in:0  
thread.get_id:140658209871616  
Argument passed in:1  
thread.get_id:140658209871616  
Argument passed in:2  
thread.get_id:140658209871616  
Argument passed in:3  
thread.get_id:140658209871616  
Argument passed in:4  
thread.get_id:140658209871616  
Argument passed in:5  
thread.get_id:140658209871616  
Argument passed in:6  
thread.get_id:140658209871616  
Argument passed in:7  
thread.get_id:140658209871616  
Argument passed in:8  
thread.get_id:140658209871616  
Argument passed in:9  
Hello from the main thread!
```

The result seems a little strange...anyone see the problem?

```
8 // This time create a lambda function  
9 auto lambda = [](int x){  
10     std::cout << "thread.get_id:" << std::this_thread::get_id() << std::endl;  
11     std::cout << "Argument passed in:" << x << std::endl;  
12 };  
13  
14 std::vector<std::thread> threads;  
15 // Create a collection of threads  
16 for(int i=0; i < 10; i++){  
17     threads.push_back(std::thread(lambda,i));  
18     threads[i].join();  
19 }  
20  
21 // Continue executing the main thread  
22 std::cout << "Hello from the main thread!" << std::endl;  
23  
24 return 0;  
25 }  
26 }
```

# Now how about if we

- So here we create each of our threads and join them

```
mike:concurrency$ g++ -std=c++17 thread3.cpp -o prog -lpthread
mike:concurrency$ ./prog
thread.get_id:140658209871616
Argument passed in:0
thread.get_id:140658209871616
Argument passed in:1
thread.get_id:140658209871616
Argument passed in:2
thread.get_id:140658209871616
Argument passed in:3
thread.get_id:140658209871616
Argument passed in:4
thread.get_id:140658209871616
Argument passed in:5
thread.get_id:140658209871616
Argument passed in:6
thread.get_id:140658209871616
Argument passed in:7
thread.get_id:140658209871616
Argument passed in:8
thread.get_id:140658209871616
Argument passed in:9
Hello from the main thread!
```

By joining our threads immediately after launching our code, we've effectively made our program sequential (i.e. no performance gain)

```
8
9     // Create a lambda function
10    auto lambda = [=]{ 
11        std::cout << "thread.get_id:" << std::this_thread::get_id() << std::endl;
12        std::cout << "Argument passed in:" << x << std::endl;
13    };
14
15    std::vector<std::thread> threads;
16    // Create a collection of threads
17    for(int i=0; i < 10; i++){
18        threads.push_back(std::thread(lambda,i));
19        threads[i].join();
20    }
21
22    // Continue executing the main thread
23    std::cout << "Hello from the main thread!" << std::endl;
24
25    return 0;
26 }
```

# Now how about if we wanted 1

- So here we create each of our threads and join them

```
mike:concurrency$ g++ -std=c++17 thread3_fix.cpp -o prog -lpthread
mike:concurrency$ ./prog
thread.get_id:139995667298048
Argument passed in:0
thread.get_id:139995507902208
Argument passed in:3
thread.get_id:thread.get_id:139995642119936
Argument passed in:4
thread.get_id:139995633727232
Argument passed in:5
139995650512640
Argument passed in:2
thread.get_id:139995658905344
Argument passed in:1
thread.get_id:139995608549120
Argument passed in:8
thread.get_id:139995532752640
Argument passed in:9
thread.get_id:139995616941824
Argument passed in:7
thread.get_id:139995625334528
Argument passed in:6
Hello from the main thread!
```

```
1 // @file thread3_fix.cpp
2 // g++ -std=c++17 thread3_fix.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // In <thread> header
5 #include <vector>
6
7 int main() {
8
9     // This time we're using a lambda function
10    auto lambda = [] (int x) {
11        std::cout << "thread.get_id:" << std::this_thread::get_id() << std::endl;
12        std::cout << "Argument passed in:" << x << std::endl;
13    };
14
15    std::vector<std::thread> threads;
16    // Create a collection of threads
17    for(int i=0; i < 10; i++){
18        threads.push_back(std::thread(lambda,i));
19    }
20    // Join all of our threads here--
21    // one or more may have launched, but we'll have
22    // to wait in main until ALL threads finish.
23    for(int i=0; i < 10; i++){
24        threads[i].join();
25    }
26
27    // Continue executing the main thread
28    std::cout << "Hello from the main thread!" << std::endl;
29
30    return 0;
31 }
```

Here's the fix and new output

Observe, the thread execution is out of order now (which is expected when 10 threads are simultaneously executed, the threads are scheduled according to OS)

# C++ 20 - std::jthread

- std::jthread launches a thread and joins the thread on destruction
  - This may be more useful (especially for beginners) as we don't forget to join!
    - If you need more control on when to join, then prefer std::thread and join explicitly
  - (Note: This code does the right thing--threads are immediately launched and not sequentially waited upon)

```
1 // @file thread4.cpp
2 // g++-10 -std=c++20 thread4.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6
7 int main() {
8
9     // This time create a lambda function
10    auto lambda = [](int x){
11        std::cout << "thread.get_id:" << std::this_thread::get_id() << std::endl;
12        std::cout << "Argument passed in:" << x << std::endl;
13    };
14
15    // Note: We now have a jthread
16    //       No joins in the program
17    std::vector<std::jthread> threads;
18    // Create a collection of threads
19    for(int i=0; i < 10; i++){
20        threads.push_back(std::jthread(lambda,i));
21    }
22
23    // Continue executing the main thread
24    std::cout << "Hello from the main thread!" << std::endl;
25
26    return 0;
27 }
```

---

# What about threads working together?

How do threads work with data?

# 1000 threads working together (1/4)

- Assume you have some shared task:
  - Perhaps 1000 threads counting the sizes of files in a directory to find the total bytes of the files in a directory
    - i.e., a parallel wc tool
  - Each of those threads then needs to write to some ‘shared value’ to sum the total.
- On the example on the right, let’s simplify and just have 1000 threads work together to increment a value

```
1 // @file thread5.cpp
2 // g++ -std=c++17 thread5.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6
7 // Some shared value
8 // A little ugly that it's a static global but it's a toy example.
9 static int shared_value = 0;
10
11 void increment_shared_value(){
12     shared_value = shared_value + 1;
13 }
14
15 int main() {
16     std::vector<std::thread> threads;
17     // Create a collection of threads
18     for(int i=0; i < 1000; i++){
19         threads.push_back(std::thread(increment_shared_value));
20     }
21     // Join our threads
22     for(int i=0; i < 1000; i++){
23         threads[i].join();
24     }
25     // Retrieve our result
26     std::cout << "Result = " << shared_value << std::endl;
27
28     return 0;
29 }
```

# 1000 threads working together (2/4)

- Assume you have some shared task:
  - Perhaps 1 sizes of file total bytes
    - i.e.,
  - Each of them write to some shared value to sum the total.
- On the example on the right, let's simplify and just have 1000 threads work together to increment a value

Here we launch  
and join our 1000  
threads

```
1 // @file thread5.cpp
2 // g++ -std=c++17 thread5.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6
7 // Some shared value
8 // A little ugly that it's a static global but it's a toy example.
9 static int shared_value = 0;
10
11 void increment_shared_value(){
12     shared_value = shared_value + 1;
13 }
14
15 int main() {
16     std::vector<std::thread> threads;
17     // Create a collection of threads
18     for(int i=0; i < 1000; i++){
19         threads.push_back(std::thread(increment_shared_value));
20     }
21     // Join our threads
22     for(int i=0; i < 1000; i++){
23         threads[i].join();
24     }
25     // Retrieve our result
26     std::cout << "Result = " << shared_value << std::endl;
27
28     return 0;
29 }
```

# 1000 threads working together (3/4)

- Assume you have some shared task:
  - Each thread calls 'increment\_shared\_value' and adds 1 to the shared\_value
  - write to some shared value to sum the total.
- On the example on the right, let's simplify and just have 1000 threads work together to increment a value

```
1 // @file thread5.cpp
2 // g++ -std=c++17 thread5.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6
7 // Some shared value
8 // A little ugly that it's a static global but it's a toy example.
9 static int shared_value = 0;
10
11 void increment_shared_value(){
12     shared_value = shared_value + 1;
13 }
14
15 int main() {
16     std::vector<std::thread> threads;
17     // Create a collection of threads
18     for(int i=0; i < 1000; i++){
19         threads.push_back(std::thread[increment_shared_value]);
20     }
21     // Join our threads
22     for(int i=0; i < 1000; i++){
23         threads[i].join();
24     }
25     // Retrieve our result
26     std::cout << "Result = " << shared_value << std::endl;
27
28     return 0;
29 }
```

# 1000 threads working together (4/4)

- Assume tasks  
It's not deterministic!

(We probably want the correct answer more often than 2/6 tries...)

## Why?

- i.e., a parallel wc tool
- Each of those threads write to some total.
- On the example, simplify and just work together to

```
thread5.cpp
cd=c++17 thread5.cpp -o prog -lpthread
#include <iostream>
#include <thread> // Include the thread library
#include <vector>

shared_value
Le ugly that it's a static global but it's a toy example.
static shared_value = 0;

void increment_shared_value(){
    shared_value = shared_value + 1;
}

int main(int argc, char* argv[])
{
    if(argc != 2)
    {
        std::cout << "Usage: " << argv[0] << " <file>" << std::endl;
        return 1;
    }

    std::vector<std::thread> threads;
    for(int i = 0; i < 1000; ++i)
    {
        threads.push_back(std::thread{increment_shared_value});
    }

    for(auto& t : threads)
    {
        t.join();
    }

    std::cout << "Result = " << shared_value << std::endl;
}
```

```
mike:concurrency$ g++-10 -std=c++17 thread5.cpp -o prog -lpthread
mike:concurrency$ ./prog
Result = 998
mike:concurrency$ ./prog
Result = 1000
mike:concurrency$ ./prog
Result = 1000
mike:concurrency$ ./prog
Result = 997
mike:concurrency$ ./prog
Result = 998
mike:concurrency$ ./prog
Result = 996
```

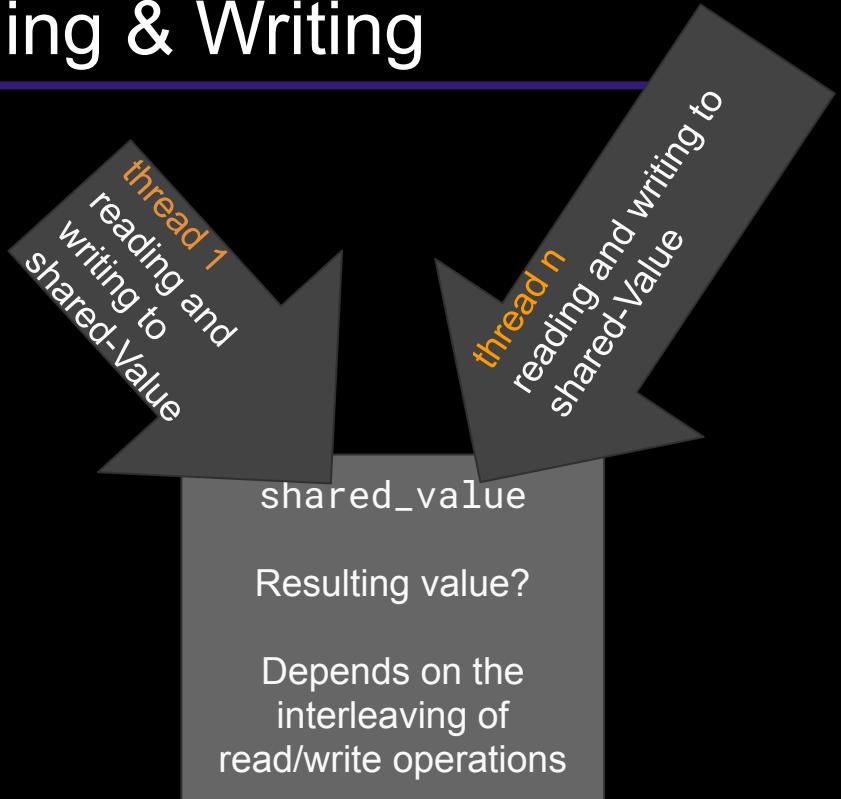
---

# Pitfalls of Concurrent Programming (Specifically with threads)

Data Races, Starvation, and Deadlock

# Problem with Threads -- Reading & Writing

- Data Race (or race condition)
  - Because data is shared--one or more thread could be writing to the same piece of memory at the same time
    - One thread may have read a ‘stale’ value right before the new ‘write’ to the value
      - The thread that then writes will update +1 to a stale value, overwriting the other threads update
  - This makes the operation non-deterministic
    - i.e. We may get unexpected or undefined results regarding the final value based on a non-deterministic order of operations



# Solving a Data Race

---

- Fixing Data Races
  - We can use a ‘lock’ to protect data, so that only one thread at a time can thus access the memory.
  - In C++, we have these available, and they are known as ‘[mutexes](#)’
- A mutex, allows ‘mutual exclusion’ to a block of code.
  - Thus, the operation is ‘atomic’ in the sense that only 1 operation can happen while the lock is held.
  - *Analogy:*
    - *Think about having exactly 1 key to your home, and you always carry the key with you.*
    - *Only the person who has the key can access the house.*
    - *When the person enters, they lock the door*
    - *When the person leaves, they can pass on the key to someone else to enter, who will also lock the door when they enter.*

# std::mutex (1/2)

- Four new lines of code added
  - The mutex library
  - A global lock
  - A lock and unlock call on our global lock

```
1 // @file thread6.cpp
2 // g++ -std=c++17 thread6.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 # include <mutex> // New library
7
8 // Some shared value
9 static int shared_value = 0;
10 std::mutex gLock; // A global lock
11
12 void increment_shared_value(){
13     gLock.lock();
14     shared_value = shared_value + 1;
15     gLock.unlock();
16 }
17
18 int main() {
19     std::vector<std::thread> threads;
20     // Create a collection of threads
21     for(int i=0; i < 1000; i++){
22         threads.push_back(std::thread(increment_shared_value));
23     }
24     // Join our threads
25     for(int i=0; i < 1000; i++){
26         threads[i].join();
27     }
28     // Retrieve our result
29     std::cout << "Result = " << shared_value << std::endl;
30
31     return 0;
32 }
```

# std::mutex (2/2)

- Four new lines of code added
  - The mutex library
  - A global lock
  - A lock and unlock call on our global lock
- Results look good to me!

```
mike:concurrency$ g++-10 -std=c++17 thread6.cpp -o prog -lpthread
mike:concurrency$ ./prog
Result = 1000
```

```
1 // @file thread6.cpp
2 // g++ -std=c++17 thread6.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread> // Include the thread library
5 #include <vector>
6 #include <mutex> // New library
7
8 // Some shared value
9 static int shared_value = 0;
10 std::mutex gLock; // A global lock
11
12 void increment_shared_value(){
13     gLock.lock();
14     shared_value = shared_value + 1;
15     gLock.unlock();
16 }
17
18 int main() {
19     std::vector<std::thread> threads;
20     // Create a collection of threads
21     for(int i=0; i < 1000; i++){
22         threads.push_back(std::thread(increment_shared_value));
23     }
24     // Join our threads
25     for(int i=0; i < 1000; i++){
26         threads[i].join();
27     }
28     // Retrieve our result
29     std::cout << "Result = " << shared_value << std::endl;
30
31     return 0;
32 }
```

# std::mutex and mutual exclusion (1/4)

- So what our lock is doing is providing access to only one thread at a time (mutually exclusive access).

```
12 void increment_shared_value(){  
13     gLock.lock();  
14     shared_value = shared_value + 1;  
15     gLock.unlock();  
16 }
```

# std::mutex and mutual exclusion (2/4)

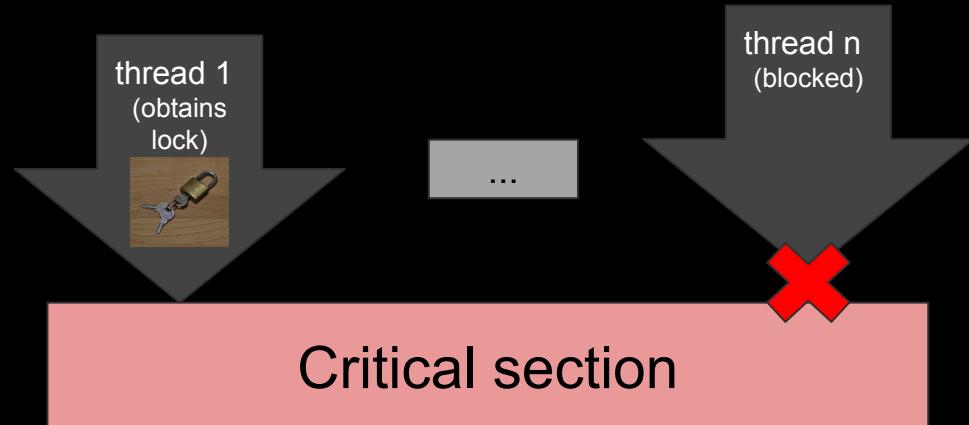
- So what our lock is doing is providing access to only one thread at a time (mutually exclusive access).
  - This region is called the ‘critical section’ that is protected by the lock.
  - Critical because we only want one thread at a time to enter and modify the shared state in the program.

```
12 void increment_shared_value(){  
13     gLock.lock();  
14     shared_value = shared_value + 1;  
15     gLock.unlock();  
16 }
```

# std::mutex and mutual exclusion (3/4)

- So what our lock is doing is providing access to only one thread at a time (mutually exclusive access).
  - This region is called the ‘critical section’ that is protected by the lock.
  - Critical because we only want one thread at a time to enter and modify the shared state in the program.

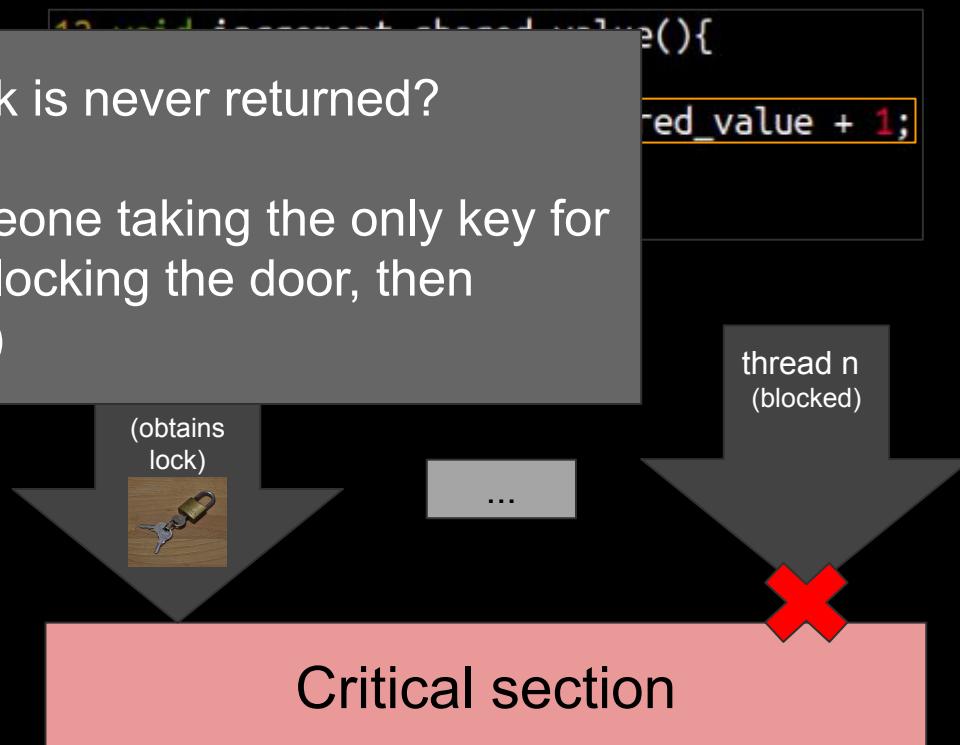
```
12 void increment_shared_value(){  
13     gLock.lock();  
14     shared_value = shared_value + 1;  
15     gLock.unlock();  
16 }
```



Critical section

# std::mutex and mutual exclusion (4/4)

- So what happens if a thread providing mutual exclusion never returns the lock? (e.g., equivalent to someone taking the only key for your house, walking in, locking the door, then flushing key down toilet)
  - This is a **deadlock**.
  - Critical because we only want one thread at a time to enter and modify the shared state in the program.



Critical section

# Deadlock - lack of *any* progress for a thread (1/2)

- Deadlock
  - Is the prevention of a thread from ever acquiring a resource
    - Thus, no forward progress can be made (the thread waits forever)
  - This typically happens when a thread does not release a lock, and goes out of scope or otherwise terminates before releasing the lock

```
1 // @file thread7_deadlock.cpp
2 // g++ -std=c++17 thread7_deadlock.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <mutex>    // New library
7
8 // Some shared value
9 static int shared_value = 0;
10 std::mutex gLock; // A global lock
11
12 void increment_shared_value(){
13     gLock.lock();
14     shared_value = shared_value + 1;
15 //     gLock.unlock(); // Oops, never return lock
16 }
```

# Deadlock - lack of *any* progress for a thread (2/2)

---

- Fixing deadlock
  - Re-run code, and see if you are missing a pair of lock/unlock
  - Static analysis techniques (i.e. thread sanitizers) may detect deadlock before compilation.
  - Otherwise deadlock has to be carefully detected at run-time and fixed.
- Note: Deadlock is the most extreme form of starvation
  - Starvation is when a thread cannot fairly acquire access to a resource

```
1 // @file thread7_deadlock.cpp
2 // g++ -std=c++17 thread7_deadlock.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <mutex>    // New library
7
8 // Some shared value
9 static int shared_value = 0;
10 std::mutex gLock; // A global lock
11
12 void increment_shared_value(){
13     gLock.lock();
14     shared_value = shared_value + 1;
15 //     gLock.unlock(); // Oops, never return lock
16 }
```

# Careful with std::mutex (1/2)

- So let's make sure we have a lock for every unlock
  - Our code is fixed right?
  - (I agree this looks correct)
- The problem is if another programmer comes and updates line 14

```
1 // @file thread7_deadlock.cpp
2 // g++ -std=c++17 thread7_deadlock.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <mutex>   // New library
7
8 // Some shared value
9 static int shared_value = 0;
10 std::mutex gLock; // A global lock
11
12 void increment_shared_value(){
13     gLock.lock();
14     shared_value = shared_value + 1;
15     gLock.unlock(); // Fixed right?
16 }
```

## Careful with std::mutex (2/2)

- So let's make sure we have a lock for every unlock
  - Our code is fixed right?
  - (I agree this looks correct)
- The problem is if another programmer comes and

```
1 // @file thread7_deadlock.cpp
2 // g++ -std=c++17 thread7_deadlock.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <mutex>     // New library
7
8 // Some shared value
9 static int shared_value = 0;
10 std::mutex gLock; // A global lock
11
12 void increment_shared_value(){
13     gLock.lock();
14     shared_value = shared_value + 1;
15     gLock.unlock(); // Fixed right?
16 }
```

```
12 void increment_shared_value(){
13     gLock.lock();
14     try{
15         shared_value = shared_value + 1;
16         throw "Dangerous exception abort";
17     }catch(...){
18         std::cout << "handle exception by returning from thread\n";
19         return;
20     }
21     gLock.unlock(); // Oops, never return lock
22 }
```

Maybe our object can throw an exception, or a programmer updates to the following

So this code will also deadlock! Consider the more complex case where some ‘exception’ is thrown and we ‘forget’ to also release the lock in catch.

You *could still* remember to use a lock, but we have a better tool

## Prefer lock\_guard (C++11) over lock/unlock (1/2)

- We instead of a lockGuard that can ‘wrap’ an individual std::mutex
    - The destructor of lock\_guard will take care of releasing the lock
  - std::lock\_guard is a good example of RAII
    - lock\_guard takes ownership of the lock, and when we leave scope the mutex is released (and the lock\_guard destroyed)

```
1 // @file thread9.cpp
2 // g++ -std=c++17 thread9.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <mutex>   // New library
7
8 // Some shared value
9 static int shared_value = 0;
10 std::mutex gLock; // A global lock
11
12 void increment_shared_value(){
13     // lock_guard follows RAII principles and will
14     // release lock after leaving scope.
15     // This includes if an exception is thrown.
16     std::lock_guard<std::mutex> lockGuard(gLock);
17     try{
18         shared_value = shared_value + 1;
19         throw "Dangerous exception abort";
20     }catch(...){
21         std::cout << "handle exception by returning from thread\n";
22         return;
23     }
24 }
```

# Prefer lock\_guard (C++11) over lock/unlock (2/2)

- We instead of a lockGuard that can ‘wrap’ an individual std::mutex
  - The destructor of lock\_guard will take care of releasing the lock
- std::lock\_guard is a good example of

```
1 // @file thread9.cpp
2 // g++ -std=c++17 thread9.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <mutex>    // New library
7
8 // Some shared value
9 static int shared_value = 0;
10 std::mutex gLock; // A global lock
11
12 void increment_shared_value(){
13     // lock_guard follows RAII principles and will
14     // release lock after leaving scope.
15     // This includes if an exception is thrown.
16     std::lock_guard<std::mutex> lockGuard(gLock);
17     try{
18         shared_value = shared_value + 1;
19         throw "Dangerous exception abort";
20     }
21 }
```

std::lock\_guard is only 3 member functions

## Member functions

(constructor)	constructs a <b>lock_guard</b> , optionally locking the given mutex <b>(public member function)</b>
(destructor)	destructs the <b>lock_guard</b> object, unlocks the underlying mutex <b>(public member function)</b>
<b>operator=</b> [deleted]	not copy assignable <b>(public member function)</b>

# std::scoped\_lock - Other mechanisms

- std::scoped\_lock (C++17) -
  - An update to lock\_guard, but can acquire multiple locks at once
    - i.e., `std::scoped_lock scoped_lock(mutex1, mutex2);`
    - Prefer scoped\_lock (over lock\_guard) if you are able to utilize C++17.

## Member functions

(constructor)	constructs a scoped_lock, optionally locking the given mutexes <small>(public member function)</small>
(destructor)	destructs the scoped_lock object, unlocks the underlying mutexes <small>(public member function)</small>
<b>operator=</b> [deleted]	not copy-assignable <small>(public member function)</small>

# std::atomic (1/3)

- For the problem we were trying to solve (incrementing a shared\_value with multiple threads), we could actually use atomics.

- atomics

```
1 // @file atomics.cpp
2 // g++ -std=c++17 atomics.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <atomic>   // New library
7
8 // Some shared value
9 static std::atomic<int> shared_value=0;
10
11 void increment_shared_value(){
12     shared_value++;
13 }
14
15 int main() {
16     std::vector<std::thread> threads;
17     // Create a collection of threads
18     for(int i=0; i < 1000; i++){
19         threads.push_back(std::thread(increment_shared_value));
20     }
21     // Join our threads using a ranged-based loop
22     // (Our intent is to iterate through all threads in container)
23     for(auto& th: threads){
24         threads.join();
25     }
26     // Retrieve our result
27     std::cout << "Result = " << shared_value << std::endl;
28
29     return 0;
30 }
```

# std::atomic (2/3)

- For the problem we were trying to solve (incrementing a shared\_value with multiple threads), we could actually use atomics.

- [atomics](#)

```
1 // @file atomics.cpp
2 // g++ -std=c++17 atomics.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <atomic>    // New library
```

## std::atomic

Defined in header `<atomic>`

<code>template&lt; class T &gt;</code>	(1) <small>(since C++11)</small>
<code>struct atomic;</code>	
<code>template&lt; class U &gt;</code>	(2) <small>(since C++11)</small>
<code>struct atomic&lt;U*&gt;;</code>	
Defined in header <code>&lt;memory&gt;</code>	
<code>template&lt; class U &gt;</code>	(3) <small>(since C++20)</small>
<code>struct atomic&lt;std::shared_ptr&lt;U&gt;&gt;;</code>	
<code>template&lt; class U &gt;</code>	(4) <small>(since C++20)</small>
<code>struct atomic&lt;std::weak_ptr&lt;U&gt;&gt;;</code>	
Defined in header <code>&lt;stdatomic.h&gt;</code>	
<code>#define _Atomic(T) /* see below */</code>	(5) <small>(since C++23)</small>

Each instantiation and full specialization of the `std::atomic` template defines an atomic type. If one thread writes to an atomic object while another thread reads from it, the behavior is well-defined (see [memory model](#) for details on data races).

In addition, accesses to atomic objects may establish inter-thread synchronization and order non-atomic memory accesses as specified by `std::memory_order`.

`std::atomic` is neither copyable nor movable.

Works on primitive types (or any type that is trivially copyable)

atomic variables ensure any read or write synchronizes

```
2     std::cout << "RESULT = " << shared_value << std::endl,
3
4     return 0;
5 }
```

# std::atomic (3/3)

- Here's the example, with minimal changes
  - No need for std::mutex on shared\_value
  - If two or more threads are writing (or one reading and one writing) to the object, the operations are well-defined.
    - (Note: I would probably rename variable with a suffix of \_atomic for readability)
    - Accidentally putting a lock, (or excessive locking) is bad for performance

```
1 // @file atomics.cpp
2 // g++ -std=c++17 atomics.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <atomic>    // New library
7
8 // Some shared value
9 static std::atomic<int> shared_value=0;
10
11 void increment_shared_value(){
12     shared_value++;
13 }
14
15 int main() {
16     std::vector<std::thread> threads;
17     // Create a collection of threads
18     for(int i=0; i < 1000; i++){
19         threads.push_back(std::thread(increment_shared_value));
20     }
21     // Join our threads using a ranged-based loop
22     // (Our intent is to iterate through all threads in container)
23     for(auto& th: threads){
24         threads.join();
25     }
26     // Retrieve our result
27     std::cout << "Result = " << shared_value << std::endl;
28
29     return 0;
30 }
```

# (Brief Aside) static in C++11 (and beyond)

- static local variable are guaranteed by the C++11 standard to only be initialized once (See also [std::call\\_once](#))

- Intuitively this makes sense, because the static variable in some way needs to be initialized exactly 1 time
- Takeaway: No lock needed

## Static local variables

Variables declared at block scope with the specifier `static` or `thread_local` (since C++11) have static or thread (since C++11) storage duration but are initialized the first time control passes through their declaration (unless their initialization is zero- or constant-initialization, which can be performed before the block is first entered). On all further calls, the declaration is skipped.

If the initialization throws an exception, the variable is not considered to be initialized, and initialization will be attempted again the next time control passes through the declaration.

If the initialization recursively enters the block in which the variable is being initialized, the behavior is undefined.

If multiple threads attempt to initialize the same static local variable concurrently, the initialization occurs exactly once (similar behavior can be obtained for arbitrary functions with `std::call_once`).

(since C++11)

Note: usual implementations of this feature use variants of the double-checked locking pattern, which reduces runtime overhead for already-initialized local statics to a single non-atomic boolean comparison.

The destructor for a block-scope static variable is called at program exit, but only if the initialization took place successfully.

Function-local static objects in all definitions of the same inline function (which may be implicitly inline) all refer to the same object defined in one translation unit, as long as the function has external linkage.

[https://en.cppreference.com/w/cpp/language/storage\\_duration#Static\\_local\\_variables](https://en.cppreference.com/w/cpp/language/storage_duration#Static_local_variables)

---

# Condition Variables

A way to signal an event between 2 or more threads

# Condition Variables (1/4)

- Perhaps a *stranger* thing when working with threads

- Idea:

- Given 2 threads (a worker and reporter), we want one thread to wait on the result of the other
  - Condition\_variable allows us to in a sense, take one thread off the queue, until it is notified that it should start working
    - (See example on right)

```
1 // @file cv.cpp
2 // g++ -std=c++17 cv.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <chrono>
7 #include <mutex>
8 #include <condition_variable> // New library
9
10 std::mutex gLock; // A global lock
11 std::condition_variable gConditionVariable;
12
13 int main() {
14     std::cout << "main() starts" << std::endl;
15
16     int result = 0; // The result we're trying to compute
17     bool notified = false;
18     // This thread reports the result (consumes the result)
19     // Could call this a 'consumer' or 'reporter' thread
20     std::thread reporter([&]{
21         std::unique_lock<std::mutex> lock(gLock);
22         // Wait here until notified
23         if(!notified){ // We need some variable to wait on
24             gConditionVariable.wait(lock);
25         }
26         std::cout << "\tReporter Result = " << result << std::endl;
27     });
28
29     // This thread does the work
30     std::thread worker([&]{
31         std::unique_lock<std::mutex> lock(gLock);
32         // Do our work
33         result = 42;
34         // Update our variable
35         notified = true;
36         // Artificial pause to show that reporter will indeed wait
37         std::this_thread::sleep_for(std::chrono::seconds(2));
38         // Output a message
39         std::cout << "\twork complete" << std::endl;
40         // Notify one thread about a change
41         // Could also notify_all
42         gConditionVariable.notify_one();
43     });
44
45
46     // Don't forget to join!
47     reporter.join();
48     worker.join();
49     std::cout << "main end " << std::endl;
50     return 0;
51 }
```

# Condition Variables (2/4)

- Perhaps a *stranger* thing when working with threads

- Idea:

- Given 2 threads (a worker and reporter), we want one thread to wait on the result of the other
  - Condition\_variable allows us to in a sense, take one thread off the queue, until it is notified that it should start working

Here's our reporter thread

reporter will wait on a 'variable' to be updated (notified)  
Then we'll be unblocked (gConditionVariable.wait) once we have been notified we can proceed

```
1 // @file cv.cpp
2 // g++ -std=c++17 cv.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <chrono>
7 #include <mutex>
8 #include <condition_variable> // New library
9
10 std::mutex gLock; // A global lock
11 std::condition_variable gConditionVariable;
12
13 int main() {
14     std::cout << "main() starts" << std::endl;
15
16     int result = 0; // The result we're trying to compute
17     bool notified = false;
18     // This thread reports the result (consumes the result)
19     // Could call this a 'consumer' or 'reporter' thread
20     std::thread reporter([&]{
21         std::unique_lock<std::mutex> lock(gLock);
22         // wait here until notified
23         if(!notified){ // We need some variable to wait on
24             gConditionVariable.wait(lock);
25         }
26         std::cout << "\tReporter Result = " << result << std::endl;
27     });
28
29     // This thread does the work
30     std::thread worker([&]{
31         std::unique_lock<std::mutex> lock(gLock);
32         // Do our work
33         result = 42;
34         // Update our variable
35         notified = true;
36         // artificial pause to show that reporter will indeed wait
37         std::this_thread::sleep_for(std::chrono::seconds(2));
38         // Output a message
39         std::cout << "\twork complete" << std::endl;
40         // Notify one thread about a change
41         // Could also notify_all
42         gConditionVariable.notify_one();
43     });
44
45
46     // Don't forget to join!
47     reporter.join();
48     worker.join();
49     std::cout << "main end " << std::endl;
50     return 0;
51 }
```

# Condition Variables (3/4)

- Perhaps a *stranger* thing when working with threads

- Idea:

- Given 2 threads (a worker and reporter), we want one thread to wait on the result of the other
  - Condition\_variable allows us to in a sense, take one thread off the queue, until it is notified that it should start working

Here's our worker thread

We update our result (i.e., do our work), and then change notified.

At the end, we notify\_one() to signal to our reporter that the reporter thread no longer needs to wait

```
1 // @file cv.cpp
2 // g++ -std=c++17 cv.cpp -o prog -lpthread
3 #include <iostream>
4 #include <thread>    // Include the thread library
5 #include <vector>
6 #include <chrono>
7 #include <mutex>
8 #include <condition_variable> // New library
9
10 std::mutex gLock; // A global lock
11 std::condition_variable gConditionVariable;
12
13 int main() {
14     std::cout << "main() starts" << std::endl;
15
16     int result = 0; // The result we're trying to compute
17     bool notified = false;
18     // This thread reports the result (consumes the result)
19     // Could call this a 'consumer' or 'reporter' thread
20     std::thread reporter([&]{
21         std::unique_lock<std::mutex> lock(gLock);
22         // wait here until notified
23         if(!notified){ // We need some variable to wait on
24             gConditionVariable.wait(lock);
25         }
26         std::cout << "\tReporter Result = " << result << std::endl;
27     });
28
29     // This thread does the work
30     std::thread worker([&]{
31         std::unique_lock<std::mutex> lock(gLock);
32         // Do our work
33         result = 42;
34         // Update our variable
35         notified = true;
36         // artificial pause to show that reporter will indeed wait
37         std::this_thread::sleep_for(std::chrono::seconds(2));
38         // Output a message
39         std::cout << "\twork complete" << std::endl;
40         // Notify one thread about a change
41         // Could also notify_all
42         gConditionVariable.notify_one();
43     });
44
45
46     // Don't forget to join!
47     reporter.join();
48     worker.join();
49     std::cout << "main end " << std::endl;
50     return 0;
51 }
```

```
mike:concurrency$ g++-10 -std=c++17 cv.cpp -o prog -lpthread  
mike:concurrency$ ./prog  
main() starts  
    work complete  
    Reporter Result = 42
```

- **main end**  
with threads

- Idea:
  - Given 2 threads (a worker and reporter), we want one thread to wait on the result of the other
  - Condition\_variable allows us to in a sense, take one thread off the queue, until it is notified that it should start working
    - (See example on right)

```
1 // @file cv.cpp  
2 //  
3 // This file shows how to use std::condition_variable to compute  
4 // results.  
5 //  
6 // This program has two threads:  
7 // 1. A reporter thread which waits for a result to be computed.  
8 // 2. A worker thread which does the work.  
9 //  
10 // The reporter thread calls std::condition_variable::wait() to  
11 // wait for the worker thread to signal that the result is ready.  
12 // The worker thread calls std::condition_variable::notify_one()  
13 // to wake up the reporter thread.  
14 //  
15 // The reporter thread then outputs the result.  
16 //  
17 bool notified = false;  
18 // This thread reports the result (consumes the result)  
19 // Could call this a 'consumer' or 'reporter' thread  
20 std::thread reporter([&]{  
21     std::unique_lock<std::mutex> lock(gLock);  
22     // wait here until notified  
23     if(!notified){ // We need some variable to wait on  
24         gConditionVariable.wait(lock);  
25     }  
26     std::cout << "\tReporter Result = " << result << std::endl;  
27 });  
28  
29 // This thread does the work  
30 std::thread worker([&]{  
31     std::unique_lock<std::mutex> lock(gLock);  
32     // Do our work  
33     result = 42;  
34     // Update our variable  
35     notified = true;  
36     // artificial pause to show that reporter will indeed wait  
37     std::this_thread::sleep_for(std::chrono::seconds(2));  
38     // Output a message  
39     std::cout << "\twork complete" << std::endl;  
40     // Notify one thread about a change  
41     // Could also notify_all  
42     gConditionVariable.notify_one();  
43 });  
44  
45  
46 // Don't forget to join!  
47 reporter.join();  
48 worker.join();  
49 std::cout << "main end " << std::endl;  
50 return 0;  
51 }
```

# std::unique\_lock - Other mechanisms

- std::unique\_lock (C++11) -
  - A bit more powerful than lock\_guard and scoped\_guard in that we can control locking and unlocking
  - Used in condition\_variable
  - Also follows RAII so we can use it safely.

Member functions	
(constructor)	constructs a unique_lock, optionally locking (public member function)
(destructor)	unlocks (i.e., releases ownership) (public member function)
<code>operator=</code>	unlocks (i.e., releases ownership) (public member function)
Locking	
<code>lock</code>	locks (i.e., takes ownership of) the mutex (public member function)
<code>try_lock</code>	tries to lock (i.e., takes ownership) (public member function)
<code>try_lock_for</code>	attempts to lock (i.e., takes ownership) if the mutex has been unavailable for the specified duration (public member function)
<code>try_lock_until</code>	tries to lock (i.e., takes ownership) until the specified time (public member function)
<code>unlock</code>	unlocks (i.e., releases ownership) (public member function)
Modifiers	
<code>swap</code>	swaps state with another std::unique_lock (public member function)
<code>release</code>	disassociates the associated mutex (public member function)
Observers	
<code>mutex</code>	returns a pointer to the associated mutex (public member function)
<code>owns_lock</code>	tests whether the lock owns (i.e., is locked) (public member function)
<code>operator bool</code>	tests whether the lock owns (i.e., is locked) (public member function)

# There exist several other primitives you can find here

- [https://en.cppreference.com/w/cpp/thread#Mutual\\_exclusion](https://en.cppreference.com/w/cpp/thread#Mutual_exclusion)

Mutual exclusion	
Mutual exclusion algorithms prevent multiple threads from simultaneously accessing shared resources. This prevents data races and provides support for synchronization between threads.	
Defined in header < <code>mutex</code> >	
<a href="#"><code>mutex</code> (C++11)</a>	provides basic mutual exclusion facility (class)
<a href="#"><code>timed_mutex</code> (C++11)</a>	provides mutual exclusion facility which implements locking with a timeout (class)
<a href="#"><code>recursive_mutex</code> (C++11)</a>	provides mutual exclusion facility which can be locked recursively by the same thread (class)
<a href="#"><code>recursive_timed_mutex</code> (C++11)</a>	provides mutual exclusion facility which can be locked recursively by the same thread and implements locking with a timeout (class)
Defined in header < <code>shared_mutex</code> >	
<a href="#"><code>shared_mutex</code> (C++17)</a>	provides shared mutual exclusion facility (class)
<a href="#"><code>shared_timed_mutex</code> (C++14)</a>	provides shared mutual exclusion facility and implements locking with a timeout (class)
Generic mutex management	
Defined in header < <code>mutex</code> >	
<a href="#"><code>lock_guard</code> (C++11)</a>	implements a strictly scope-based mutex ownership wrapper (class template)
<a href="#"><code>scoped_lock</code> (C++17)</a>	deadlock-avoiding RAII wrapper for multiple mutexes (class template)
<a href="#"><code>unique_lock</code> (C++11)</a>	implements movable mutex ownership wrapper (class template)
<a href="#"><code>shared_lock</code> (C++14)</a>	implements movable shared mutex ownership wrapper (class template)
<a href="#"><code>defer_lock_t</code> (C++11)</a> <a href="#"><code>try_to_lock_t</code> (C++11)</a> <a href="#"><code>adopt_lock_t</code> (C++11)</a>	tag type used to specify locking strategy (class)
<a href="#"><code>defer_lock</code> (C++11)</a> <a href="#"><code>try_to_lock</code> (C++11)</a> <a href="#"><code>adopt_lock</code> (C++11)</a>	tag constants used to specify locking strategy (constant)

---

# Quick Recap

(Of the C++ we have learned)

# The C++ What we have so far

---

- std::thread used with ‘join’
- std::jthread (C++20 feature)
- std::mutex
  - lock and unlock (to create a critical section that is atomically executed)
- std::lock\_guard (and also a note on scoped\_lock and unique\_lock)
  - Which helps us from forgetting to ‘unlock’ a lock (Uses RAII to release lock)
- condition\_variable
  - A way to ‘signal’ between multiple threads when work is done
- std::atomic<T>
  - Useful for primitive types we want to synchronize
- And 2 caveats to worry about with concurrency
  - data races (non-deterministic behavior with reads and write order)
  - deadlock (lack of progress)

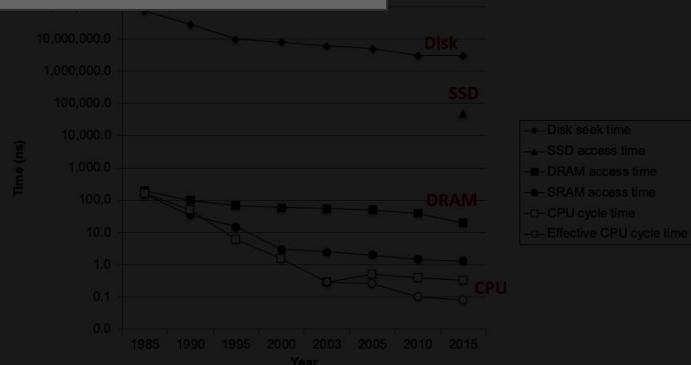
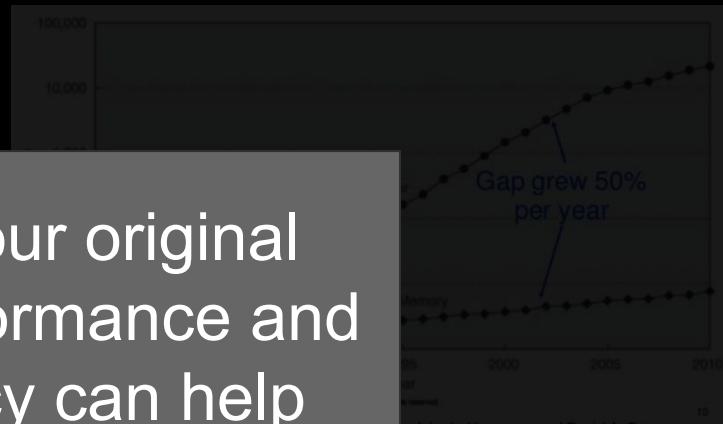
# Our original motivation was about performance (1/2)

- An interesting reality is that many applications we write are I/O bound

- That means many operations
- The figure shows that processor access time is increasing in magnitude

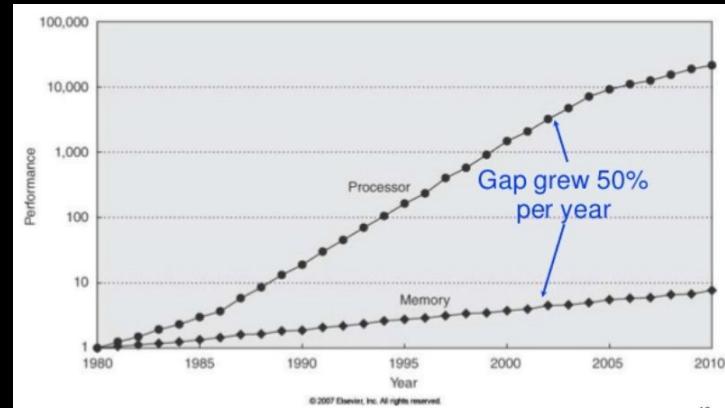
Let's try to answer our original motivation about performance and see how concurrency can help

- Thus, we often cannot fetch data at the rate that we process it.

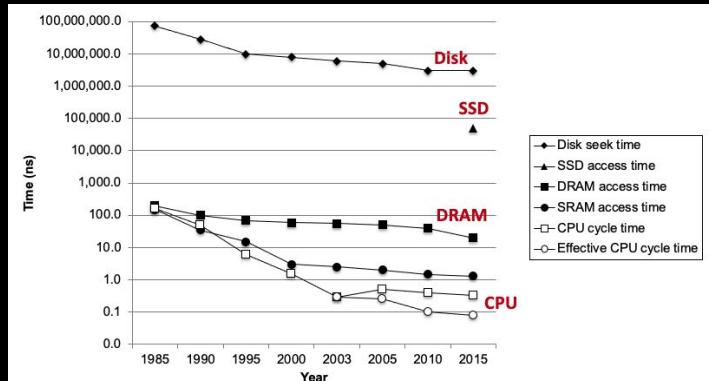


# Our original motivation was about performance (2/2)

- An interesting reality is that many applications we write are I/O bound
  - That means that we are waiting on memory operations
  - The figure to the right shows that while our processors have gotten faster over time, accessing memory remains orders of magnitude slower
    - Thus, we often cannot fetch data at the rate that we process it.



Source: Computer Architecture, A Quantitative Approach by John L. Hennessy and David A. Patterson



---

# Asynchronous Programming

Another form concurrency where execution can happen independently of the main program flow

Asynchronous means that events happen ‘without synchronicity’ or ‘without order’.

# std::async example

- `#include <future>`
- `std::async`
  - Promise and Future
    - Promise - Will hold the result
    - Future - Where the future result will be stored
  - We are blocked at `a.get()` until the value is returned.

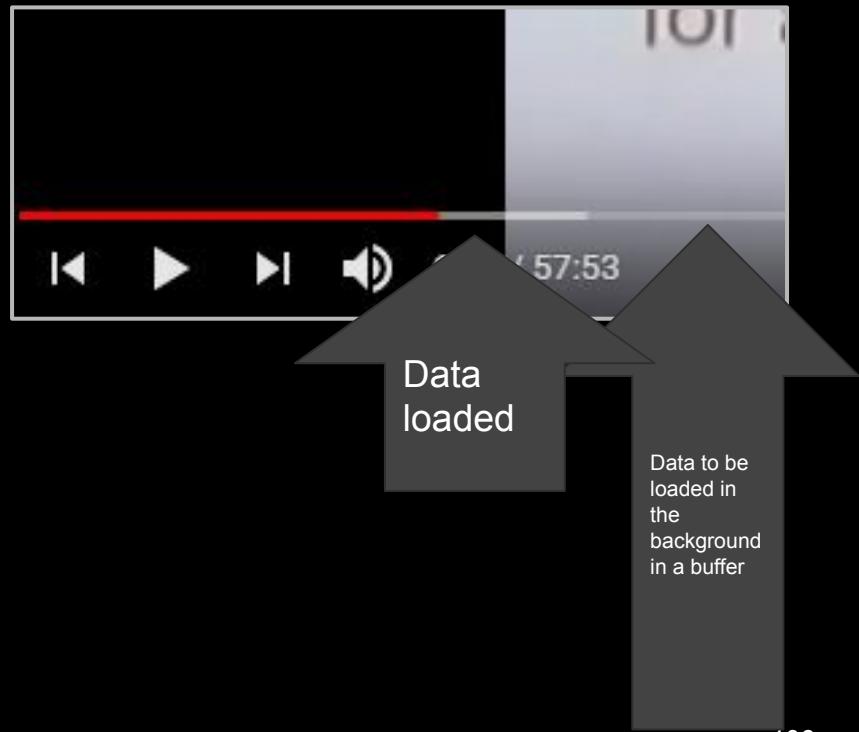
```
1 // @file thread10.cpp
2 // g++ -std=c++17 thread10.cpp -o prog -lpthread
3 #include <iostream>
4 #include <future>
5
6 int square(int x){
7     return x*x;
8 }
9
10 int main() {
11
12     // asyncFunction is a 'future'
13     // type = std::future<int>
14     auto asyncFunction = std::async(&square,12);
15     // .... some time passes
16     int result = asyncFunction.get(); // We are blocked here if
17                                         // our value has not been
18                                         // computed. Otherwise, the
19                                         // value from get() (which
20                                         // is wrapped in a promise
21                                         // is returned.
22     std::cout << "The async thread has returned! " << result
23                                         << std::endl;
24
25     return 0;
26 }
27
```

1,1

```
mike:concurrency$ g++-10 -std=c++17 thread10.cpp -o prog -lpthread
mike:concurrency$ ./prog
The async thread has returned! 144
```

# A Concrete Example for std::async

- Blocking Input/Output (I/O)
  - I/O is any task where we are reading or writing data.
    - e.g. network connection (e.g. downloading data), disk load (e.g. opening a file)
  - We can use a ‘background thread’ (i.e., std::async) execute to start loading that data.
    - The application can then proceed unblocked until it needs that data.
    - If we do not have the data ready when we need, we are thus ‘blocked’ -- hence the term Blocking I/O



# Async I/O Simulation (1/5)

- “mocked” version of using an **async** thread to load data
  - We spawn a ‘background thread’ asynchronously using `std::async`
  - Then in our ‘main loop’ we continuously query to see if our function has returned
  - (I have added a few artificial sleeps to make it more interesting)

```
1 // @file async.cpp
2 // g++ -std=c++17 async.cpp -o prog -lpthread
3 #include <iostream>
4 #include <future>
5 #include <thread>
6 #include <chrono>
7
8 // Toy example - This will be our 'background thread'
9 // that executes asynchronously
10 bool bufferedFileLoading(){
11     size_t bytesLoaded= 0;
12     while(bytesLoaded < 20000){
13         std::cout << "Loading File..." << std::endl;
14         std::this_thread::sleep_for(std::chrono::milliseconds(250));
15         bytesLoaded+= 1000;
16     }
17     return true;
18 }
19
20 int main() {
21     // Launch thread asynchronously, and this will execute in background
22     std::future<bool> backgroundThread = std::async(std::launch::async,
23                                                     bufferedFileLoading);
24
25     std::future_status status;
26     // Meanwhile, we have our main thread of execution
27     while(true){
28         std::cout << "Main thread running" << std::endl;
29         // artificial pause
30         std::this_thread::sleep_for(std::chrono::milliseconds(50));
31         // Check if our
32         status = backgroundThread.wait_for(std::chrono::milliseconds(1));
33         // If our data is ready--do something
34         // we'll just terminate for now
35         if(status == std::future_status::ready){
36             std::cout << "data ready..." << std::endl;
37             break;
38         }
39     }
40
41     return 0;
42 }
```

# As

Here we'll create a background thread that will execute with std::async

I've been explicit in setting up the parameters and types.

Also, there is a 'status' that we'll keep track of so we know when a value has been returned

```
1 // @file async.cpp
2 // g++ -std=c++17 async.cpp -o prog -lpthread
3 #include <iostream>
4 #include <future>
5 #include <thread>
6 #include <chrono>
7
8 // Toy example - This will be our 'background thread'
9 // that executes asynchronously
10 bool bufferedFileLoading(){
11     size_t bytesLoaded= 0;
12     while(bytesLoaded < 20000){
13         std::cout << "Loading File..." << std::endl;
14         std::this_thread::sleep_for(std::chrono::milliseconds(250));
15         bytesLoaded+= 1000;
16     }
17     return true;
18 }
19
20 int main() {
21     // Launch thread asynchronously, and this will execute in background
22     std::future<bool> backgroundThread = std::async(std::launch::async,
23                                                     bufferedFileLoading);
24
25     // Store status of our future
26     std::future_status status;
27
28     // Meanwhile, we have our main thread of execution
29     while(true){
30         std::cout << "Main thread running" << std::endl;
31         // artificial pause
32         std::this_thread::sleep_for(std::chrono::milliseconds(50));
33         // Check if our
34         status = backgroundThread.wait_for(std::chrono::milliseconds(1));
35         // If our data is ready--do something
36         // we'll just terminate for now
37         if(status == std::future_status::ready){
38             std::cout << "data ready..." << std::endl;
39             break;
40         }
41     }
42 }
```

# Here we are reading in 'bytes' from a file.

Perhaps we are 'streaming'  
in some # of bytes from a  
data source

- asynchronously using std::async
- Then in our 'main loop' we continuously query to see if our function has returned
- (I have added a few artificial sleeps to make it more interesting)

```
1 // @file async.cpp
2 // g++ -std=c++17 async.cpp -o prog -lpthread
3 #include <iostream>
4 #include <future>
5 #include <thread>
6 #include <chrono>
7
8 // Toy example - This will be our 'background thread'
9 // that executes asynchronously
10 bool bufferedFileLoading(){
11     size_t bytesLoaded= 0;
12     while(bytesLoaded < 20000){
13         std::cout << "Loading File..." << std::endl;
14         std::this_thread::sleep_for(std::chrono::milliseconds(250));
15         bytesLoaded+= 1000;
16     }
17     return true;
18 }
19
20 int main() {
21     // Launch thread asynchronously, and this will execute in background
22     std::future<bool> backgroundThread = std::async(std::launch::async,
23                                                     bufferedFileLoading);
24
25     std::future_status status;
26     // Meanwhile, we have our main thread of execution
27     while(true){
28         std::cout << "Main thread running" << std::endl;
29         // artificial pause
30         std::this_thread::sleep_for(std::chrono::milliseconds(50));
31         // Check if our
32         status = backgroundThread.wait_for(std::chrono::milliseconds(1));
33         // If our data is ready--do something
34         // we'll just terminate for now
35         if(status == std::future_status::ready){
36             std::cout << "data ready..." << std::endl;
37             break;
38         }
39     }
40
41     return 0;
42 }
```

# Async I/O Simulation (4/5)

- “mocked” version of using an **async** thread to load data
  - We spawn a ‘background thread’ asynchronously using `std::async`

In our main loop we will check every 1 millisecond the status of our future value (which is wrapped in a promise object)

```
1 // @file async.cpp
2 // g++ -std=c++17 async.cpp -o prog -lpthread
3 #include <iostream>
4 #include <future>
5 #include <thread>
6 #include <chrono>
7
8 // Toy example - This will be our 'background thread'
9 // that executes asynchronously
10 bool bufferedFileLoading(){
11     size_t bytesLoaded= 0;
12     while(bytesLoaded < 20000){
13         std::cout << "Loading File..." << std::endl;
14         std::this_thread::sleep_for(std::chrono::milliseconds(250));
15         bytesLoaded+= 1000;
16     }
17     return true;
18 }
19
20 int main() {
21     // Launch thread asynchronously, and this will execute in background
22     std::future<bool> backgroundThread = std::async(std::launch::async,
23                                                     bufferedFileLoading);
24
25     // Store status of our future
26     std::future_status status;
27     // Meanwhile, we have our main thread of execution
28     while(true){
29         std::cout << "Main thread running" << std::endl;
30         // artificial pause
31         std::this_thread::sleep_for(std::chrono::milliseconds(50));
32         // Check if our
33         status = backgroundThread.wait_for(std::chrono::milliseconds(1));
34         // If our data is ready--do something
35         // we'll just terminate for now
36         if(status == std::future_status::ready){
37             std::cout << "data ready..." << std::endl;
38             break;
39     }
40
41     return 0;
42 }
```

# Async I/O Simulation (5/5)

```
mike:concurrency$
```

1

```
1 // @file async.cpp
2 // g++ -std=c++
3 #include <iostream>
4 #include <future>
5 #include <thread>
6 #include <chrono>
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
example
bufferedFileLoading()
{
    std::atomic<bool> bytesLoaded{false};
    std::atomic<bool> ready{false};

    std::thread backgroundThread{[&] {
        std::cout << "Background thread started" << std::endl;
        std::this_thread::sleep_for(std::chrono::milliseconds(50));
        bytesLoaded.store(true);
    }};
    std::future<bool> backgroundFuture = backgroundThread.get_future();

    std::thread mainThread{[&] {
        std::cout << "Main thread started" << std::endl;
        std::this_thread::sleep_for(std::chrono::milliseconds(10));
        if(backgroundFuture.wait_for(std::chrono::milliseconds(1)) == std::future_status::ready) {
            std::cout << "Background thread finished" << std::endl;
            bytesLoaded.store(true);
        }
    }};
    std::future<bool> mainFuture = mainThread.get_future();

    std::cout << "Waiting for both threads to finish" << std::endl;
    std::cout << "Main thread finished" << std::endl;
    std::cout << "Background thread finished" << std::endl;
}
```

(Again--code available from github)

---

# Other Odds and Ends

Coming to the finale--a rapid fire section of topics not covered

# Thread Sanitizer

---

- Since we have talked quite a bit about threads in the first portion--thread sanitizer is a useful tool for detecting concurrency violations (i.e., data races).
  - See an example with clang: <https://clang.llvm.org/docs/ThreadSanitizer.html>

# Performance

---

- I haven't touched on this--but launching 10000 threads does not make your application **10000 times faster**.
  - There is some overhead with threads (creation of the threads, context switching, managing the threads, etc.)
  - Likely you'll want to have some fixed number of threads (perhaps as a function of how many CPU cores you have).
    - Folks will create something called a 'thread pool' to structure their concurrent programs.

# C++20 Coroutines [[reference](#)]

---

- Coroutines are another way to ‘collaboratively’ write code.
  - They’re new enough in C++ I personally need to learn/think more about them.
  - Coroutines to my understanding allowed you to yield and save state (which could be done asynchronously) and give you more control of your program (i.e., lazy-computation).
  - Nice explanation of coroutines:  
<https://stackoverflow.com/questions/553704/what-is-a-coroutine>  
<https://en.cppreference.com/w/cpp/language/coroutines>

# C++ Memory Model [[reference](#)]

---

- The memory model in C++ is what has enabled std::thread and other concurrent programming techniques
  - This defines the order of ‘reads’ and ‘writes’ to memory
  - Having this order defined, helps stabilize the language, and ensure that whoever implements <thread> is able to do so safely.
- For a beginner, it’s worth running the examples provided here a few times, then perhaps viewing the memory model as you move further.
  - (For more advanced users, please take a look)

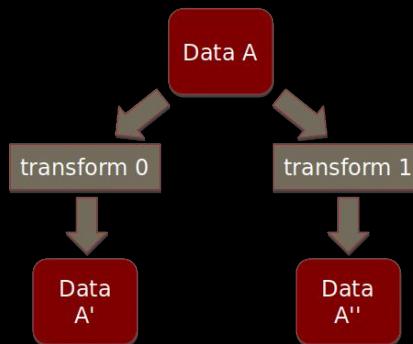
---

# Common Concurrency Patterns

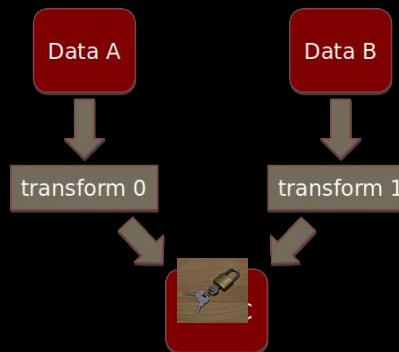
You will know enough to implement these (They take a little thought!)

# Reader/Writers Problem

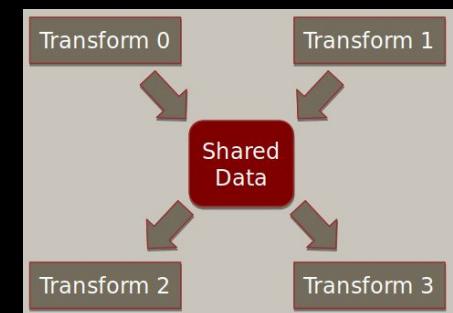
- Consider this pattern where you have to either read to a location, write to a location, or read and write to a location
  - We have `std::shared_mutex` which allows multiple threads to read a value
  - Utilizing `std::shared_mutex` with `std::unique_lock` can help solve the readers/writers problem



Here multiple reader threads are 'reading' Data A and performing some transformation.



Here multiple writer threads are generating and combining their result into Data 'C'.



Here multiple writer threads are creating shared data, and multiple readers than take that data as input.

# Other popular concurrency patterns

---

- Producer Consumer pattern
  - Pattern in which ‘producer’ threads generate data
    - (Usually adding to some queue)
  - Consumer threads then consume any data from the queue, and process it.
  - [https://en.wikipedia.org/wiki/Producer%20consumer\\_problem](https://en.wikipedia.org/wiki/Producer%20consumer_problem)
- Barrier
  - Allow a fixed number of threads to enter a critical section
  - (C++20) [std::counting\\_semaphore](#)
    - Allow multiple threads to enter
    - Related ideas
      - (C++20) [std::latch](#)
      - (C++20) [std::barrier](#)

---

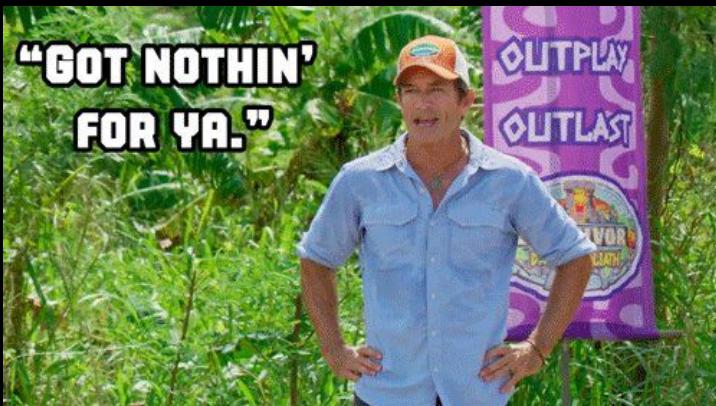
# Brief Notes on Parallelism

Parallelism being the simultaneous execution of threads (often with less thought to think about synchronization (i.e., locks) for correctness\*)

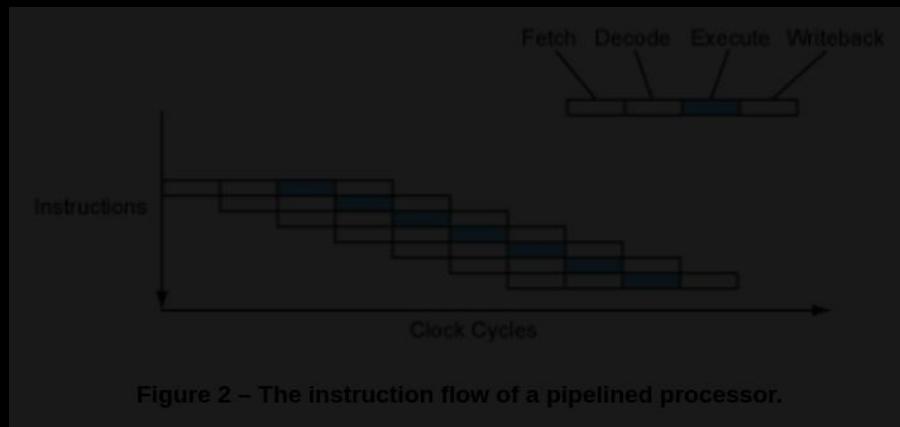
\*Still need to think about synchronization points in a sense for performance

# Some parallelism for free (implicit parallelism)

- No programming required--hardware does the work



Reference to American show *Survivor* in which contestants receive no prize for losing a challenge



# Some parallelism for free (implicit parallelism)

- No programming required--hardware does the work
  - CPU Pipelining is an example of *parallelism* we typically get for free
    - (i.e. implicit parallelism)
  - Potential compiler optimizations to automatically vectorize code.

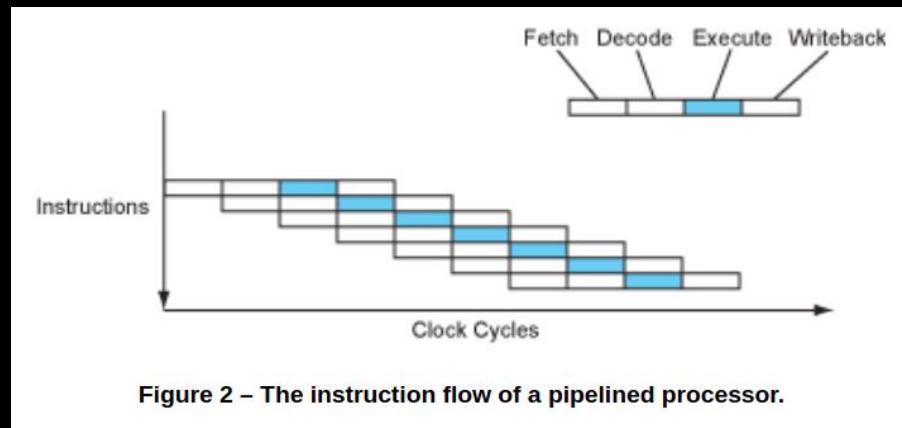
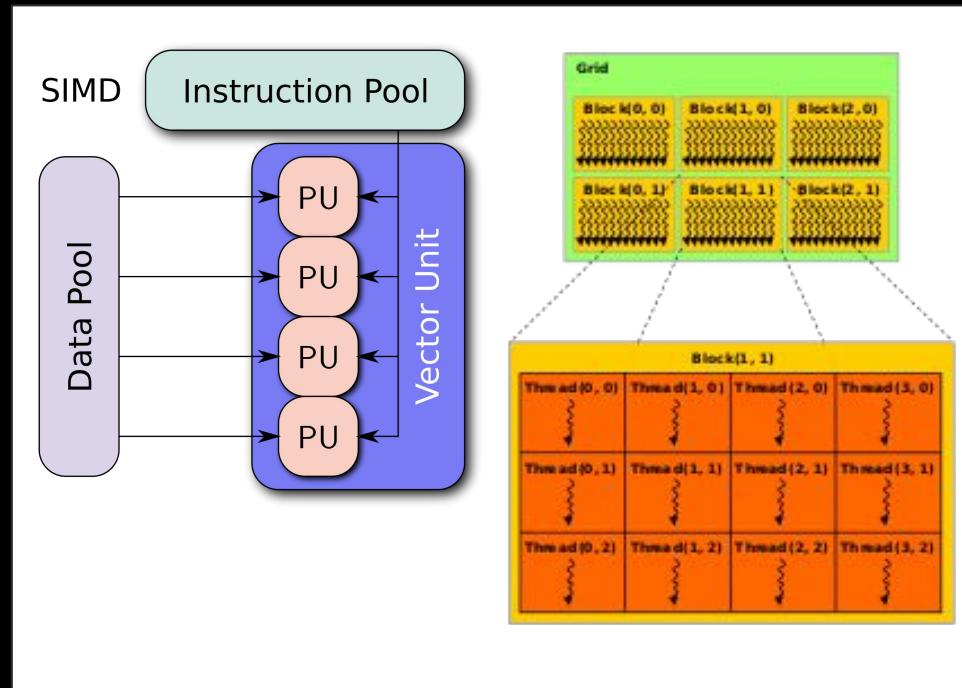


Figure 2 – The instruction flow of a pipelined processor.

<https://www.lighterra.com/papers/modernmicroprocessors/>

# Some parallelism not for free (explicit parallelism)

- Requires programming CPU or GPU explicitly
  - Using the GPU
    - Whether CUDA or OpenCL for general purpose GPU programming
    - Or perhaps a shader language like GLSL or HLSL
    - Or any other API
  - SIMD Instructions
    - Our SSE or AVX instructions



# Examples of Parallelism

- Pretty much everything in graphics
  - <https://software.intel.com/en-us/articles/performance-methods-and-practices-of-directx-11-multithreaded-rendering>
- Pretty much everywhere in machine learning...
- Pretty much everywhere you need to increase performance! :)



# C++ std::par and std::par\_unseq

---

- [CppCon 2016: Bryce Adelstein Lelbach “The C++17 Parallel Algorithms Library and Beyond”](#)
  - Implementation notes
- [CppCon 2017: Dietmar Kühl “C++17 Parallel Algorithms”](#)
  - Usage examples

---

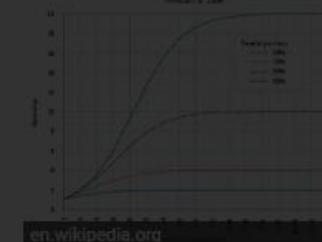
# Conclusion and Further Resources

And some analogies (for teaching) and to wrap up what we've learned

# Why Parallel?

- Performance (execution speed)
- But how much performance?

**Amdahl's law** is a formula used to find the maximum improvement possible by improving a particular part of a system. In parallel computing, **Amdahl's law** is mainly used to predict the theoretical maximum speedup for program processing using multiple processors. ... This term is also known as **Amdahl's argument**.



en.wikipedia.org

What is Amdahl's Law? - Definition from Techopedia

<https://www.techopedia.com/definition/17035/amdaahls-law>

I like to include Amdahl's law somewhere so we can even think if we have the opportunity to make part of our program concurrent or parallel

$$(1 - P) + s$$

s = speedup of task that benefits from improved resources

p = portion of execution time benefiting from improved speedup

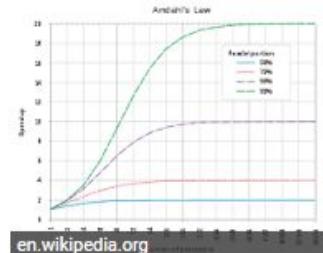
[https://en.wikipedia.org/wiki/Amdahl%27s\\_law](https://en.wikipedia.org/wiki/Amdahl%27s_law)

Applied example: <http://web.cs.iastate.edu/~prabhu/Tutorial/CACHE/CompPerf.pdf>

# Why Parallel?

- Performance (execution speed)
- But how much performance?

**Amdahl's law** is a formula used to find the maximum improvement possible by improving a particular part of a system. In parallel computing, **Amdahl's law** is mainly used to predict the theoretical maximum speedup for program processing using multiple processors. ... This term is also known as **Amdahl's argument**.



What is Amdahl's Law? - Definition from Techopedia  
<https://www.techopedia.com/definition/17035/amdalhs-law>

$$S_{\text{latency}}(s) = \frac{1}{(1-p) + \frac{p}{s}}$$

s = speedup of task that benefits from improved resources

p = portion of execution time benefiting from improved speedup

[https://en.wikipedia.org/wiki/Amdahl%27s\\_law](https://en.wikipedia.org/wiki/Amdahl%27s_law)

Applied example: <http://web.cs.iastate.edu/~prabhu/Tutorial/CACHE/CompPerf.pdf>

# Concurrency Caveat #1 | Data Race

---

- When two (or more) processes contending for one shared resource.
  - One parking spot 2 cars want to acquire
    - (Do both acquire, do neither, does only 1? Depends on the driver I suppose!)



# Data race is not always as obvious...(1/4)

---

- Imagine you check your fridge and find there is no milk
  - So you run to the store



Googling the “Too much milk” problem has an interesting lineage -- <https://blog.regehr.org/archives/1145>  
(Not certain the origin, but it’s used for decades)

# Data race is not always as obvious...(2/4)

---

- Imagine you check your fridge and find there is no milk
  - So you run to the store
- Then moments later your roommate checks the fridge and finds it is empty
  - So they run to the store



Googling the “Too much milk” problem has an interesting lineage -- <https://blog.regehr.org/archives/1145>  
(Not certain the origin, but it’s used for decades)

# Data race is not always as obvious...(3/4)

- Imagine you check your fridge and find there is no milk
  - So you run to the store
- Then moments later your roommate checks the fridge and finds it is empty
  - So they run to the store
- Roommate # 3 comes and notices the same
  - ....



Googling the “Too much milk” problem has an interesting lineage -- <https://blog.regehr.org/archives/1145>  
(Not certain the origin, but it's used for decades)

# Data race is not always as obvious...(4/4)

- You get the idea when you then find out you have 3 times as much milk as your house needs when everyone returns.
  - Then you can talk to your students about using locking mechanisms for safety
  - After they learn the foundations, then provide further guidance for lock\_guard, etc.

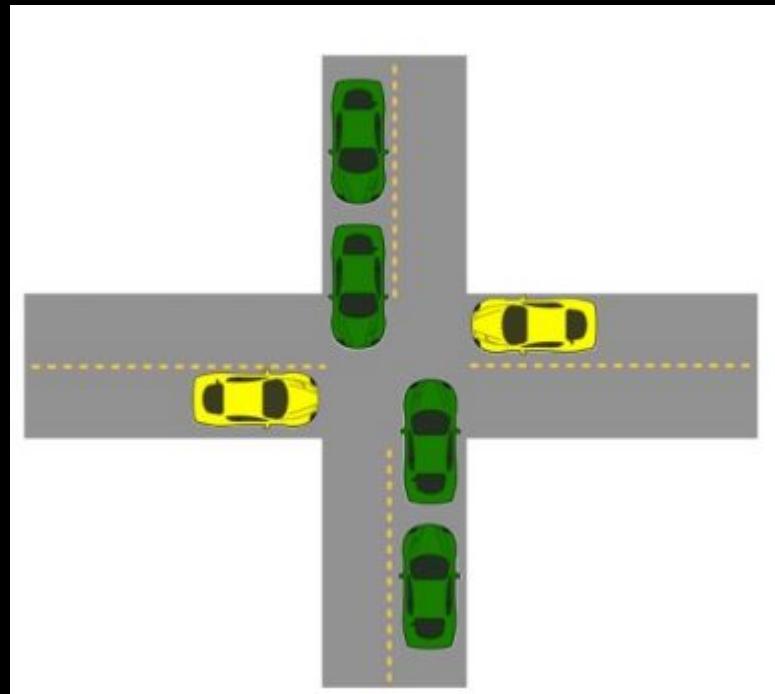


Googling the “Too much milk” problem has an interesting lineage -- <https://blog.regehr.org/archives/1145>  
(Not certain the origin, but it's used for decades)

# Concurrency Caveat #2 | Starvation

---

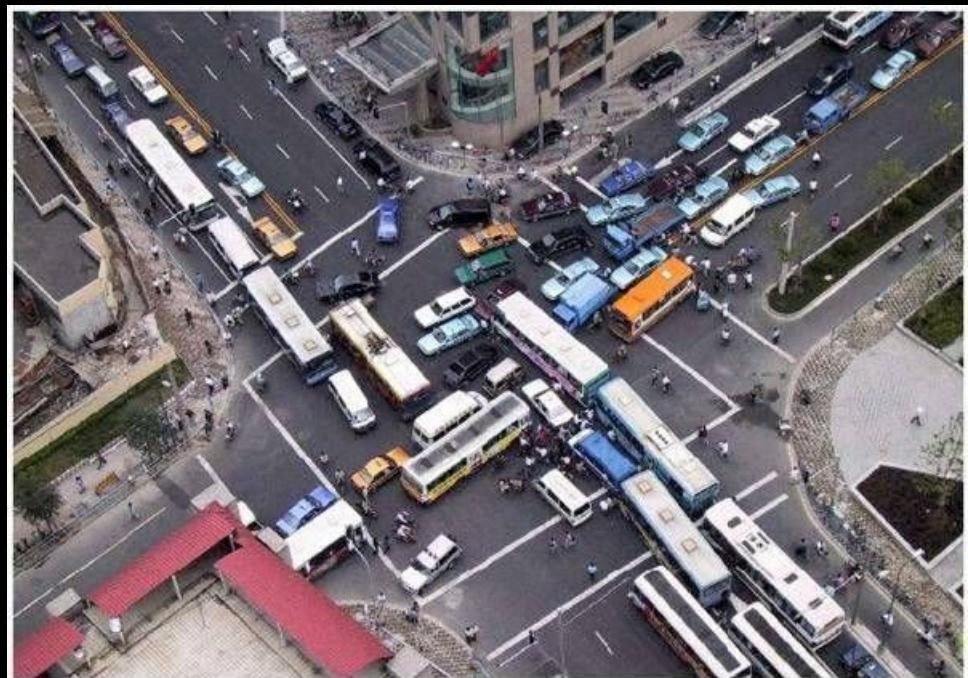
- Imagine a constant stream of green cars
  - Progress is always made by the green cars
- The yellow cars can never make progress to get across the street.
  - They are resource starved of a shared resource (again, they cannot cross the intersection)



# Concurrency Caveat #3 | Deadlock

---

- Grid lock in a traffic jam
  - Each car prevents others from going through a shared resource (the intersection).
- (One car needs mutually exclusive access to a portion of the intersection in order to move forward)
- (deadlock is an extreme form of starvation)



# Concurrent Programming takes some extra care

---

1. Races: Outcome depends on the arbitrary scheduling decisions elsewhere in the system
    - e.g. Who gets the last seat on the airplane.
  2. Starvation/Fairness: External events and/or scheduling decisions can prevent sub-task progress
    - e.g. Someone jumping in front of you in line
  3. Deadlock: Improper resource allocation prevents forward progress
    - e.g. traffic gridlock
- But regardless, concurrent programming is important and necessary!

# Concurrent Programming takes some extra care

---

1. Races: Outcome depends on the arbitrary scheduling decisions elsewhere in the system
    - { The examples shown today are good ‘toy’ examples to learn from.
  2. S{
    - Si Consider what layers of abstraction you can add to make them
    - { safer if your domain allows (i.e., avoiding managing lots of little
  3. D{
    - global std::mutex. )
- But regardless, concurrent programming is important and necessary!

# Further Resources

---

- Code from this talk: <https://github.com/MikeShah/cppcon2021>
- [Back to Basics: Concurrency - Arthur O'Dwyer - CppCon 2020](#)
  - Great resource diving further into concurrency primitives
- <https://www.modernescpp.com/>
  - Rainer Grimm (Several nice articles on website, and also several cppcon topics on beginner and advanced concurrency topics)
- [CppCon 2017: Fedor Pikus “C++ atomics, from basic to advanced. What do they really do?”](#)
  - Full look at atomics
- Little Book of Semaphores <https://greenteapress.com/wp/semaphores/>
  - Useful puzzles to reason about common concurrency patterns (reader/writer, producer/consumer, etc.)
- The ‘comet book’ <https://pages.cs.wisc.edu/~remzi/OSTEP/>
  - (Free/donation based) Operating Systems book Insight into how foundational constructs work

# Closing--Why do we need concurrency and parallelism?

---

- Again--the main motivator is performance
  - Doing multiple things at once *things* at once is faster than doing one thing at a time.
  - Sometimes however, it makes our programs cleaner (`std::async`) and otherwise is a logical way to solve a problem.
- Sometimes concurrency is necessary to model how things truly work in the real world
  - e.g. An orchestra operates concurrently
  - e.g. A traffic light that controls an intersection
  - e.g. pragmatically it may just help us structure our problem and break work into smaller chunks



# Back to Basics: Concurrency

Mike Shah, Ph.D.

[@MichaelShah](https://@MichaelShah) | [mshah.io](https://mshah.io) | [www.youtube.com/c/MikeShah](https://www.youtube.com/c/MikeShah)

Thank you Cppcon attendees, reviewers, chairs!

---

Thank you!