

# Mutex Conclusion

## Mutex Conclusion

- We have covered
  - Basic mutex class (`std::mutex`)
  - Mutex wrappers (`std::lock_guard` and `std::unique_lock`)
  - Mutexes and time-outs (`std::timed_mutex`, `std::unique_lock`)
  - Shared mutexes (`std::shared_mutex` and `std::shared_lock`)
  - Static and thread-local data
  - Deadlock
  - Deadlock avoidance (`std::scoped_lock`, `std::lock()`)
  - Livelock and livelock avoidance

## Locking Guidelines

- Locking impacts on other threads
  - They will have to wait longer for a resource they need
  - This affects performance
- Always hold a lock for the shortest possible time
- Avoid locking lengthy operations if possible
  - e.g. input/output

## Recommendations for Reading Shared Data

- Reading
  - Lock
  - Make a copy of the shared data
  - Unlock and process the copy

## Locking Guidelines for Data Structures

- Do not lock any more elements than necessary
  - e.g. Accessing a single element in a linked list
  - Do not lock the entire list
  - This will block other threads from accessing unrelated elements
- Do not make locking too fine-grained
  - Do not lock individual elements when inserting and deleting
  - Another thread may need to access a neighbouring element
  - Data race
  - Need to lock neighbouring elements in a single operation

## Pros and Cons of Mutexes

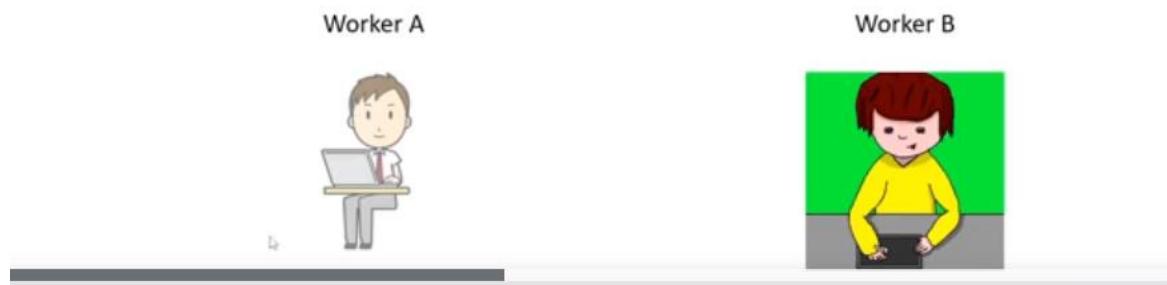
- Fairly straightforward way to protect shared data
  - Prevent data races and race conditions
- Locking and unlocking are slow operations
- Low-level
  - The programmer must remember to use a mutex
  - The programmer must use the right mutex
  - The programmer must understand how different threads can modify the data
- Most real-world programs use higher-level structures
  - Mutex wrapper classes
  - Classes from the next section!

# Thread Coordination

## Coordination Between Workers



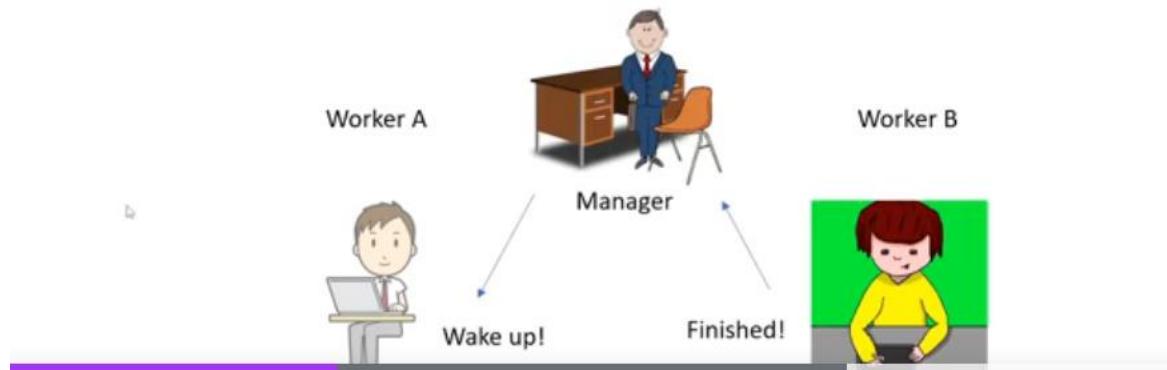
- Problem:
  - Worker A is writing a document
  - Worker A needs an image that Worker B is producing
  - Worker A cannot continue until Worker B has finished



## Coordination Between Workers



- One solution is to introduce a manager
  - The manager coordinates the two workers



## Coordination Between Workers

- Solution:
- Worker B tells the manager when they have finished
- The manager tells Worker A to resume work
  - B is working
  - A is waiting
  - B finishes their work
  - B tells manager
  - Manager tells A to resume
  - A resumes work

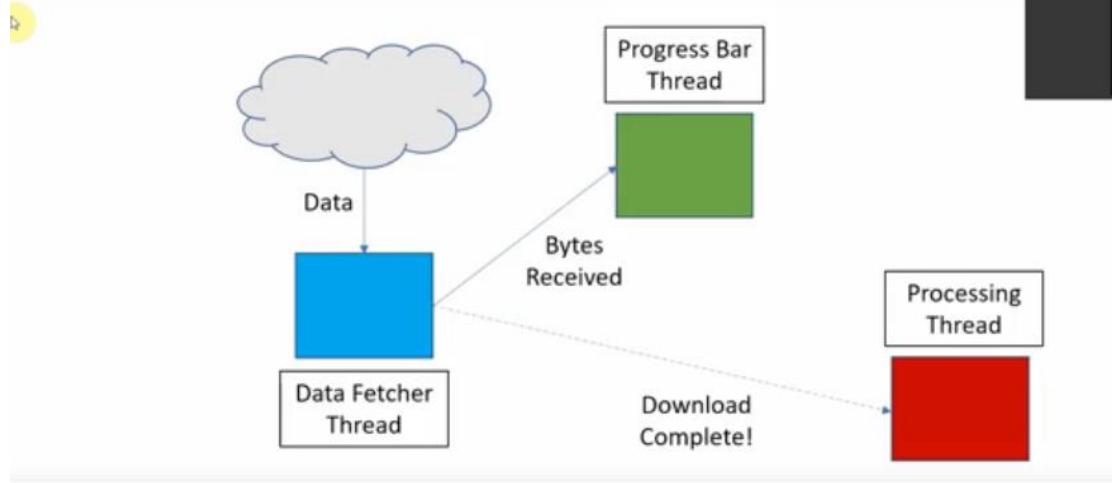
## Coordination Between Threads



- Similar situations arise with threads
- e.g. A program which is performing a download
  - One thread fetches the data over the network
  - Another thread displays a progress bar
  - A third thread will process the data when the download is complete

## Thread Coordination Practical

# Program Performing a Download



## Communication Between Threads

- The threads run concurrently
  - The data fetching thread runs continually
  - The progress bar thread waits for information
  - The processor thread waits until all the data has been received
- When the download is complete
  - The fetching thread terminates
  - The progress bar thread terminates
  - The processor thread runs

## Data Sharing Between Threads

- The downloaded data is shared by all three threads
  - The data fetching thread appends to it
  - The progress bar thread calculates its size
  - The processor thread uses the data
- Potential Data Race
  - Multiple threads
  - Modification

## Coordination of Threads

- We will use two bools to coordinate the threads
- "progress" flag
  - The fetching thread sets this when it has new data
  - The progress bar thread checks this flag
- "completed" flag
  - The fetching thread sets this when it finishes
  - The other two threads check this flag
- Potential Data Race
  - Multiple threads
  - Modification
- Use mutexes

## Hot Loop

- We need to lock the mutex while checking a bool

```
// In progress bar task function  
std::lock_guard data_lck(data_mutex);  
while (!update_progress) {}
```

- The thread will run flat out
  - The processor core will run at 100%
  - Other threads cannot do useful work
  - Uses a lot of electricity
- The fetcher thread cannot set the flag

## Hot Loop Avoidance

- To avoid this, unlock the mutex inside the loop

```
std::unique_lock<std::mutex> data_lck(data_mutex);
```

- ```
while (!update_progress) {  
    data_lck.unlock();  
    std::this_thread::sleep_for(10ms);  
    data_lck.lock();  
}
```
- Sleeping allows other threads to use the core
  - The fetcher thread can set the flag

```
Source.cpp 114 Thread Coordination Practical 01 Progress [Global Scope] main() 95 completed_ 96 std::lock_guard<std::mutex> lock{completed_}; 97 std::cout <> "Processing thread waiting for data..."; 98 sdata = Block1; 99 std::cout <> "Fetcher thread waiting for data..."; 100 Received 6 bytes so far 101 // Process 102 } 103 std::cout <> "Progress bar thread waiting for data..."; 104 int main() 105 { 106 // Start th 107 std::thread fetcher(fetcher_func); 108 std::thread prog(prog_func); 109 std::thread processor(processor_func); 110 fetcher.join(); 111 std::cout <> "Fetcher thread waiting for data..."; 112 prog.join(); 113 std::cout <> "Received 12 bytes so far"; 114 processor.join(); 115 std::cout <> "Progress bar thread has ended"; 116 std::cout <> "Processing sdata: Block1Block2Block3Block4Block5"; 117 std::cout <> "Press any key to continue . . .";
```

## Implementation with Mutex

- This is not ideal
  - Too many loops
  - Too much explicit locking and unlocking
  - How do we choose the sleep duration?
- Better solution
  - Thread A indicates that it is waiting for something
  - Thread B does the "something"
  - Thread A is woken up and resumes

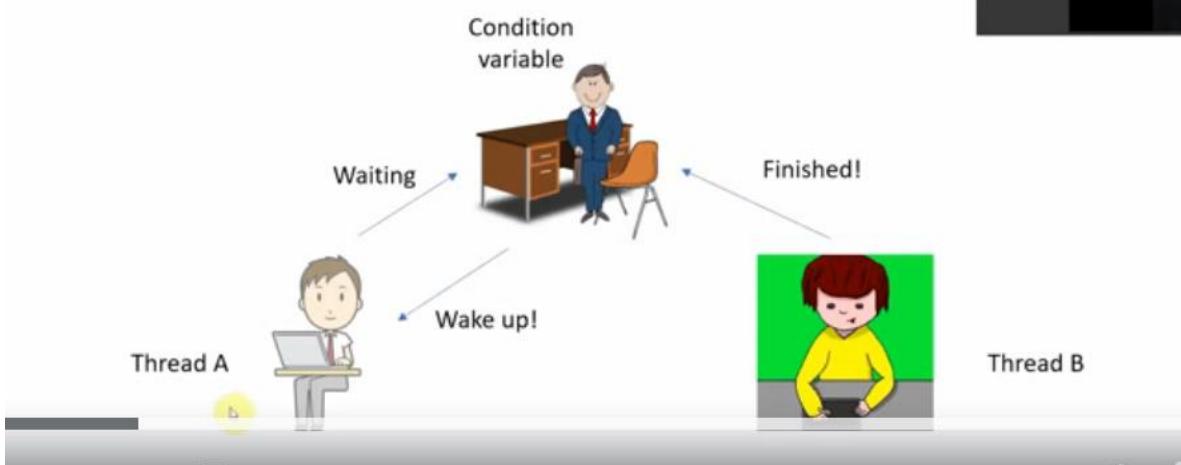
## Condition Variables

# Condition Variable Motivation

- Suppose we have two threads
  - "Writer" thread modifies a shared string
  - "Reader" thread uses the modified string
- The two threads need to be coordinated
- We must also avoid a data race

## Coordination Between Threads

- One solution is to introduce a condition variable



## Condition Variable Overview

- Scenario
  - Thread A tells the condition variable it is waiting
  - Thread B notifies the condition variable when it updates the string
  - The condition variable wakes thread A up
  - Thread A then uses the string

# Condition Variable and Thread Safety

- We use a mutex to protect critical sections
- The condition variable also uses the same mutex
  - Thread coordination
  - No data race

## std::condition\_variable

- Defined in <condition\_variable>
- `wait()`
  - Takes an argument of type `std::unique_lock`
  - It unlocks its argument and blocks the thread until a notification is received
- `wait_for()` and `wait_until()`
  - Re-lock their argument if a notification is not received in time
- `notify_one()`
  - Wake up one of the waiting threads
  - The scheduler decides which thread is woken up
- `notify_all()`
  - Wake up all the waiting threads

## Condition Variable Scenario

- Thread A locks the mutex
  - It calls the condition variable's wait() member function
  - The condition variable unlocks the mutex
  - The condition variable blocks this thread
- Thread B locks the mutex
  - It modifies the string and unlocks the mutex
  - It calls notify\_one()
- The condition variable wakes thread A up
  - The wait() call returns with the mutex locked
  - Thread A resumes execution and uses the string

## Reader thread

```
// Waiting thread
void reader()
{
    // Lock the mutex
    std::unique_lock<std::mutex> uniq_lck(mut);
    // Call wait() on the condition variable
    // Unlocks the mutex and makes this thread sleep
    cond_var.wait(uniq_lck);

    // The condition variable wakes this thread up and locks the mutex

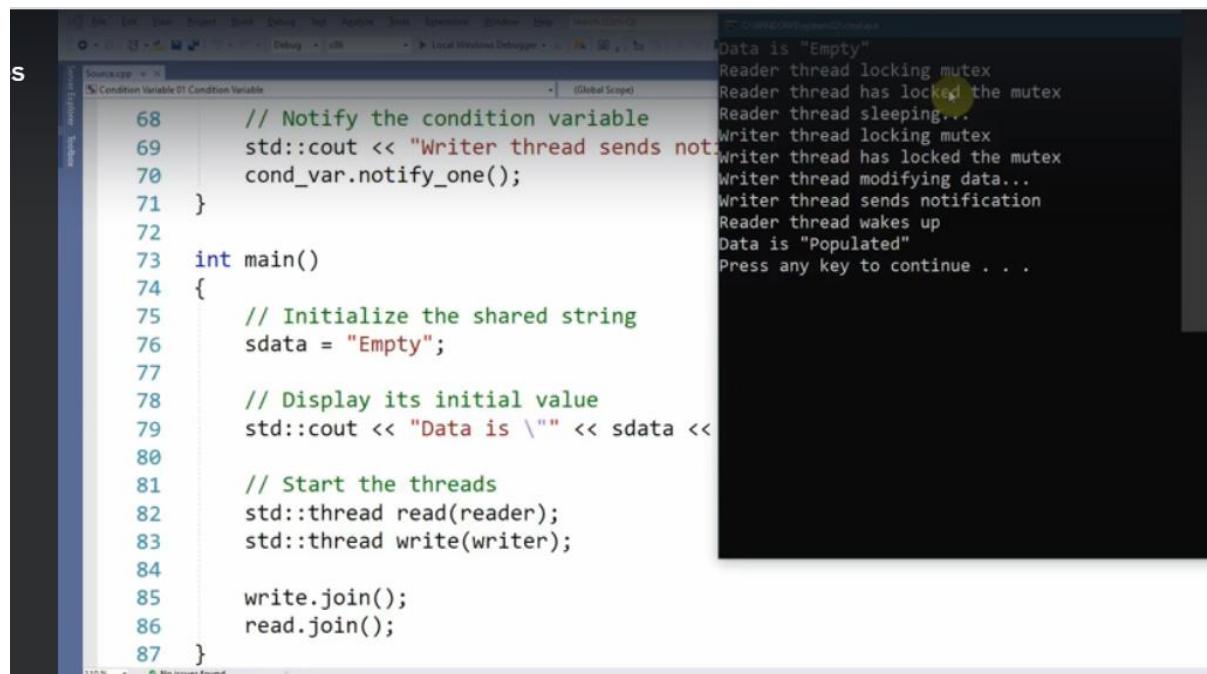
    // Use the shared data
    // ...
}
```

# Writer thread

```
// Notifying thread
void writer()
{
{
    // Lock the mutex
    std::lock_guard<std::mutex> lck_guard(mut);

    // Modify the shared data
    sdata = "Populated";
    } // Release the lock

    // Notify the condition variable
    cond_var.notify_one();
}
```



The screenshot shows a debugger interface with two panes. The left pane displays the C++ code for the writer thread. The right pane shows the console output with the following text:

```
Data is "Empty"
Reader thread locking mutex
Reader thread has locked the mutex
Reader thread sleeping...
Writer thread locking mutex
Writer thread has locked the mutex
Writer thread modifying data...
Writer thread sends notification
Reader thread wakes up
Data is "Populated"
Press any key to continue . . .
```

The code in the left pane highlights the line `cond\_var.notify\_one();` with a yellow circle and arrow, indicating it is the current instruction being executed by the writer thread.

```
1 Source.cpp - [Condition Variable]
2 Condition Variable 01 Condition Variable
3
4
5
6
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
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
```

## std::condition\_variable\_any

- std::condition\_variable only works with std::mutex
  - Does not work with std::timed\_mutex
- There is also std::condition\_variable\_any
  - Works with any mutex-like object
  - Including our own types
  - May have more overhead than std::condition\_variable

## Condition Variables with Predicate

### Lost Wakeup

- The example in the previous lecture has a problem
- wait() will block until the condition variable is notified
- If the writer calls notify() before the reader calls wait()
  - The condition variable is notified when there are no threads waiting
  - The reader will never be woken up
  - The reader could be blocked forever
- This is known as a "lost wakeup"

The screenshot shows a debugger interface with a stack trace and a code editor. The stack trace on the right lists several events:

- Data is "Empty"
- Writer thread locking mutex
- Writer thread has locked the mutex
- Reader thread locking mutex
- Writer thread modifying data...
- Writer thread sends notification
- Reader thread has locked the mutex
- Reader thread sleeping...

The code editor on the left contains C++ code for a condition variable predicate test:

```
64
65     // Notify the condition variable
66     std::cout << "Writer thread sends notification" << endl;
67     cond_var.notify_one();
68 }
69
70 int main()
71 {
72     // Initialize the shared string
73     sdata = "Empty";
74
75     // Display its initial value
76     std::cout << "Data is \\" << sdata << endl;
77
78     // Start the threads
79     // If the writer thread finishes before the reader
80     std::thread write(writer);
81     std::this_thread::sleep_for(500ms);
82     std::thread read(reader);
```

## Spurious Wakeup

- Occasionally, the reader will be "spuriously" woken up
  - The reader thread has called wait()
  - The writing thread has not called notify()
  - The condition variable wakes the reader up anyway
- This is due to the way that std::condition\_variable is implemented
  - Avoiding spurious wakeups adds too much overhead
- Fortunately, there is a way to solve both spurious and lost wakeups



## wait() with Predicate

- wait() takes an optional second argument
  - A predicate
- Typically, the predicate checks a shared bool
  - The bool is initialized to false
  - It is set to true when the writer sends the notification
- The reader thread will call this predicate
- It will only call wait() if the predicate returns false
  - Also available with wait\_for() and wait\_until()

## Using wait() with Predicate

- Add a shared boolean flag, initialized to false
- In the wait() call, provide a callable object that checks the flag

```
// bool flag for predicate  
bool condition = false;  
  
// Waiting thread  
void reader()  
{  
    // Lock the mutex  
    std::unique_lock<std::mutex> uniq_lck(mut);  
  
    // Lambda predicate that checks the flag  
    cond_var.wait(uniq_lck, [] {return condition;});  
  
    // ...  
}
```

## Using wait() with Predicate

- Add a shared boolean flag, initialized to false
- In the wait() call, provide a callable object that checks the flag

```
// bool flag for predicate  
bool condition = false;  
  
// Waiting thread  
void reader()  
{  
    // Lock the mutex  
    std::unique_lock<std::mutex> uniq_lck(mut);  
  
    while (!condition)  
        cond_var.wait();  
  
    // ...  
}
```

## Using wait() with Predicate

- In the writer thread, set the flag to true

```
{  
    std::lock_guard<std::mutex> lck_guard(mut);  
    sdata = "Populated";  
  
    // Set the flag  
    condition = true;  
}  
    b  
  
    // Notify the condition variable  
    cv.notify_one();  
}
```

The screenshot shows a debugger interface with two panes. The left pane displays the source code for a C++ program named 'Source.cpp'. The code contains a main function that initializes a shared string 'sdata' to "Empty", prints its initial value, performs a sleep operation, and then writes to it. The right pane shows a series of log messages indicating the execution flow: the writer thread sends a notification, locks the mutex, modifies the data, unlocks the mutex, and sends a notification. The reader thread then wakes up, unlocks the mutex, and sleeps again. A yellow arrow points to the message 'Reader thread wakes up'.

```
75     // Notify the condition variable
76     std::cout << "Writer thread sends notification" << endl;
77     cond_var.notify_one();
78 }
79
80 int main()
81 {
82     // Initialize the shared string
83     sdata = "Empty";
84
85     // Display its initial value
86     std::cout << "Data is \\" << sdata << endl;
87
88     // The notification is not lost,
89     // even if the writer thread finishes
90     // or there is a "spurious wakeup" (writer)
91     std::thread write(writer);
92     std::this_thread::sleep_for(500ms);
93     std::thread read(reader);
94 }
```

## Lost Wakeup Avoidance

- The writer notifies the condition variable
- The reader thread locks the mutex
- The reader thread calls the predicate
- If the predicate returns true
  - Lost wakeup scenario - the writer has already sent a notification
  - The reader thread continues, with the mutex locked
- If the predicate returns false
  - Normal scenario
  - The reader thread calls wait() again

## Spurious Wakeup Avoidance



- The writer notifies the condition variable (or not)
- The reader thread locks the mutex
- The reader thread calls the predicate
- If the predicate returns true
  - Genuine wakeup - the writing thread really has sent a notification
  - The reader thread continues, with the mutex locked
- If the predicate returns false
  - Spurious wakeup scenario - the writing thread has not sent a notification
  - The reader thread calls wait() again

## Multiple Threads



- Condition variables are particularly useful here
  - Multiple threads are waiting for the same event
- **notify\_all()**
  - The condition variable wakes up all the threads which have called wait()
  - The threads could wake up in any order
  - All the reader threads process the data
- **notify\_one()**
  - Only one of the threads which called wait() will be woken up
  - The other waiting threads will remain blocked
  - A different reader thread processes the data each time

The screenshot shows a debugger interface with two windows. The left window displays the C++ code for a function named `multiple_writers`. The right window shows the output of the program's execution, which includes logs of thread activities such as locking and unlocking mutexes and sending notifications.

```
88 // Display its initial value
89 std::cout << "Data is \\" << sdata <<
90
91 // The notification is not lost,
92 // even if the writer thread finishes
93 // or there is a "spurious wakeup" (w
94
95 std::thread write(writer);
96 std::thread read1(reader);
97 std::this_thread::sleep_for(10ms);
98 std::thread read2(reader);
99 std::this_thread::sleep_for(10ms);
100 std::thread read3(reader);
101 std::this_thread::sleep_for(10ms);
102
103 write.join();
104 read1.join();
105 read2.join();
106 read3.join();
107 }
```

```
Writer thread locking mutex
Writer thread has locked the mutex
Reader thread locking mutex
Reader thread locking mutex
Reader thread locking mutex
Writer thread modifying data...
Writer thread unlocks the mutex
Writer thread sends notification
Reader thread has locked the mutex
Reader thread sleeping...
Reader thread sleeping...
Reader thread 18152 wakes up
Data is "Populated"
Reader thread unlocks the mutex
Reader thread has locked the mutex
Reader thread sleeping...
Reader thread 18976 wakes up
Data is "Populated"
Reader thread unlocks the mutex
Reader thread has locked the mutex
Reader thread sleeping...
Reader thread 18644 wakes up
Data is "Populated"
Reader thread unlocks the mutex
Press any key to continue . . .
```

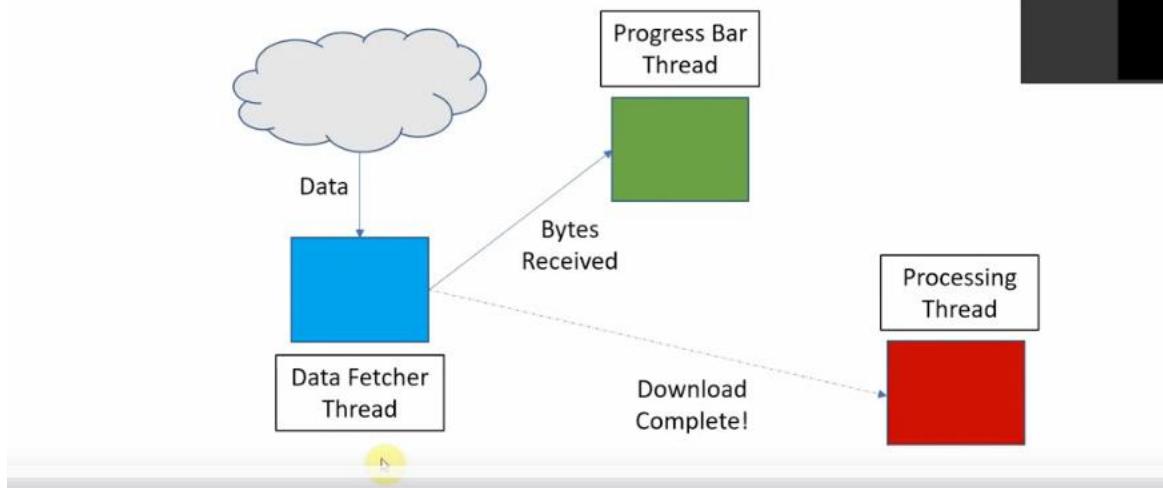
The screenshot shows a debugger interface with two windows. The left window displays the C++ code for a function named `main`. The right window shows the output of the program's execution, which includes logs of thread activities such as locking and unlocking mutexes and sending notifications.

```
68 condition = true;
69
70 std::cout << "Writer thread unlock
71 }
72
73 // Notify the condition variable
74 std::cout << "Writer thread sends not
75
76 //cond_var.notify_all();
77
78 for (int i = 0; i < 2; ++i)
79     cond_var.notify_one();
80
81 }
82
83 int main()
84 {
85     // Initialize the shared string
86     sdata = "Empty";
87 }
```

```
Writer thread locking mutex
Writer thread has locked the mutex
Reader thread locking mutex
Reader thread locking mutex
Reader thread locking mutex
Writer thread modifying data...
Writer thread unlocks the mutex
Writer thread sends notification
Reader thread has locked the mutex
Reader thread sleeping...
Reader thread sleeping...
Reader thread 31240 wakes up
Data is "Populated"
Reader thread unlocks the mutex
Reader thread has locked the mutex
Reader thread sleeping...
Reader thread 30508 wakes up
Data is "Populated"
Reader thread unlocks the mutex
Reader thread has locked the mutex
Reader thread sleeping...
Reader thread 28364 wakes up
Data is "Populated"
Reader thread unlocks the mutex
Press any key to continue . . .
```

# Condition Variable Practical

# Program Performing a Download



## Coordination of Threads

- We will use two condition variables
- "data\_cv"
  - The fetching thread notifies this when it has new data
  - The progress bar waits on it and updates itself
- "completed\_cv"
  - The fetching thread notifies this when the download completes
  - The progress bar waits on it and exits
  - The processing thread waits on it and processes the data
- We use predicates with the condition variables
  - Avoid lost and spurious wake-ups

# Progress Bar

- Implemented as a loop:
  - Wait on data\_cv
  - Update progress
  - Wait on completed\_cv
  - If the download is complete, exit
- Use blocking wait() on data\_cv
- Use non-blocking wait\_for() on completed\_cv

The screenshot shows a debugger interface with a code editor and a terminal window. The code editor displays a C++ file named 'Condition Variable Practical 01 Progress'. The terminal window shows the output of the program, which is a sequence of messages indicating the progress of data fetching and processing by three threads: fetcher, progress\_bar, and processor.

```
96     completed_cv.wait(compl_lck, [] { return true; });
97     compl_lck.unlock();
98
99     std::lock_guard<std::mutex> data_lck(data_mutex);
100    std::cout << "Processing sdata: " << sdata;
101
102    // Process the data...
103 }
104
105 int main()
106 {
107     // Start the threads
108     std::thread fetcher(fetch_data);
109     std::thread prog(progress_bar);
110     std::thread processor(process_data);
111
112     fetcher.join();
113     prog.join();
114     processor.join();
115 }
```

C:\WINDOWS\system32\cmd.exe  
Progress bar thread waiting for data...  
Processing thread waiting for data...  
Fetched sdata: Block1  
Fetcher thread waiting for data...  
Received 6 bytes so far  
Progress bar thread waiting for data...  
Fetched sdata: Block1Block2  
Fetcher thread waiting for data...  
Received 12 bytes so far  
Progress bar thread waiting for data...  
Fetched sdata: Block1Block2Block3  
Fetcher thread waiting for data...  
Received 18 bytes so far  
Progress bar thread waiting for data...  
Fetched sdata: Block1Block2Block3Block4  
Fetcher thread waiting for data...  
Received 24 bytes so far  
Progress bar thread waiting for data...  
Fetched sdata: Block1Block2Block3Block4Block5  
Fetcher thread has ended  
Received 30 bytes so far  
Progress bar thread has ended  
Processing sdata: Block1Block2Block3Block4Blocks  
Press any key to continue . . .

```
Fetcher thread waiting for data...
Progress bar thread waiting for data...
Processing thread waiting for data...
Fetched sdata: Block1
Fetcher thread waiting for data...
Received 6 bytes so far
Fetched sdata: Block1Block2
Fetcher thread waiting for data...
Fetched sdata: Block1Block2Block3
Fetcher thread waiting for data...
Fetched sdata: Block1Block2Block3Block4
Fetcher thread waiting for data...
Fetched sdata: Block1Block2Block3Block4Block5
Fetcher sdata has ended
Progress bar thread has ended
Processing sdata: Block1Block2Block3Block4Block5
Press any key to continue . . .
```

# Futures and Promises

## Overview

### Transferring Data Between Threads

- `std::thread` does not provide a way to return a value
  - So far, we have used a shared variable
  - Access to the shared variable needs to be protected by locks
- Condition variables allow us to coordinate threads
  - A thread can signal to another thread that shared data has been modified
  - Cannot directly transfer data from one thread to another

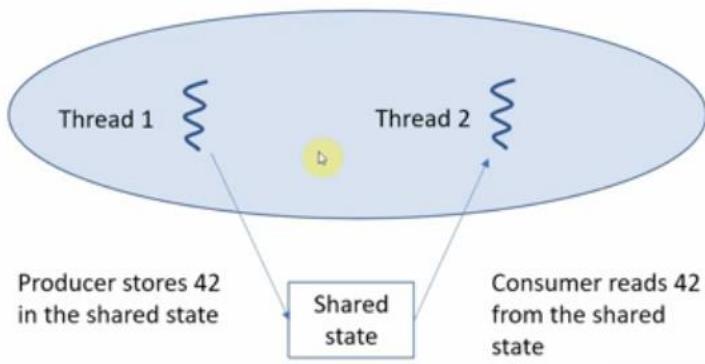
## std::future and std::promise

- Classes for transferring data between threads
- Together, these set up a "shared state" between threads
- The shared state can transfer data from one thread to another
  - No shared data variables
  - No explicit locking

## Producer-Consumer Model

- Futures and promises use a producer-consumer model
  - Reader/writer threads are an example of this model
- A "producer" thread will generate a result
- A "consumer" thread waits for the result
- The producer thread generates the result
- The producer thread stores the result in the shared state
- The consumer thread reads the result from the shared state

## Producer-Consumer Transfer of Data



## Transfer of Data using Future and Promise

- An `std::promise` is associated with the producer
- An `std::future` object is associated with the consumer
  - The consumer calls a member function of the future object
  - The function blocks until the result becomes available
- The producer thread sends the result
  - The promise object stores the result in the shared state
- The consumer thread receives the result
  - The member function reads the result from the shared state
  - The member function returns the result

## Exception Handling

- Futures and promises also work with exceptions
  - The promise stores the exception in the shared state
- This exception will be rethrown in the consumer thread
  - By the future's blocking function
- The producer thread "throws" the exception to the consumer

## std::future and std::promise Classes

## Future and Promise

- A promise object is associated with a future object
- Together, they create a "shared state"
  - The promise object stores a result in the shared state
  - The future object gets the result from the shared state

## std::future

- Represents a result that is not yet available
- One of the most important classes in C++ concurrency
  - Works with many different asynchronous objects and operations
  - Not just std::promise
- An std::future object is not usually created directly
  - Obtained from an std::promise object
  - Or returned by an asynchronous operation

## std::future

- Template class defined in <future>
  - The parameter is the type of the data that will be returned
- get() member function
  - Obtains the result when ready
  - Blocks until the operation is complete
  - Fetches the result and returns it
- wait() and friends
  - Block but do not return a result
  - wait() blocks until the operation is complete
  - wait\_for() and wait\_until() block with a timeout

## std::promise

- Template class defined in <future>
  - The parameter is the type of the result
- Constructor
  - Creates an associated std::future object
  - Sets up the shared state with it
- `get_future()` member function
  - Returns the associated future

```
std::promise<int> prom;
std::future<int> fut = prom.get_future();
```

## std::promise Interface

- `get_future()` member function
  - Returns the std::future object associated with this promise
- `set_value()`
  - Sets the result to its argument
- `set_exception()`
  - Indicates that an exception has occurred
  - This can be stored in the shared state

## Producer-Consumer Model

- Parent thread
  - Creates an std::promise object
- Producer task function
  - Takes the std::promise object as argument
  - Calls set\_value()
    - Or set\_exception()
- Consumer task function
  - Takes the associated std::future object as argument
  - Calls get()
    - Or wait() and friends

## Futures and Promises Examples

### Producer-Consumer Model

- Parent thread
  - Creates an std::promise object
- Producer task function
  - Takes the std::promise object as argument
  - Calls set\_value()
    - Or set\_exception()
- Consumer task function
  - Takes the associated std::future object as argument
  - Calls get()
    - Or wait() and friends

## Producer Thread Example

```
// The producer's task function takes an std::promise as argument
void produce(std::promise<int>& px)
{
    // Produce the result
    int x = 42;

    // Store the result in the shared state
    px.set_value(x);
}
```

## Consumer Thread Example

```
// The consumer's task function takes an std::future as argument
void consume(std::future<int>& fx)
{
    // Get the result from the shared state
    int x = fx.get();
}
```

# Parent Thread Example

```
// Create an std::promise object
std::promise<int> prom;

// Get the associated future
std::future<int> fut = prom.get_future();

// The producer task function takes the promise as argument
std::thread thr_producer(produce, std::ref(prom));

// The consumer task function takes the future as argument
std::thread thr_consumer(consume, std::ref(fut));
```

The screenshot shows a C++ development environment with the following details:

- IDE:** Visual Studio (VS Code) interface.
- File:** Source.cpp -> Futures and Promises 01 Producer Consumer
- Code (Visible Lines):**

```
13     // Produce the result
14     int x = 42;
15     std::this_thread::sleep_for(2s);
16
17     // Store the result in the shared state
18     std::cout << "Promise sets shared state to 42"
19     px.set_value(x);
20 }
```

```
23 // The consumer's task function takes an std::future<int>
24 void consume(std::future<int>& fx)
25 {
26     // Get the result from the shared state
27     std::cout << "Future calling get()..." I
28     int x = fx.get();
29     std::cout << "Future returns from calling get()\n";
30     std::cout << "The answer is " << x << '\n';
31 }
```

```
32
```
- Output Terminal:**

```
C:\WINDOWS\system32\cmd.exe
Future calling get()...
Promise sets shared state to 42
Future returns from calling get()
The answer is 42
Press any key to continue . . .
```

## Producer-Consumer with Exception Handling

- In the producer thread
  - Put a try block around code that might throw
  - In the catch block, call `set_exception()` on the promise
  - This captures the active exception
- `set_exception()` takes a pointer to the exception object
  - We can use a catch-all handler
  - Pass the return value from `std::current_exception()`

## Producer-Consumer with Exception Handling

- In the consumer thread
  - Put a try block around the call to `get()` or `wait()`
  - Write a catch block to handle the exception

## Producer with Exception Handling

```
void produce(std::promise<int>& px)
{
    try {
        // Code that may throw
        ...
        // If no exception, store the result in the shared state
        px.set_value(x);
    }
    catch (...) {
        // Exception caught - store it in the shared state
        px.set_exception(std::current_exception());
    }
}
```

# Consumer with Exception Handling

```
void consume(std::future<int>& fx)
{
    try {
        // Get the result from the shared state - may throw
        int x = fx.get(); I
    }
    catch (...) {
        // Exception thrown - get it from the shared state
    }
}
```

The screenshot shows a debugger interface with two panes. The left pane is a code editor for a file named 'Source.cpp' containing C++ code related to futures and promises. The right pane is a terminal window showing the output of the program's execution.

**Code Editor (Source.cpp):**

```
27     // Exception thrown - store it in
28     px.set_exception(std::current_exception());
29 }
30 }
31
32 // The consumer's task function takes an &reference
33 void consume(std::future<int>& fx)
34 {
35     std::cout << "Future calling get()..." I
36     try {
37         // Get the result from the shared state
38         int x = fx.get(); I
39         std::cout << "Future returns from "
40         std::cout << "The answer is " << x;
41     }
42     catch (std::exception& e) {
43         // Exception thrown - get it from the shared state
44         std::cout << "Exception caught: " << e.what() << '\n';
45     }
46 }
```

**Terminal Output:**

```
Future calling get()... I
Exception caught: Oops
Press any key to continue . . .
```

```
C:\WINDOWS\system32\cmd.exe
Future calling get()...
Promise sets shared state to 42
Future returns from calling get()
The answer is 42
Press any key to continue . . .

Source.cpp  x  Futures and Promises 02 Exception
7 #include <exception>
8 #include <chrono>
9
10 // The producer's task function takes a std::promise<int> parameter
11 void produce(std::promise<int>& px)
12 {
13     try {
14         using namespace std::literals;
15         int x = 42;
16         std::this_thread::sleep_for(2s);
17
18         // Code that may throw
19         if (0)
20             throw std::out_of_range("Oops");
21
22         // No exception - store the result
23         std::cout << "Promise sets shared state to " << x << '\n';
24         px.set_value(x);
25     }
26     catch (...) {
```

The screenshot shows a Windows command prompt window titled 'cmd.exe' with the path 'C:\WINDOWS\system32\cmd.exe'. The output of the program is displayed:  
Future calling get()...  
Promise sets shared state to 42  
Future returns from calling get()  
The answer is 42  
Press any key to continue . . .

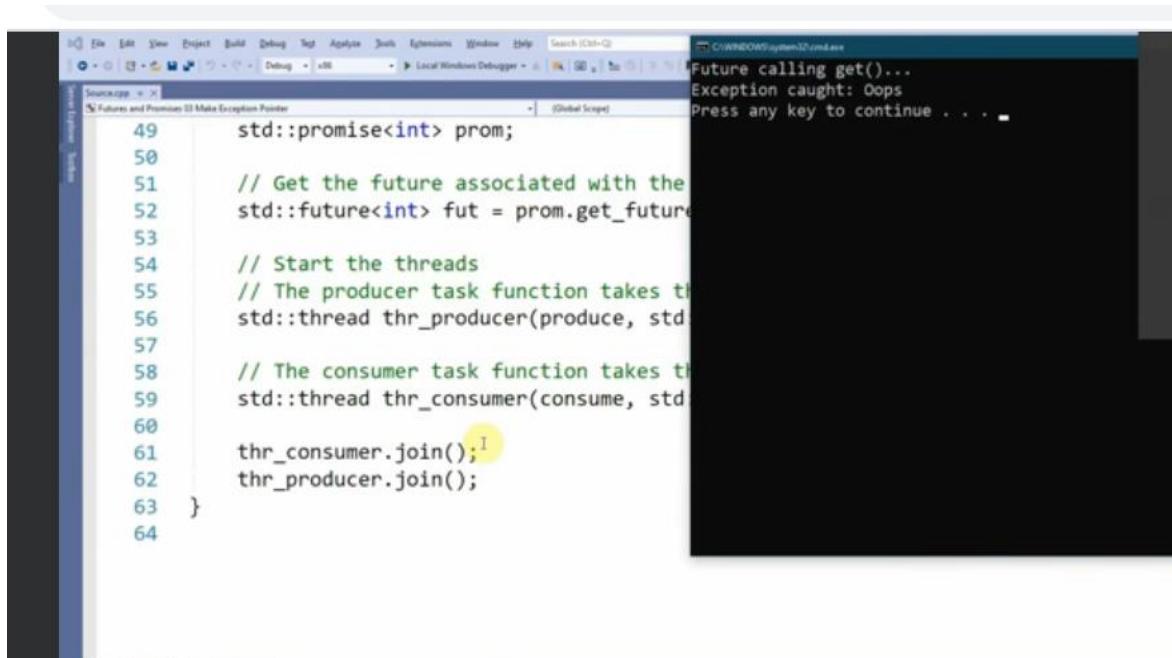
## Producer with std::make\_exception\_ptr()

- To throw an exception ourselves, we could
  - Throw the exception inside a try block
  - Write a catch block that calls set\_exception()
- C++11 has std::make\_exception\_ptr()
  - Takes the exception object we want to throw
  - Returns a pointer to its argument
  - Pass this pointer to set\_exception()
  - Avoid "boilerplate" code
  - Better code generation

# Producer with std::make\_exception\_ptr()

```
void produce(std::promise<int>& px)
{
    ...
    if (...) {
        px.set_exception(std::make_exception_ptr(std::out_of_range("Oops")));
        return;
    }

    // Store the result in the shared state
    px.set_value(x);
}
```



A screenshot of a Windows debugger interface. The title bar says "Future calling get()...". The status bar at the bottom right says "Exception caught: Oops Press any key to continue . . .". The main window shows a C++ code editor with the file "Source.cpp" open. The code is as follows:

```
49     std::promise<int> prom;
50
51     // Get the future associated with the
52     std::future<int> fut = prom.get_future();
53
54     // Start the threads
55     // The producer task function takes the
56     std::thread thr_producer(produce, std::ref(prom));
57
58     // The consumer task function takes the
59     std::thread thr_consumer(consume, std::ref(fut));
60
61     thr_consumer.join();I
62     thr_producer.join();
63 }
64
```

The line "thr\_consumer.join();<sup>I</sup>" is highlighted with a yellow circle, indicating the point where the exception was caught.

# Promises and Multiple Waiting Threads

## Single Producer with Multiple Consumers

- Single producer thread
  - Produces a result or an event
- Multiple consumer threads
  - Use the result
  - Or wait for the event to occur
- Used in many applications

## std::future and Multiple Waiting Threads

- Designed for use with a single consumer thread
  - Assumes it has exclusive read access to the shared state
- Cannot be safely shared between threads
  - Data race
- Cannot be copied
  - Move-only class

The screenshot shows a Microsoft Visual Studio interface. In the center, a code editor displays a file named 'Source.cpp' with the following content:

```
41 std::future<int> fut = prom.get_future();
42
43 // Start the threads
44 // The producer task function takes the future and returns a shared state
45 std::thread thr_producer(produce, std::move(fut));
46
47 // The consumer task function takes the shared state and returns the final answer
48 std::thread thr_consumer(consume, std::move(prom));
49 std::thread thr_consumer2(consume, std::move(prom));
50
51 thr_consumer.join();
52 thr_consumer2.join();
53 thr_producer.join();
54 }
```

To the right of the code editor, a command prompt window shows the output of the program:

```
Future calling get()...
Future calling get()...
Promise sets shared state to 42
Future returns from calling get()
The answer is 42
```

Below the command prompt, a 'Microsoft Visual C++ Runtime Library' dialog box is open, showing a 'Debug Error' message:

Program: ..\00 Future\Debug\Promise with Multiple Futures
00 Future.exe  
abort() has been called  
(Press Retry to debug the application)

Buttons for Abort, Retry, and Ignore are visible at the bottom of the dialog.

## std::shared\_future

- Can be copied
  - Each thread has its own object
  - They all share the same state with the std::promise
  - Calling get() or wait() from different copies is safe

## Obtaining an std::shared\_future object

- Normally, we do not create a shared future directly
- We can move from an existing std::future

```
std::shared_future<int> shared_fut1 = std::move(fut);
```
- We can call share() on the std::future

```
std::shared_future<int> shared_fut2 = fut.share();
```
- We can also obtain a shared\_future directly from a promise

```
shared_future<int> shared_fut3 = prom.get_future();
```

## std::shared\_future Example

- The producer will be the same as before
- The consumer now takes an std::shared\_future

```
void consume(std::shared_future<int>& fx);
```

## std::shared\_future Example

```
// Parent thread
std::promise<int> prom;

// Move the promise's future into a shared future
std::shared_future<int> shared_fut1 = prom.get_future();

// Copy the shared future object
std::shared_future<int> shared_fut2 = shared_fut1;

// Start two consumer threads, each with its own shared future object
std::thread thr_consumer1(consume, std::ref(shared_fut1));
std::thread thr_consumer2(consume, std::ref(shared_fut2));
std::thread thr_producer(produce, std::ref(prom));
```

The screenshot shows the Microsoft Visual Studio IDE interface. The top menu bar includes File, Edit, View, Project, Build, Debug, Test, Analyze, Tools, Extensions, Window, Help, and Search (Ctrl+Q). The title bar says "Source.cpp - e:\..." and "Promise with Multiple Futures 01 Shared Future". The left sidebar has "Server Explorer", "Toolbox", and "Solution Explorer" sections. The main code editor window contains the C++ code provided above. To the right of the editor is the "Immediate Window" which shows the following output:

```
C:\WINDOWS\system32\cmd.exe
Thread 13772 calling get()...
Thread 5924 calling get()...
Promise sets shared state to 42
Thread 5924 returns from calling get()
Thread 5924 has answer 42
Thread 13772 returns from calling get()
Thread 13772 has answer 42
Press any key to continue . . .
```

# Integer Operations and Threads

## Integer Operations and Threads

- Integer operations are usually a single instruction
  - True on x64
  - Provided the data is correctly aligned and fits into a single word
- A thread cannot be interrupted while performing integer operations
- Do we still need to lock a shared integer?



```
Source.cpp + x
Integer Operations and Threads 01 Data Race
224594
Press any key to continue . . .

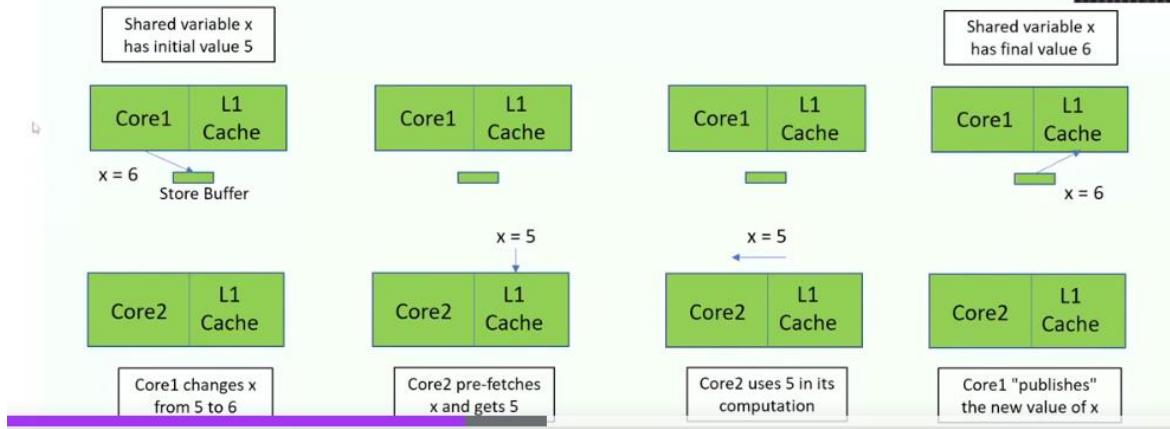
11  {
12      for (int i = 0; i < 100'000; ++i) {
13          ++counter;
14      }
15  }
16
17 int main()
18 {
19     std::vector<std::thread> tasks;
20
21     for (int i = 0; i < 10; ++i)
22         tasks.push_back(std::thread(task));
23
24     for (auto& thr: tasks)
25         thr.join();
26
27     std::cout << counter << '\n';
28 }
29
```

# Integer Operation Discussion



- The output is a number like 258413
  - Data race
- The `++` operation is a single instruction
- However, `++count` involves three operations
  - Pre-fetch the value of count
  - Increment the value in the processor core's register
  - Publish the new value of count
- The thread could use a stale value in its calculation
- The thread could publish its result after a thread which ran later

# Synchronization Issues



## Data Race Scenario

- Threads A and B run on different processor cores
  - Thread A's core pre-fetches the value 5 for count
  - Thread B's core pre-fetches the value 5 for count
  - Thread A's core increments the value to 6
  - Thread B's core increments the value to 6
  - Thread A publishes the value 6 for count
  - Thread B publishes the value 6 for count
- count initially had the value 5
  - It was incremented twice
  - It has the final value 6, when it should be 7

## Thread Synchronization

- We need to make sure that
  - Thread B uses the latest value for count
  - Thread A publishes its result immediately
- A mutex does this internally when we call lock() and unlock()
- Here, we can do this by declaring count as "atomic"

The screenshot shows the Microsoft Visual Studio IDE. In the center, there is a code editor window titled "Source.cpp" containing the following C++ code:

```
// A shared variable is modified by multiple threads
// Use a mutex to prevent a data race
#include <thread>
#include <iostream>
#include <vector>
#include <mutex>

std::mutex mut;
int counter = 0;

void task()
{
    for (int i = 0; i < 100'000; ++i) {
        std::lock_guard<std::mutex> lck(mut);
        ++counter;
    }
}

int main()
```

To the right of the code editor is a terminal window titled "C:\WINDOWS\system32\cmd.exe" showing the output of the program:

```
1000000
Press any key to continue . . .
```

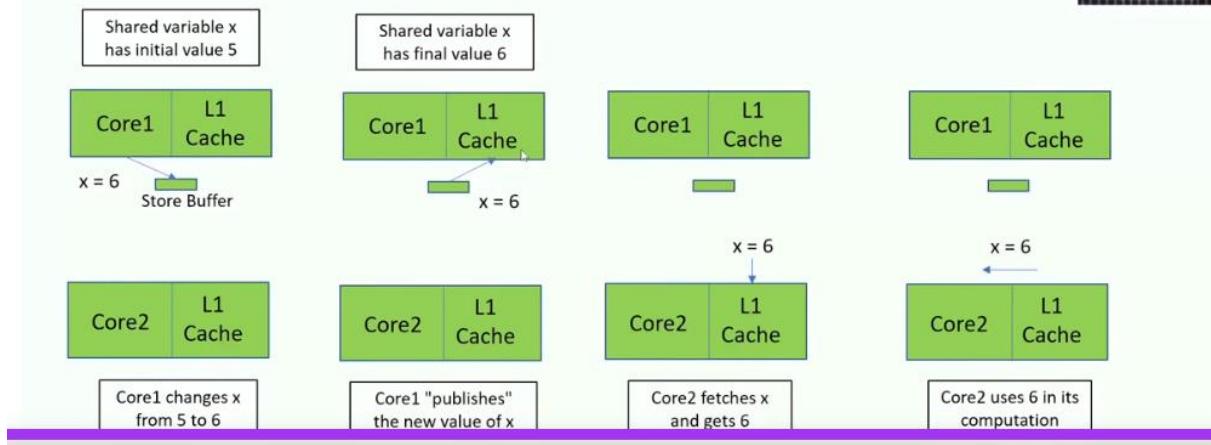
## Thread Synchronization

- We need to make sure that
  - Thread B uses the latest value for count
  - Thread A publishes its result immediately
- A mutex does this internally when we call lock() and unlock()
- Here, we can do this by declaring count as "atomic"

# Atomic Keyword

- The compiler will generate special instructions which
  - Disable pre-fetch for count
  - Flush the store buffer immediately after doing the increment
- This also avoids some other problems
  - Hardware optimizations which change the instruction order
  - Compiler optimizations which change the instruction order
- The result is that only one thread can access count at a time
- This prevents the data race
  - It also makes the operation take much longer

## Synchronization with Atomic Variable



## Scenario without Data Race

- Threads A and B run on different processor cores
  - Thread A's core fetches the value 5 for count
  - Thread A's core increments the value to 6
  - Thread A publishes the value 6
  - Thread B's core fetches the value 6 for count
  - Thread B's core increments the value to 7
  - Thread B publishes the value 7
- count initially had the value 5
  - It was incremented twice
  - It has the final value 7, which is what it should be

## Atomic Types

### Atomic Types

- All operations on the variable will be atomic
- C++11 defines an atomic template
  - In the `<atomic>` header
  - The parameter is the type of the object
  - The object must be initialized

```
// Atomic int, initialized to 0
atomic<int> x = 0;
```

## Atomic Types

- The parameter must be a type which is "trivially copyable"
  - Scalar type
  - Class where all the copy and move constructors are trivial
- Normally only integer types and pointers are used
- For more complex types, locks may be silently added
  - Locking a mutex takes longer
  - To avoid this, use a pointer to the type

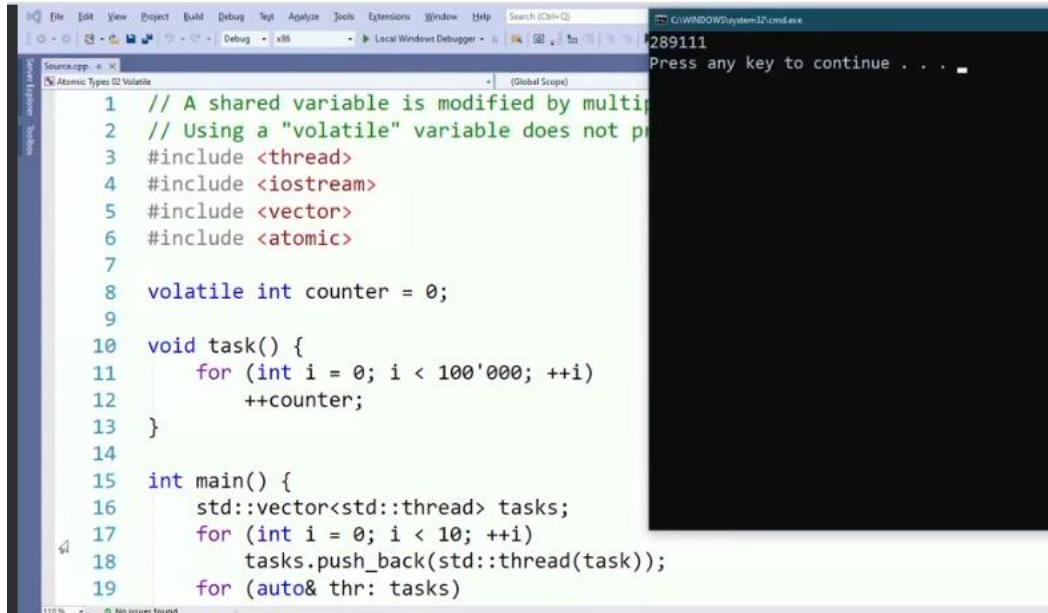
## Using an std::atomic<T> Object

- We can assign to and from the object
  - ↳ 

```
x = 2;           // Atomic assignment to x
y = x;           // Atomic assignment from x. y can be non-atomic
```
- These are two distinct atomic operations
  - Other threads could interleave between them
- Operations such as ++ are atomic
  - Fetch old value
  - Increment
  - Store new value

# Volatile Keyword

- May change without explicit modification
  - Prevents some compiler optimizations
  - Typically used when accessing hardware
- Often mistakenly used in threading
  - Some programmers expect the Java/C# behaviour
  - Has no impact on thread safety
- At one point, Visual Studio supported this in C++
  - Removed before C++11



The screenshot shows a Visual Studio interface. On the left is the Source Editor with the file 'Source.cpp' open, containing C++ code. The code uses the 'volatile' keyword to ensure thread safety for a shared variable. On the right is a terminal window showing the output of the program, which is the number '289111' followed by the instruction 'Press any key to continue . . .'. The terminal window title is 'C:\WINDOWS\system32\cmd.exe'.

```
// A shared variable is modified by multiple threads
// Using a "volatile" variable does not prevent
// the compiler from optimizing the code.
#include <thread>
#include <iostream>
#include <vector>
#include <atomic>
volatile int counter = 0;
void task() {
    for (int i = 0; i < 100'000; ++i)
        ++counter;
}
int main() {
    std::vector<std::thread> tasks;
    for (int i = 0; i < 10; ++i)
        tasks.push_back(std::thread(task));
    for (auto& thr: tasks)
```

## Double-checked Locking

- One solution is to make the initialized object atomic  

```
atomic<Test *> ptest = nullptr;
```
- Atomic types do not support the . or -> operators
- We must copy to a non-atomic pointer before we can use it  

```
Test *ptr = ptest;
ptr->func();
```

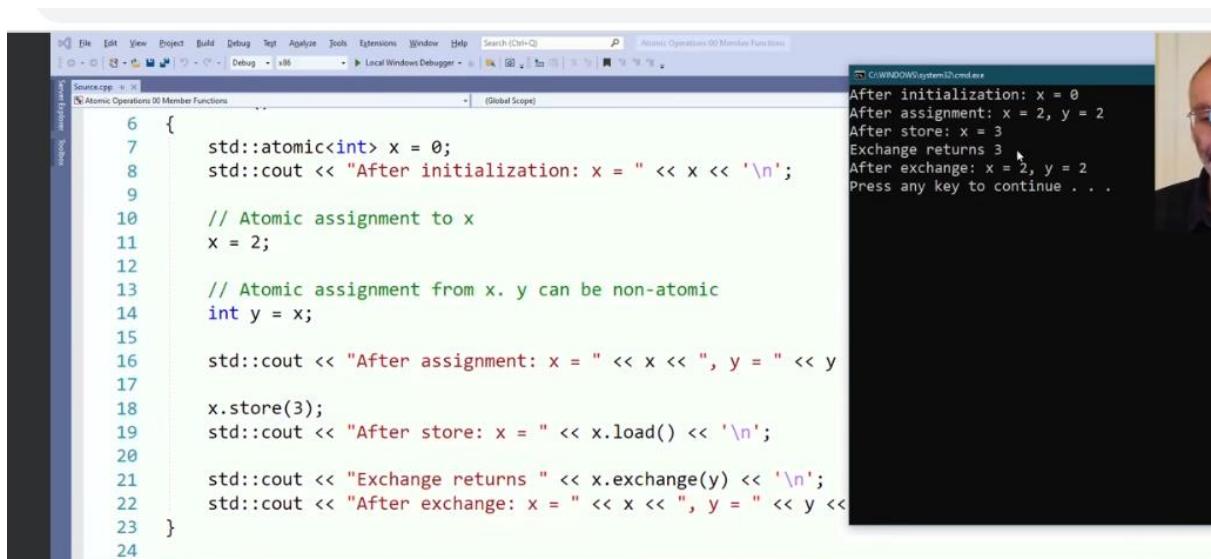
## Atomic Operations

### Member Functions for Atomic Types

- **store()**
  - Atomically replace the object's value with its argument
- **load()**
  - Atomically return the object's value
- **operator =()**
- **operator T()**
  - Synonyms for store() and load()
- **exchange()**
  - Atomically replace the object's value with its argument
  - Returns the previous value

# Member Functions for Specializations

- Atomic pointers support pointer arithmetic
  - increment and decrement operators
  - fetch\_add() synonym for x++
  - fetch\_sub() synonym for x--
  - += and -= operators
- Integer specializations have these, plus
  - Atomic bitwise logical operations &, | and ^



```
Source.cpp: 6
6  {
7      std::atomic<int> x = 0;
8      std::cout << "After initialization: x = " << x << '\n';
9
10     // Atomic assignment to x
11     x = 2;
12
13     // Atomic assignment from x. y can be non-atomic
14     int y = x;
15
16     std::cout << "After assignment: x = " << x << ", y = " << y
17
18     x.store(3);
19     std::cout << "After store: x = " << x.load() << '\n';
20
21     std::cout << "Exchange returns " << x.exchange(y) << '\n';
22     std::cout << "After exchange: x = " << x << ", y = " << y <<
23 }
24
```

After initialization: x = 0  
After assignment: x = 2, y = 2  
After store: x = 3  
Exchange returns 3  
After exchange: x = 2, y = 2  
Press any key to continue . . .

## std::atomic\_flag

- std::atomic\_flag is an atomic boolean type
  - Has less overhead than std::atomic<bool>
- Only three operations
  - clear() sets flag to false
  - test\_and\_set() sets flag to true and returns previous value
  - operator =()
- Must be initialized to false

```
atomic_flag lock = ATOMIC_FLAG_INIT;
```

## Spin Lock

- A spin lock is essentially an infinite loop
  - It keeps "spinning" until a condition becomes true
- An alternative to locking a mutex or using a condition variable
- We can use `std::atomic_flag` to implement a basic spin lock
  - The loop condition is the value of the flag

## Spin Lock with `std::atomic_flag`

- Each thread calls `test_and_set()` in a loop
- If this returns true
  - Some other thread has set the flag and is in the critical section
  - Iterate again
- If it returns false
  - This thread has set the flag
  - Exit the loop and proceed into the critical section
- After the critical section, set the flag to false
  - This allows another thread to execute in the critical section

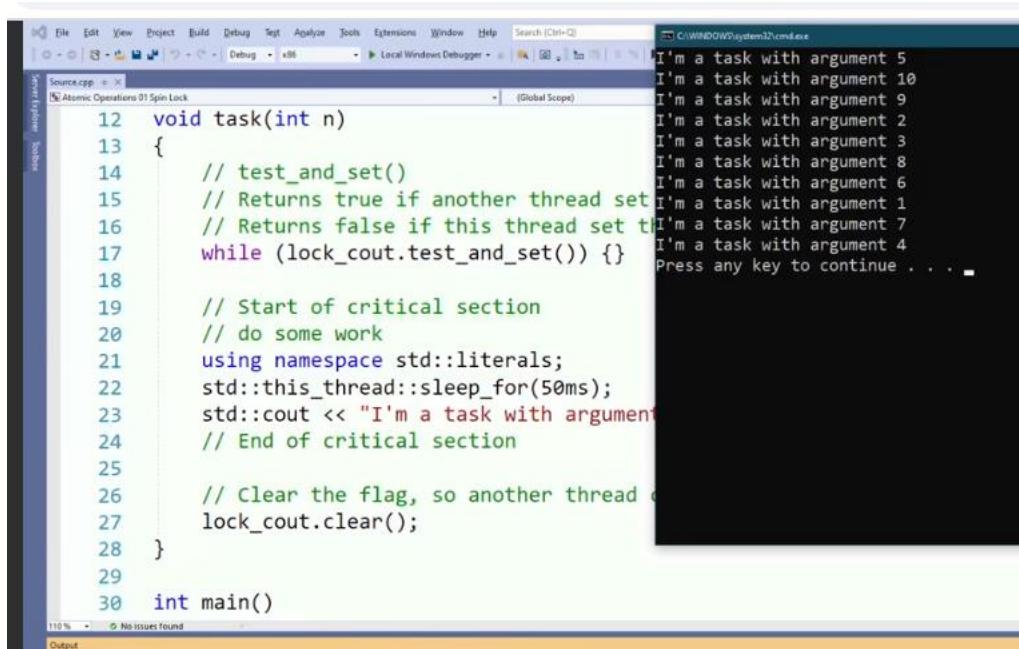
# Spin lock

```
// Initialize flag to false
std::atomic_flag flag = ATOMIC_FLAG_INIT;

void task(int n)
{
    // Loop until we can set the flag
    while (flag.test_and_set()) {}

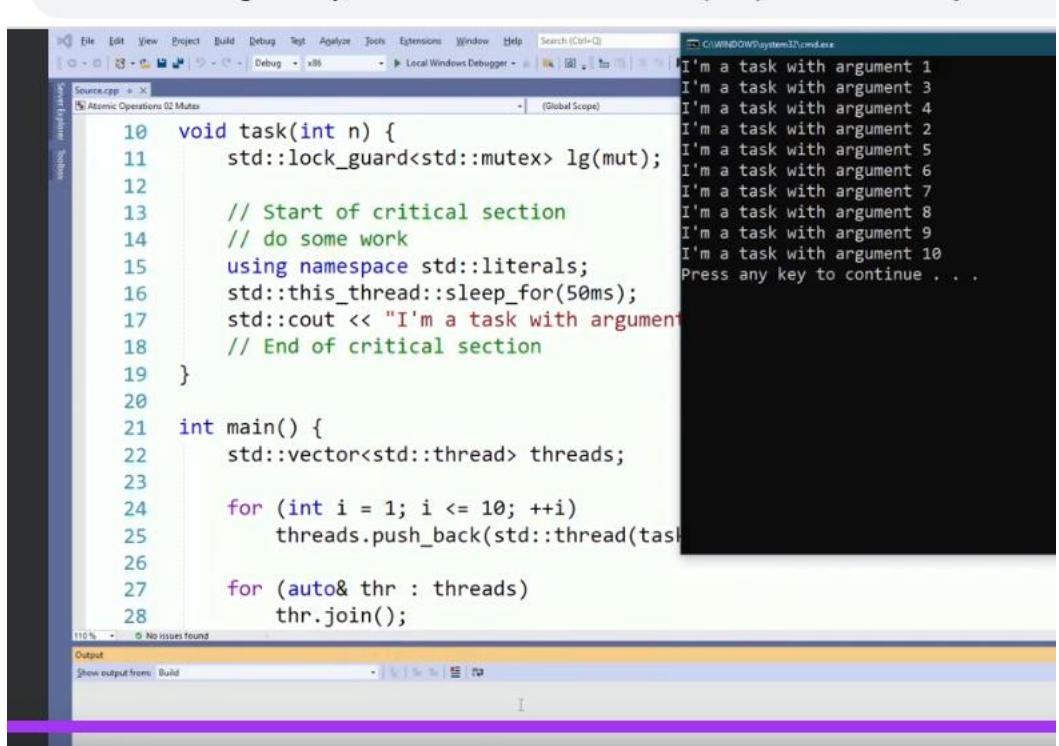
    // Critical section
    ...

    // Clear the flag
    flag.clear();
}
```



The screenshot shows a debugger interface with two panes. The left pane displays the source code for `Source.cpp`, specifically the `task` function which uses a spin lock. The right pane shows the output of the application, which is a series of messages printed by different threads, each containing an argument number from 1 to 10.

```
I'm a task with argument 5
I'm a task with argument 10
I'm a task with argument 9
I'm a task with argument 2
I'm a task with argument 3
I'm a task with argument 8
I'm a task with argument 6
I'm a task with argument 1
I'm a task with argument 7
I'm a task with argument 4
Press any key to continue . . .
```



```
Source.cpp  +  C:\WINDOWS\system32\cmd.exe
File Edit View Project Build Debug Test Analyze Tools Extensions Window Help Search (Ctrl+Q)
Debug - x86 Local Windows Debugger ...
Source.cpp (Global Scope)
10 void task(int n) {
11     std::lock_guard<std::mutex> lg(mut);
12
13     // Start of critical section
14     // do some work
15     using namespace std::literals;
16     std::this_thread::sleep_for(50ms);
17     std::cout << "I'm a task with argument " << n << endl;
18     // End of critical section
19 }
20
21 int main() {
22     std::vector<std::thread> threads;
23
24     for (int i = 1; i <= 10; ++i)
25         threads.push_back(std::thread(task, i));
26
27     for (auto& thr : threads)
28         thr.join();
}
110 %  No issues found
Output Show output from: Build
```

I'm a task with argument 1  
I'm a task with argument 3  
I'm a task with argument 4  
I'm a task with argument 2  
I'm a task with argument 5  
I'm a task with argument 6  
I'm a task with argument 7  
I'm a task with argument 8  
I'm a task with argument 9  
I'm a task with argument 10  
Press any key to continue . . .

## Pros and cons of spin lock



- A spinning thread remains active
  - A mutex may block the thread
- It can continue immediately when it "gets the lock"
  - With a mutex, the thread may need to be reloaded or woken up
- Processor-intensive
  - Only suitable for protecting very short critical sections
  - And/or very low contention
  - Performance can be heavily impacted if spinning threads interrupt each other
  - Usually only used in operating systems and libraries

## Hybrid Mutex



- Often used to implement std::mutex
- Start with a spin lock with a time out
  - If the thread sets the flag in time, enter the critical section
  - If the thread cannot set the flag in time, use the normal mutex implementation
- This gives better performance than the conventional implementation

## Lock-free Programming



## Lock-free Programming

- Threads execute critical sections concurrently
  - Without data races
  - But without using the operating system's locking facilities
- Avoids or reduces some of the drawbacks to using locks
  - Race conditions caused by forgetting to lock, or using the wrong mutex
  - Lack of composability
  - Risk of deadlock
  - High overhead
  - Lack of scalability caused by coarse-grained locking
  - Code complexity and increased overhead caused by fine-grained locking

## Locking vs Lock-free

- Both programming styles are used to manage shared state
  - Analogous to managing a traffic intersection
- Locks
  - Traffic lights control access
  - Stop and wait until able to proceed into critical section



## Locking vs Lock-free



- Lock-free
  - Motorway-style intersection
  - Traffic from different levels can go over the same section at the same time
  - Traffic from one level can merge with traffic from a different level without stopping
  - If not done carefully, collisions can occur!



## Advantages of Lock-free Programming

- If done correctly, threads can never block each other
  - No possibility of deadlock or livelock
  - If a thread is blocked, other threads can continue to execute
  - Useful if work must be completed within a time limit
  - (e.g. real time systems)

## Drawbacks of Lock-free Programming



- Very difficult to write code which is correct and efficient
- The extra complexity makes it unsuitable for many applications
  - e.g. user interface code with separation of concerns
  - May be useful in performance-critical code, such as infrastructure
- Should be used only if
  - A data structure in the program is subject to high contention
  - Which causes unacceptable performance
  - And the lock-free version brings performance up to acceptable levels

## Lock-free Programming Continued

## The Everyday World of Programming with Locks

- We can make some very useful assumptions
- Global state is consistent
  - Provided we only access shared data inside a locked region
  - No other threads will see our changes
  - Until the lock is released
- Logical consistency
  - When working inside a locked region, global state will not change
  - e.g. between evaluating an "if" statement and executing the body
- Code order
  - Statements will execute in the same order as in the source code
  - Or at least, they will appear to...

## The Strange World of Lock-free Programming

- None of these assumptions apply to lock-free programs
  - Shared data may have different values in different threads
  - The value may change between an "if" statement and its body
  - Statements may execute in a different order from the source code

## Bank Account Transfer



## Transactions

- Transactional model of lock-free programming
  - "ACID"
- Atomic/All-or-nothing
  - A transaction either completes successfully ("commit")
  - Or it fails and leaves everything as it was ("rollback")
- Consistent
  - The transaction takes the database from one consistent state to another
  - As seen by other users, the database is never in an inconsistent state

## Transactions

- Transactional model of lock-free programming
  - "ACID"
- Isolated
  - Two transactions can never work on the same data simultaneously
- Durable
  - Once a transaction is committed, it cannot be overwritten
  - ...until the next transaction sees the result of the commit
  - There is no possibility of "losing" an update

## Transactional Memory

- Put shared data in transactional memory
- All operations on shared data will be transactional
- However, there is no standard implementation in C++

## Atomic Instructions

- Atomic/All-or-nothing
  - An instruction completes successfully (we hope!)
- Consistent
  - The instruction takes the data from one consistent value to another
  - As seen by other threads, the data never has an inconsistent value
- Isolated
  - Two atomic instructions can never work on the same data simultaneously
- Durable
  - Once an instruction is completed, the data cannot be overwritten
  - ...until the next instruction sees the result
  - There is no possibility of "losing" a modification

## Lock-free Programming

- To achieve this in lock-free programming
  - Use atomic instructions
- We need to think very carefully about thread interactions
  - Other threads can interleave between each statement
  - Or between expressions within statements
  - How do concurrent writers interact with each other?
  - How do concurrent writers interact with concurrent readers?

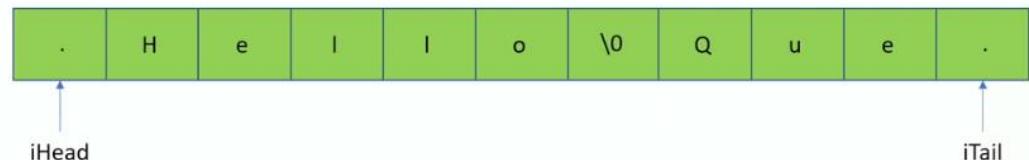
## Lock-free Programming Practical

# Lock-free Queue

- We will implement a simple queue
  - No internal or external locks
- The queue is only accessed by two threads
  - A Producer thread inserts elements into the queue
  - A Consumer thread removes elements from the queue
- The code is carefully designed
  - The Consumer and Producer threads never work on adjacent elements
  - The two threads always work on different parts of the queue

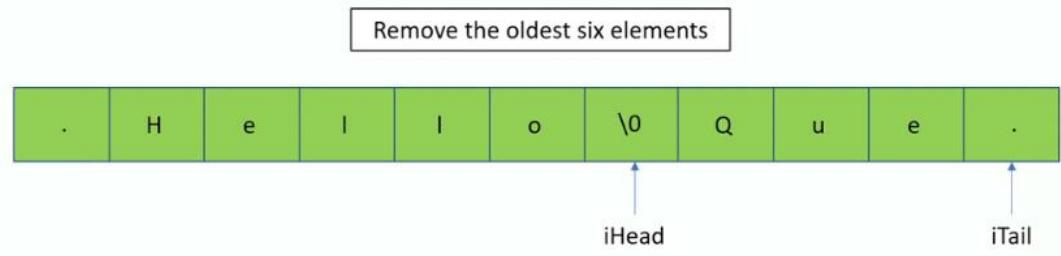
# Lock-free Queue

- The queue has two iterators, iHead and iTail
  - iHead points to the element before the oldest element
  - iTail points to the element after the newest (most recently added)



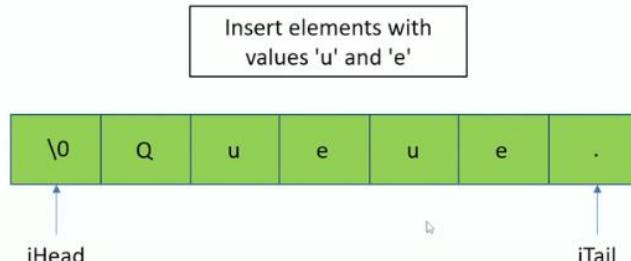
## Consumer Thread

- The Consumer thread does not modify the queue
  - It "removes" an element by incrementing iHead



## Producer Thread

- The Producer inserts elements
  - It increments iTail
  - It also erases any elements which the Consumer has removed



## Thread Separation

- Only the Producer thread can modify the queue
  - The Producer queue inserts elements
  - The Producer queue erases elements
- The two threads never overlap
  - iHead and iTail never refer to the same element
  - The Producer thread never modifies iHead
  - The Consumer thread never accesses elements after iHead

## Lock-free Queue Class

```
template <typename T>
struct LockFreeQueue
{
    private:
        std::list<T> list;
        typename std::list<T>::iterator iHead, iTail;

    public:
        LockFreeQueue() {
            list.push_back(T());           // Create a "dummy" element
            iHead = list.begin();
            iTail = list.end();
        }
}
```

## Consumer Task Member Function

```
bool Consume(T& t)
{
    auto iFirst = iHead;           // Go to the first element
    ++ iFirst;
    if (iFirst != iTail) {         // If queue is not empty
        iHead = iFirst;           // Update iHead
        t = *iHead;               // Fetch this first element
        return true;
    }
    return false;                 // No elements to fetch
}
```

## Producer Task Member Function

```
void Produce(const T& t)
{
    list.push_back(t);           // Add the new element
    iTail = list.end();          // Update iTail
    list.erase(list.begin(), iHead); // Erase the removed elements
}
```

The screenshot shows the Microsoft Visual Studio IDE. The top menu bar includes File, Edit, View, Project, Build, Debug, Test, Analyze, Tools, Extensions, Window, Help, and Search (Ctrl+Q). The toolbar has icons for New, Open, Save, Print, and others. The status bar at the bottom shows "110%" and "No issues found".

**Source.cpp:**

```
10, 10, 10,
Press any key to continue . . .

86     // Then the arguments to the member
87     // so we need to wrap them in std::
88     /*
89     lfq.Produce(i);
90     lfq.Consume(j);
91     */
92
93     std::thread produce(&LockFreeQueue::produce, lfq, i);
94     threads.push_back(std::move(produce));
95     std::thread consume(&LockFreeQueue::consume, lfq, j);
96     threads.push_back(std::move(consume));
97 }
98
99 // Wait for the threads to complete
100 for (auto& thr: threads)
101     thr.join();
102
103 lfq.Print();
104 }
```

**Output:**

```
10, 10, 10,
Press any key to continue . . .
```

## Thread Safety

- The code has a data race
  - iHead and iTail can be accessed from different threads
  - At least one thread modifies them
  - The threads are not synchronized when they access these variables

## Avoiding the Data Race

- iHead and iTail cannot be atomic types
  - std::list<T>::iterator is not trivially copyable
- We must use mutexes

# Asynchronous Programming

## Synchronous and Asynchronous Programming

- Synchronous
  - Wait for each task to complete
- Asynchronous
  - Continue without waiting for tasks to complete

## Synchronous Programming

- A task starts another task
- The current task is blocked
- Must wait until the new task completes before it can continue
  - e.g. Synchronous database access
    - Do some work
    - Request data from database
    - Wait for data
    - Receive data from database
    - Continue working

## Synchronous Tasks

- Normal function calls are synchronous tasks

```
data.save(filename);
    // Stop and wait for func to return
    // ... wait ...

    // Now we can continue with the next operation
```

- We have to stop and wait for the save operation to complete
  - Even if the next operation does not depend on it
- This reduces throughput and user satisfaction
  - GUI applications appear unresponsive
  - Clients experience slow service

## Asynchronous Programming

- A task starts another task
- The current task can continue
- The new task runs in the background
  - e.g. Asynchronous database access
  - Request data from database as a separate task
  - Do some more work in our task
  - Receive data from database

## Asynchronous Tasks

- Start off another task

```
↳ data.async_save(filename);
    // The asynchronous task runs in the background
    // We continue with the next operation
    // ... do something else ...
```

- Our thread can continue its work

```
// At some point, we may need to check if the async call has completed
// Or to get its result
```

## Advantages of Asynchronous Programming

- The current task can do other work
  - Provided it does not require the data
- The current task only blocks when it needs the data
  - If the data is already available, it can continue without stopping
- This maintains throughput and user satisfaction
  - GUI applications appear responsive
  - Clients experience normal service

## Blocking and Multi-threading Programs

- Blocking is undesirable in threaded programs
  - Blocking reduces throughput and responsiveness of the blocked thread
  - Any threads which join with this thread will also be blocked
- Particularly in a critical section
  - Any threads which are waiting to enter the critical section are also blocked
  - Possibility of deadlock, if we are using locks
- Using asynchronous programming reduces the need to block
  - But may not avoid it completely
  - e.g. if the database fetch is not complete when the data is needed

## Blocking Synchronization

- Blocking operations
- Synchronized by mutexes
  - A thread is blocked until another thread unlocks the mutex
- Or atomic operations
  - A thread is blocked until another thread completes an atomic operation

## Non-blocking Synchronization

- Non-blocking operations
- Synchronized by message queues
  - A thread pushes a message onto a concurrent queue
  - Another thread takes the message off the queue and processes it
  - The first thread continues running without waiting
- The messages often consist of callable objects
  - The message is processed by invoking the callable object
- C++ does not have a standard concurrent queue
  - Available in Boost, Microsoft's PPL, Intel's TBB

## Asynchronous Programming and Parallelism

- Asynchronous programming
- Can be used to perform parallel operations
  - Start new threads which all perform the same task
  - Collect the result from each thread as it completes its task
  - Combine the results into the final answer
- It can also be used in single-threaded programs
  - Using operating system features