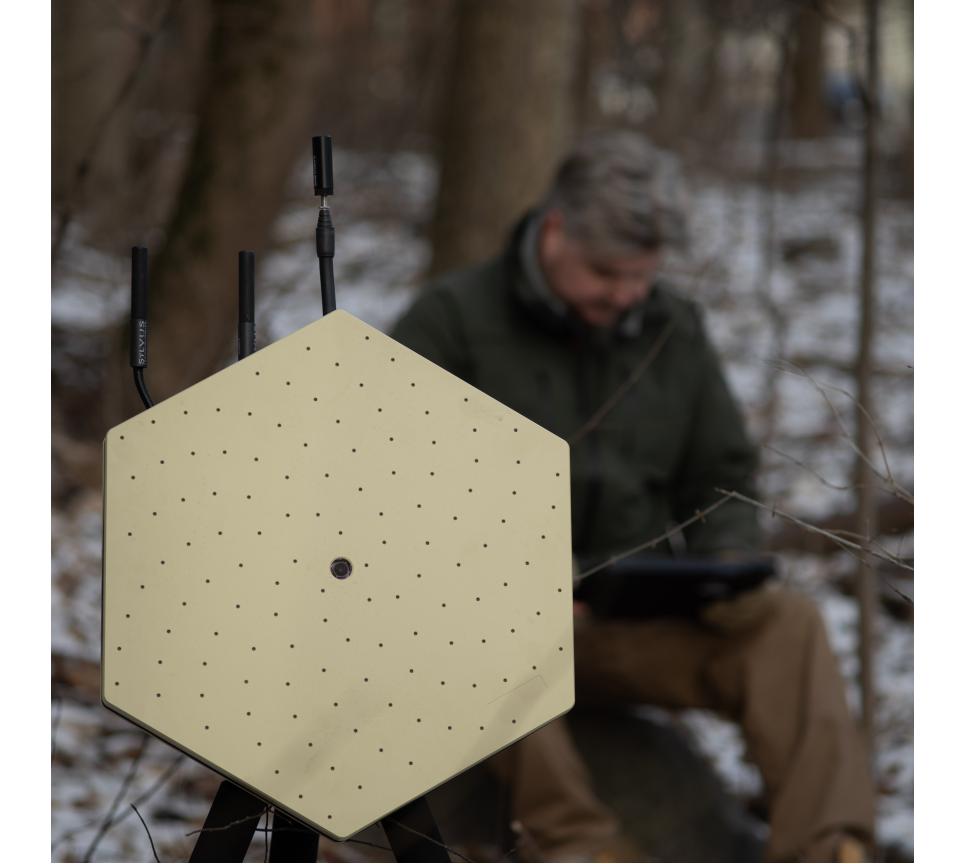


THE TWO MEMORY MODELS ANDERS SCHAU KNATTEN

C++ UNDER THE SEA 2025

Squarehead / CppQuiz.org







CPPQUIZ.ORG

C++ Quiz

You've answered 0 of 178 questions correctly. (Clear)

Question #197 Difficulty:

According to the C++23 standard, what is the output of this program?

```
#include <iostream>
int j = 1;
int main() {
  int& i = j, j;
  j = 2;
  std::cout << i << j;
}</pre>
```

Answer: The program is guaranteed to output:

Problems? View a hint or try another question.

I give up, show me the answer (make 3 more attempts first).

Mode: Training

You are currently in training mode, answering random questions. Why not Start a new quiz? Then you can boast about your score, and invite your friends.

Contribute

Create your own!

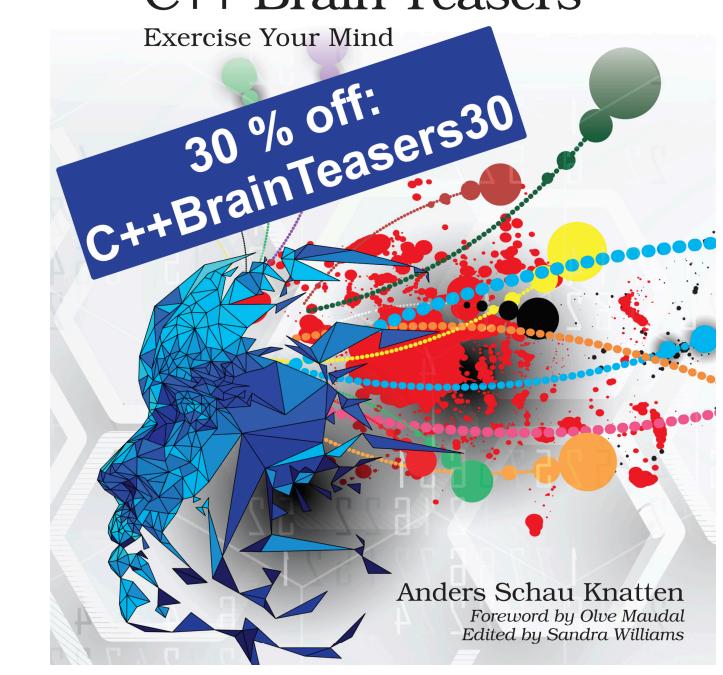
Android app

Get Sergey Vasilchenko's <u>CppQuiz Android</u> <u>app</u>.





C++ Brain Teasers







• Why?



- Why?
- Theoretical background



- Why?
- Theoretical background
- Intuition, framework for reasoning



- Why?
- Theoretical background
- Intuition, framework for reasoning
- How stuff works!





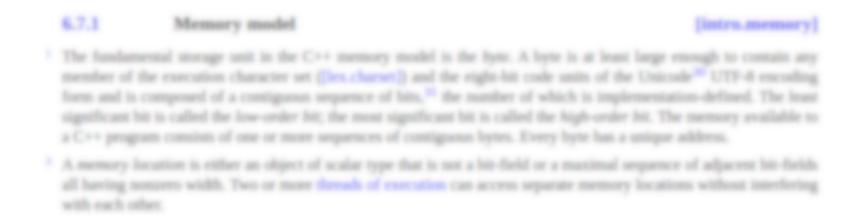
• High Level language (C++, C, Rust)



- High Level language (C++, C, Rust)
- CPU / Architecture











What's meant by "a memory location"





- What's meant by "a memory location"
- Two or more threads can access separate memory locations without interfering with each other.



• But other rules too:



- But other rules too:
 - Sequential execution [intro.execution]



- But other rules too:
 - Sequential execution [intro.execution]
 - Multi-threaded executions and data races [intro.multithread]



- But other rules too:
 - Sequential execution [intro.execution]
 - Multi-threaded executions and data races [intro.multithread]
 - Atomic operations library [atomics]



- But other rules too:
 - Sequential execution [intro.execution]
 - Multi-threaded executions and data races [intro.multithread]
 - Atomic operations library [atomics]
 - Thread support library [thread]





• X86, ARM, RISC-V, ...



- X86, ARM, RISC-V, ...
- Described in the architecture manual



- X86, ARM, RISC-V, ...
- Described in the architecture manual
- "A memory consistency model is a set of rules specifying the values that can be returned by loads of memory." - RISC-V ISA manual





• Simple if you have one thread



- Simple if you have one thread
 - A load from an address can't observe a later store to the same



- Simple if you have one thread
 - A load from an address can't observe a later store to the same
 - Program + input state = output state



- Simple if you have one thread
 - A load from an address can't observe a later store to the same
 - Program + input state = output state
- Complicated if you have many threads



- Simple if you have one thread
 - A load from an address can't observe a later store to the same
 - Program + input state = output state
- Complicated if you have many threads
 - Any number of valid interleavings of threads



- Simple if you have one thread
 - A load from an address can't observe a later store to the same
 - Program + input state = output state
- Complicated if you have many threads
 - Any number of valid interleavings of threads
 - Program + input state = many possible output states



WHAT HAS THE MEMORY MODEL EVER DONE FOR ME?



WHAT HAS THE MEMORY MODEL EVER DONE FOR ME?

- C++ memory model
 - Single threaded examples
 - Multi threaded example
 - Memory ordering



WHAT HAS THE MEMORY MODEL EVER DONE FOR ME?

- C++ memory model
 - Single threaded examples
 - Multi threaded example
 - Memory ordering
- CPU memory model:
 - Several variations, including x86 and RISC-V





```
1 int compute(const int* a, const int* b)
2 {
3    int result = *a;
4    result += 1;
5    result += *b;
6
7    return result;
8 }
```



```
int compute(const int* a, const int* b)

int result = *a;
result += 1;
result += *b;

return result;
}
```



```
int compute(const int* a, const int* b)

int result = *a;
result += 1;
result += *b;

return result;
}
```



```
int compute(const int* a, const int* b)

int result = *a;
result += 1;
result += *b;

return result;

}
```



```
int compute(const int* a, const int* b)

int result = *a;
result += 1;
result += *b;

return result;
}
```

















• Every value computation and side effect associated with a fullexpression is sequenced before every value computation and side effect associated with the next full-expression to be evaluated.



- Every value computation and side effect associated with a fullexpression is sequenced before every value computation and side effect associated with the next full-expression to be evaluated.
- Standard does not dictate how to implement C++!



- Every value computation and side effect associated with a fullexpression is sequenced before every value computation and side effect associated with the next full-expression to be evaluated.
- Standard does not dictate how to implement C++!
- Abstract machine



- Every value computation and side effect associated with a fullexpression is sequenced before every value computation and side effect associated with the next full-expression to be evaluated.
- Standard does not dictate how to implement C++!
- Abstract machine
- Only need to emulate observable behaviour



- Every value computation and side effect associated with a fullexpression is sequenced before every value computation and side effect associated with the next full-expression to be evaluated.
- Standard does not dictate how to implement C++!
- Abstract machine
- Only need to emulate observable behaviour
- Observable: I/O, files, volatile, plus more implementation defined



- Every value computation and side effect associated with a fullexpression is sequenced before every value computation and side effect associated with the next full-expression to be evaluated.
- Standard does not dictate how to implement C++!
- Abstract machine
- Only need to emulate observable behaviour
- Observable: I/O, files, volatile, plus more implementation defined
- Can reorder





```
// original
load reg_0, *a
load reg_0, reg_0, 1
load reg_1, *b
load reg_1, *b
add reg_0, reg_1
ret reg_0
// reordered
load reg_0, *a
load reg_1, *b
add reg_0, reg_0, 1
add reg_0, reg_0, reg_1
ret reg_0
```



```
// original
load reg_0, *a
load reg_0, reg_0, 1
load reg_1, *b
load reg_0, reg_0, reg_1
ret reg_0
// reordered
load reg_0, *a
load reg_1, *b
add reg_0, reg_0, 1
add reg_0, reg_0, reg_1
ret reg_0
```

Or CPU reorders it on the fly



```
// original
load reg_0, *a
load reg_0, *a
load reg_1, *b
load reg_1, *b
add reg_0, reg_1
ret reg_0
// reordered
load reg_0, *a
load reg_1, *b
add reg_0, reg_0, 1
add reg_0, reg_0, reg_1
ret reg_0
```

- Or CPU reorders it on the fly
- "All" modern CPUs are out-of-order



```
// original
load reg_0, *a
load reg_0, *eg_0, 1
load reg_1, *b
load reg_0, reg_0, reg_1
add reg_0, reg_0, reg_1
ret reg_0
// reordered
load reg_0, *a
load reg_0, *eg_1, *b
add reg_0, reg_0, 1
add reg_0, reg_0, reg_1
ret reg_0
```

- Or CPU reorders it on the fly
- "All" modern CPUs are out-of-order
- Can we swap the order of the loads too?



```
// original
load reg_0, *a
load reg_0, *a
load reg_1, *b
load reg_1, *b
add reg_0, reg_1
ret reg_0
// reordered
load reg_0, *a
load reg_1, *b
add reg_0, reg_0, 1
add reg_0, reg_0, reg_1
ret reg_0
```

- Or CPU reorders it on the fly
- "All" modern CPUs are out-of-order
- Can we swap the order of the loads too?
- Yes, can't affect the observable behaviour



SINGLE THREADED EXAMPLE (X86)

GCC 15.1 -03

```
int compute(const int* a, const int* b)
{
   int result = *a;
   result += 1;
   result += *b;
   return result;
}
```

```
1 compute(int const*, int const*):
2 mov    edx, DWORD PTR [rdi]
3 mov    eax, DWORD PTR [rsi]
4 lea    eax, [rdx+1+rax]
5 ret
```



SINGLE THREADED EXAMPLE (X86)

GCC 15.1 -03

```
int compute(const int* a, const int* b)
{
   int result = *a;
   result += 1;
   result += *b;
   return result;
}
```

```
1 compute(int const*, int const*):
2 mov    edx, DWORD PTR [rdi]
3 mov    eax, DWORD PTR [rsi]
4 lea    eax, [rdx+1+rax]
5 ret
```



SINGLE THREADED EXAMPLE (X86)

GCC 15.1 -03

```
int compute(const int* a, const int* b)
{
   int result = *a;
   result += 1;
   result += *b;
   return result;
}
```

```
1 compute(int const*, int const*):
2 mov    edx, DWORD PTR [rdi]
3 mov    eax, DWORD PTR [rsi]
4 lea    eax, [rdx+1+rax]
5 ret
```



```
1 void compute(int* a, int* b)
2 {
3     int result = get_value();
4     result++;
5     *a = 2;
6     *b = result;
7 }
```



```
1 void compute(int* a, int* b)
2 {
3     int result = get_value();
4     result++;
5     *a = 2;
6     *b = result;
7 }
```



```
1 void compute(int* a, int* b)
2 {
3     int result = get_value();
4     result++;
5     *a = 2;
6     *b = result;
7 }
```



```
1 void compute(int* a, int* b)
2 {
3     int result = get_value();
4     result++;
5     *a = 2;
6     *b = result;
7 }
```



```
1 void compute(int* a, int* b)
2 {
3     int result = get_value();
4     result++;
5     *a = 2;
6     *b = result;
7 }
```







Can we start the store to a earlier?



```
// original
move reg_0, get_value()
add reg_0, reg_0, 1
store *a, 2
store *b, reg_0
```



```
// original
move reg_0, get_value()
add reg_0, reg_0, 1
store *a, 2
store *b, reg_0
```



```
// original
move reg_0, get_value()
add reg_0, reg_0, 1
store *a, 2
store *b, reg_0
```

```
// reordered
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```



```
// original
move reg_0, get_value()
add reg_0, reg_0, 1
store *a, 2
store *b, reg_0
```

```
// reordered
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```



```
// original
move reg_0, get_value()
add reg_0, reg_0, 1
store *a, 2
store *b, reg_0
```

```
// reordered
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```



```
// original
move reg_0, get_value()
add reg_0, reg_0, 1
store *a, 2
store *b, reg_0
```

```
// reordered
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```



```
// original
move reg_0, get_value()
add reg_0, reg_0, 1
store *a, 2
store *b, reg_0
```

```
// reordered
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```

Or CPU reorders it



```
// original
move reg_0, get_value()
add reg_0, reg_0, 1
store *a, 2
store *b, reg_0
```

```
// reordered
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```

Or CPU reorders it

CPU could reorder the stores!



```
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```



```
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```

C++ to CPU: Do not store a before calling get_value!



```
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```

- C++ to CPU: Do not store a before calling get_value!
- C++ to CPU: Do not reorder stores!



```
move reg_0, get_value()
store *a, 2
add reg_0, reg_0, 1
store *b, reg_0
```

- C++ to CPU: Do not store a before calling get_value!
- C++ to CPU: Do not reorder stores!
- How do we express this to the CPU?



```
void compute(int* a, int* b)
{
    int result = get_value();
    result++;
    *a = 2;
    *b = result;
}
```

```
1 compute(int*, int*):
 2 push
          rbp
          rbp, rdi
 3 mov
          rbx
 4 push
          rbx, rsi
 5 mov
 6 sub rsp, 8
7 call get_value()
          DWORD PTR [rbp+0], 2
 8 mov
          eax, 1
 9 add
          DWORD PTR [rbx], eax
10 mov
          rsp, 8
11 add
12 pop
          rbx
13 pop
          rbp
14 ret
```



```
void compute(int* a, int* b)
{
    int result = get_value();
    result++;
    *a = 2;
    *b = result;
}
```

```
1 compute(int*, int*):
 2 push
          rbp
          rbp, rdi
 3 mov
 4 push
          rbx
          rbx, rsi
 5 mov
6 sub rsp, 8
7 call get_value()
          DWORD PTR [rbp+0], 2
 8 mov
          eax, 1
9 add
          DWORD PTR [rbx], eax
10 mov
11 add
          rsp, 8
12 pop
          rbx
13 pop
          rbp
14 ret
```



GCC 15.1 -03

```
void compute(int* a, int* b)
{
   int result = get_value();
   result++;
   *a = 2;
   *b = result;
}
```

```
1 compute(int*, int*):
2 push
          rbp
          rbp, rdi
3 mov
          rbx
4 push
          rbx, rsi
5 mov
6 sub rsp, 8
7 call get_value()
          DWORD PTR [rbp+0], 2
8 mov
          eax, 1
9 add
          DWORD PTR [rbx], eax
10 mov
          rsp, 8
11 add
12 pop
          rbx
13 pop
          rbp
14 ret
```

• X86 never reorders stores



```
void compute(int* a, int* b)
{
   int result = get_value();
   result++;
   *a = 2;
   *b = result;
}
```

```
1 compute(int*, int*):
 2 push
          rbp
          rbp, rdi
3 mov
          rbx
4 push
          rbx, rsi
5 mov
6 sub rsp, 8
7 call get_value()
          DWORD PTR [rbp+0], 2
8 mov
          eax, 1
9 add
          DWORD PTR [rbx], eax
10 mov
          rsp, 8
11 add
12 pop
          rbx
13 pop
          rbp
14 ret
```

- X86 never reorders stores
- RISC-V can reorder stores



```
void compute(int* a, int* b)
{
    int result = get_value();
    result++;
    *a = 2;
    *b = result;
}
```

```
1 compute(int*, int*):
2 push
          rbp
3 mov
          rbp, rdi
          rbx
4 push
          rbx, rsi
5 mov
6 sub rsp, 8
7 call get_value()
          DWORD PTR [rbp+0], 2
8 mov
          eax, 1
9 add
          DWORD PTR [rbx], eax
10 mov
          rsp, 8
11 add
12 pop
          rbx
13 pop
          rbp
14 ret
```

- X86 never reorders stores
- RISC-V can reorder stores to different addresses



```
void compute(int* a, int* b)
{
   int result = get_value();
   result++;
   *a = 2;
   *b = result;
}
```

```
1 compute(int*, int*):
2 push
          rbp
3 mov
        rbp, rdi
4 push
          rbx
          rbx, rsi
5 mov
6 sub rsp, 8
7 call get_value()
          DWORD PTR [rbp+0], 2
8 mov
          eax, 1
9 add
          DWORD PTR [rbx], eax
10 mov
11 add rsp, 8
12 pop
          rbx
13 pop
          rbp
14 ret
```

- X86 never reorders stores
- RISC-V can reorder stores to different addresses
- Reorder includes effects from pipeline, caches, store buffers etc







```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```

```
1 producer():
2    sub    rsp, 8
3    call initializeData()
4    mov BYTE PTR ready[rip], 1
5    add rsp, 8
6    ret
7    consumer():
8    jmp    useData()
```



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```

```
producer():
    sub rsp, 8
    call initializeData()
    mov BYTE PTR ready[rip], 1
    add rsp, 8
    ret
    consumer():
    jmp useData()
```



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```

```
1 producer():
2    sub    rsp, 8
3    call initializeData()
4    mov BYTE PTR ready[rip], 1
5    add rsp, 8
6    ret
7    consumer():
8    jmp    useData()
```



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```

```
1 producer():
2    sub    rsp, 8
3    call initializeData()
4    mov BYTE PTR ready[rip], 1
5    add rsp, 8
6    ret
7 consumer():
8    jmp    useData()
```

Informally:



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```

```
producer():
    sub rsp, 8
    call initializeData()
    mov BYTE PTR ready[rip], 1
    add rsp, 8
    ret
    consumer():
    jmp useData()
```

Informally: Always true



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```

```
1 producer():
2    sub    rsp, 8
3    call initializeData()
4    mov BYTE PTR ready[rip], 1
5    add rsp, 8
6    ret
7    consumer():
8    jmp    useData()
```

Informally: Always true / always false (UB)



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```

Informally: Always true / always false (UB) / data race (UB)



```
1 Data data;
2 bool ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```

```
producer():
    sub rsp, 8
    call initializeData()
    mov BYTE PTR ready[rip], 1
    add rsp, 8
    ret
    consumer():
    jmp useData()
```

- Informally: Always true / always false (UB) / data race (UB)
- Forward progress: Terminate, I/O, volatile, or sync/atomic



MULTI THREADED EXAMPLE (NON UB)

```
1 Data data;
2 std::atomic<bool> ready{false};
3
4 void producer()
5 {
6     initializeData();
7     ready = true;
8 }
9
10 void consumer()
11 {
12     while(!ready){}
13     useData();
14 }
```



MULTI THREADED EXAMPLE (NON UB)

```
1 Data data;
2 std::atomic<bool> ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready = true;
8 }
9
10 void consumer()
11 {
12    while(!ready){}
13    useData();
14 }
```



MULTI THREADED EXAMPLE (NON UB X86)

```
1 Data data;
2 std::atomic<bool> ready{false};
3
4 void producer()
5 {
6     initializeData();
7     ready = true;
8 }
9
10 void consumer()
11 {
12     while(!ready){}
13     useData();
14 }
```

```
1 producer():
    sub
           rsp, 8
   call
           initializeData()
  mov
           eax, 1
  xchg
           al, BYTE PTR ready[rip]
    add
           rsp, 8
    ret
  consumer():
9 .L5:
           eax, BYTE PTR ready[rip]
   movzx
           al, al
  test
   jе
            . L5
           useData()
   jmp
```



MULTI THREADED EXAMPLE (NON UB X86)

```
1 Data data;
2 std::atomic<bool> ready{false};
3
4 void producer()
5 {
6     initializeData();
7     ready = true;
8 }
9
10 void consumer()
11 {
12     while(!ready){}
13     useData();
14 }
```

```
1 producer():
    sub
           rsp, 8
   call
           initializeData()
  mov
           eax, 1
  xchg
           al, BYTE PTR ready[rip]
    add
           rsp, 8
   ret
 consumer():
9 .L5:
           eax, BYTE PTR ready[rip]
   movzx
  test
           al, al
           . L5
   je
   jmp
           useData()
```



MULTI THREADED EXAMPLE (NON UB X86)

```
1 Data data;
2 std::atomic<bool> ready{false};
3
4 void producer()
5 {
6     initializeData();
7     ready = true;
8 }
9
10 void consumer()
11 {
12     while(!ready){}
13     useData();
14 }
```

```
1 producer():
    sub
           rsp, 8
   call
           initializeData()
           eax, 1
  mov
           al, BYTE PTR ready[rip]
  xchg
    add
          rsp, 8
  ret
8 consumer():
9 .L5:
           eax, BYTE PTR ready[rip]
   movzx
           al, al
  test
   jе
           . L5
           useData()
   jmp
```



MULTI THREADED EXAMPLE (NON UB X86)

```
1 Data data;
2 std::atomic<bool> ready{false};
3
4 void producer()
5 {
6     initializeData();
7     ready = true;
8 }
9
10 void consumer()
11 {
12     while(!ready){}
13     useData();
14 }
```

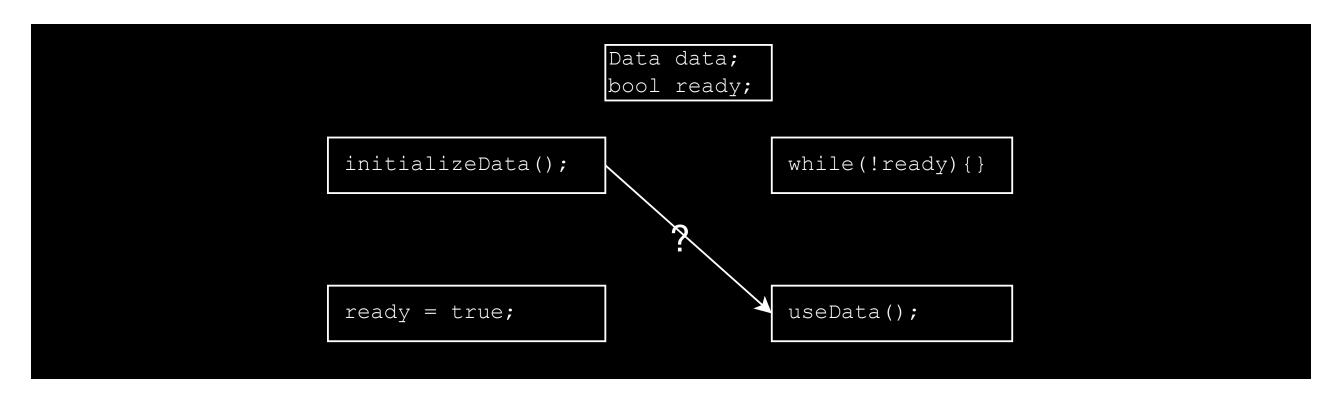
```
1 producer():
    sub
           rsp, 8
   call
           initializeData()
           eax, 1
  mov
           al, BYTE PTR ready[rip]
  xchg
    add
           rsp, 8
  ret
 consumer():
9 .L5:
           eax, BYTE PTR ready[rip]
   movzx
  test
           al, al
           . L5
   je
           useData()
   jmp
```



- Is the data initialized?
- Is the cache up to date?
- Is there data race?
- Microarchitectural details we don't know about?



- Is the data initialized?
- Is the cache up to date?
- Is there data race?
- Microarchitectural details we don't know about?







• Two expression evaluations conflict if **one of them modifies** a memory location and the **other one reads or modifies** the same memory location



• Two expression evaluations conflict if **one of them modifies** a memory location and the **other one reads or modifies** the same memory location (This is fine)



- Two expression evaluations conflict if **one of them modifies** a memory location and the **other one reads or modifies** the same memory location (This is fine)
- The execution of a program contains a data race if it contains two potentially concurrent conflicting actions,



- Two expression evaluations conflict if **one of them modifies** a memory location and the **other one reads or modifies** the same memory location (This is fine)
- The execution of a program contains a data race if it contains two potentially concurrent conflicting actions, at least one of which is not atomic,



- Two expression evaluations conflict if **one of them modifies** a memory location and the **other one reads or modifies** the same memory location (This is fine)
- The execution of a program contains a data race if it contains two potentially concurrent conflicting actions, at least one of which is not atomic, and **neither happens before the other**.

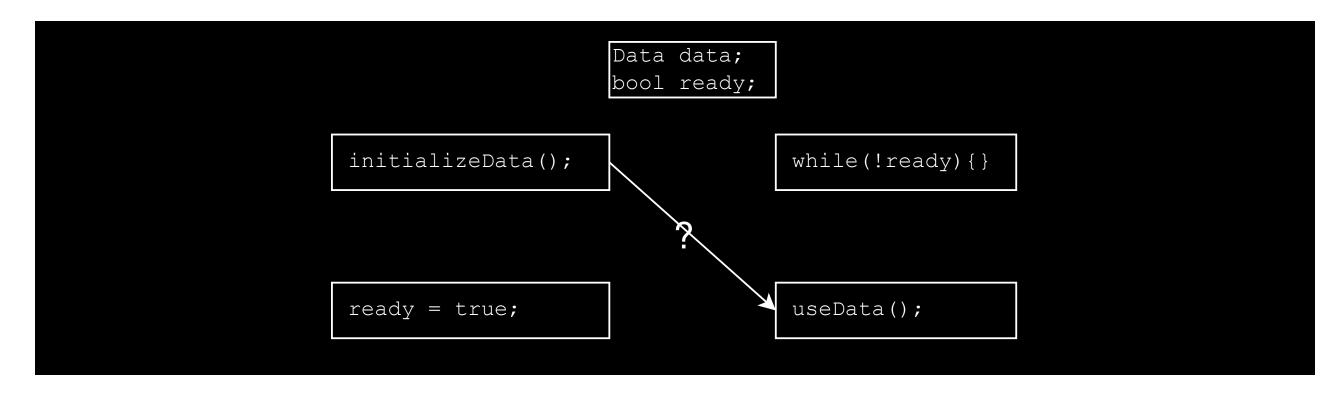


- Two expression evaluations conflict if **one of them modifies** a memory location and the **other one reads or modifies** the same memory location (This is fine)
- The execution of a program contains a data race if it contains two potentially concurrent conflicting actions, at least one of which is not atomic, and **neither happens before the other**. Any such data race results in undefined behavior.



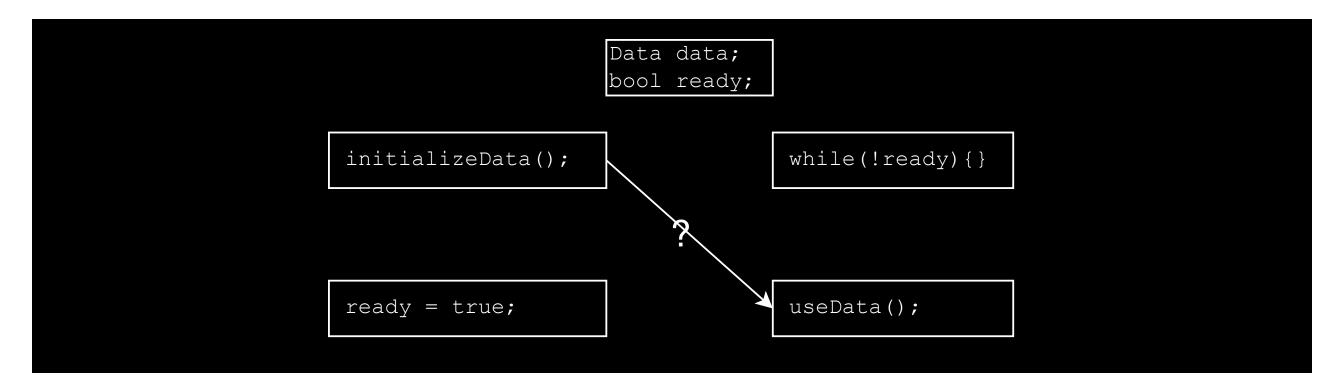
- Two expression evaluations conflict if **one of them modifies** a memory location and the **other one reads or modifies** the same memory location (This is fine)
- The execution of a program contains a data race if it contains two potentially concurrent conflicting actions, at least one of which is not atomic, and **neither happens before the other**. Any such data race results in undefined behavior.
- Does the data write happen before the read?





Does initializeData happen before useData?

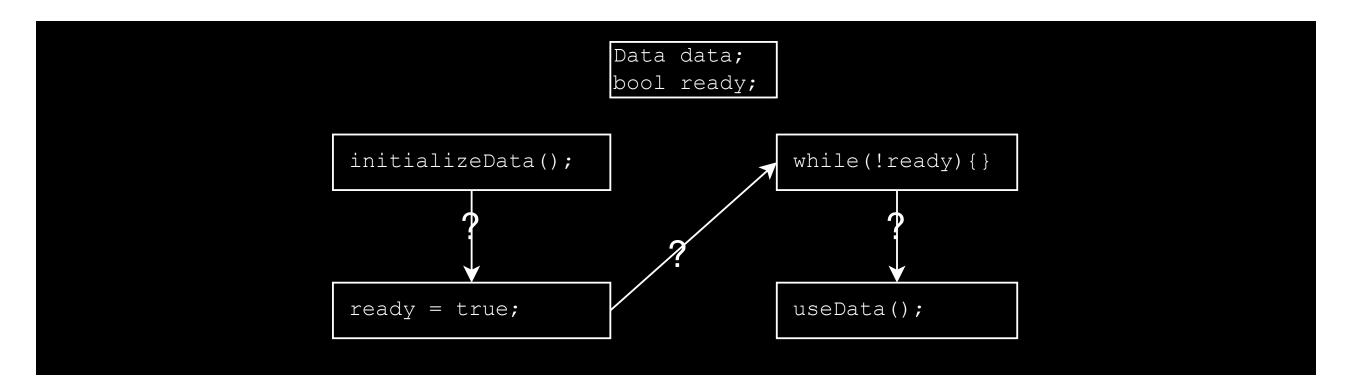




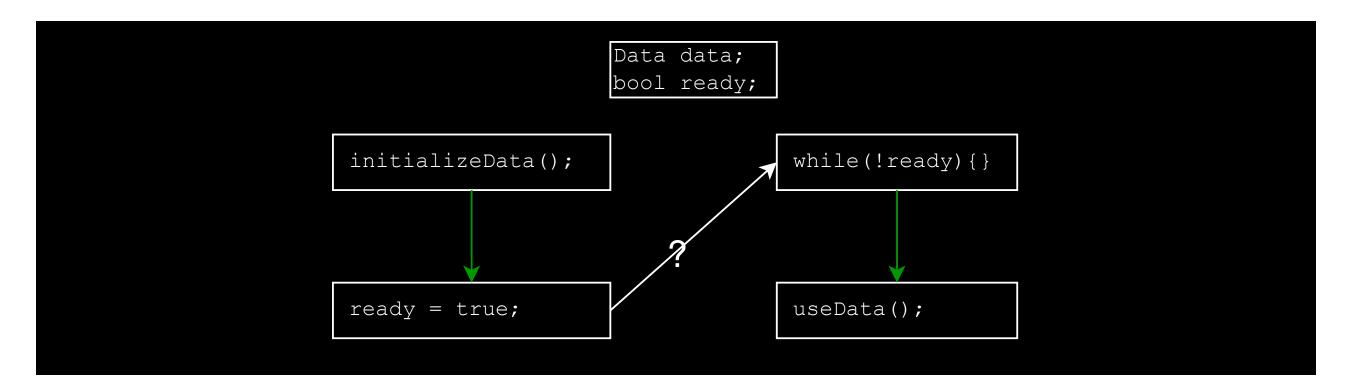
Does initializeData happen before useData?

Is there a happens-before relationship here?

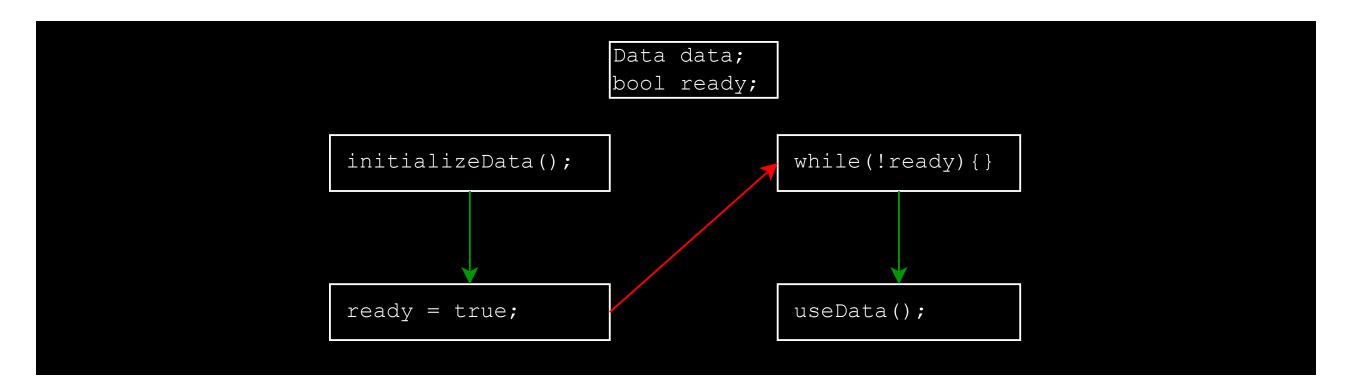




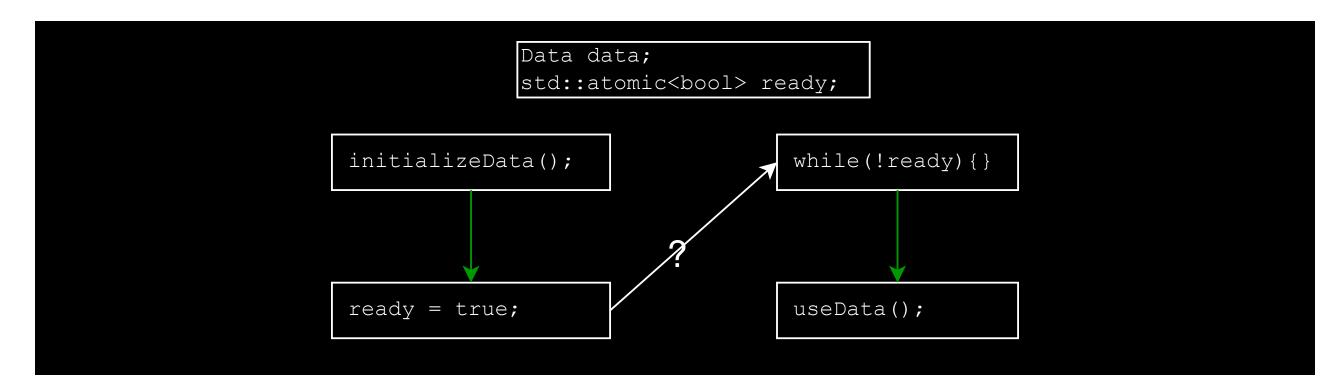




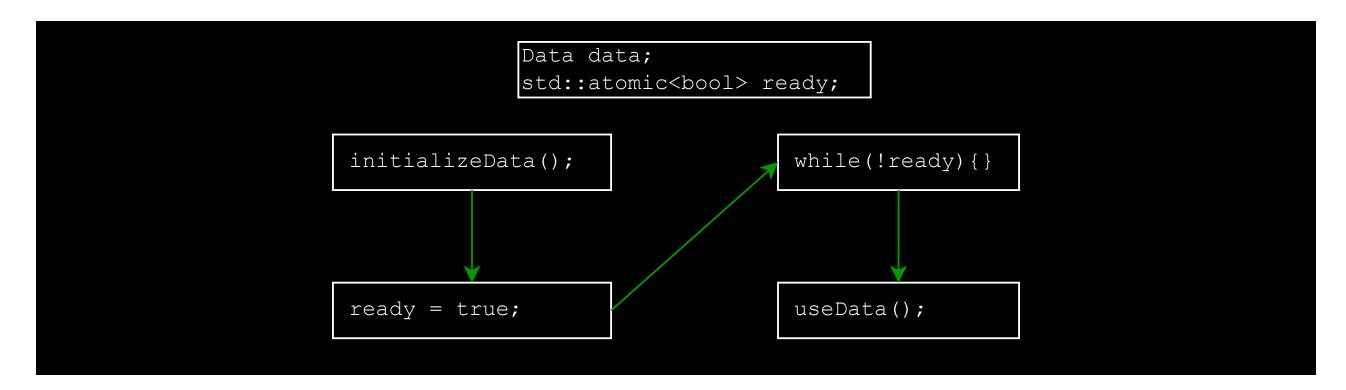
















• Atomics guarantee no half-written values



- Atomics guarantee no half-written values
- There is a total order to modifications of each individual atomic



- Atomics guarantee no half-written values
- There is a total order to modifications of each individual atomic
- Otherwise, it depends



```
Data data;
std::atomic<bool> ready{false};

void producer()
{
    initializeData();
    ready = true;
}

void consumer()
{
    while(!ready){}
    useData();
}
```



```
Data data;
std::atomic<bool> ready{false};

void producer()
{
    initializeData();
    ready.store(true, std::memory_order_seq_cst);
}

void consumer()
{
    while(!ready.load(std::memory_order_seq_cst)){}
    useData();
}
```



```
1 Data data;
2 std::atomic<bool> ready{false};
3
4 void producer()
5 {
6    initializeData();
7    ready.store(true, std::memory_order_seq_cst);
8 }
9
10 void consumer()
11 {
12    while(!ready.load(std::memory_order_seq_cst)){}
13    useData();
14 }
```



• Sequentially consistent:



- Sequentially consistent:
 - Establishes happens-before across threads



- Sequentially consistent:
 - Establishes happens-before across threads
 - Single total order for all seq_cst operations



- Sequentially consistent:
 - Establishes happens-before across threads
 - Single total order for all seq_cst operations
 - Stricter than necessary about reordering



- Sequentially consistent:
 - Establishes happens-before across threads
 - Single total order for all seq_cst operations
 - Stricter than necessary about reordering
- Relaxed:



- Sequentially consistent:
 - Establishes happens-before across threads
 - Single total order for all seq_cst operations
 - Stricter than necessary about reordering
- Relaxed:
 - No happens-before across threads



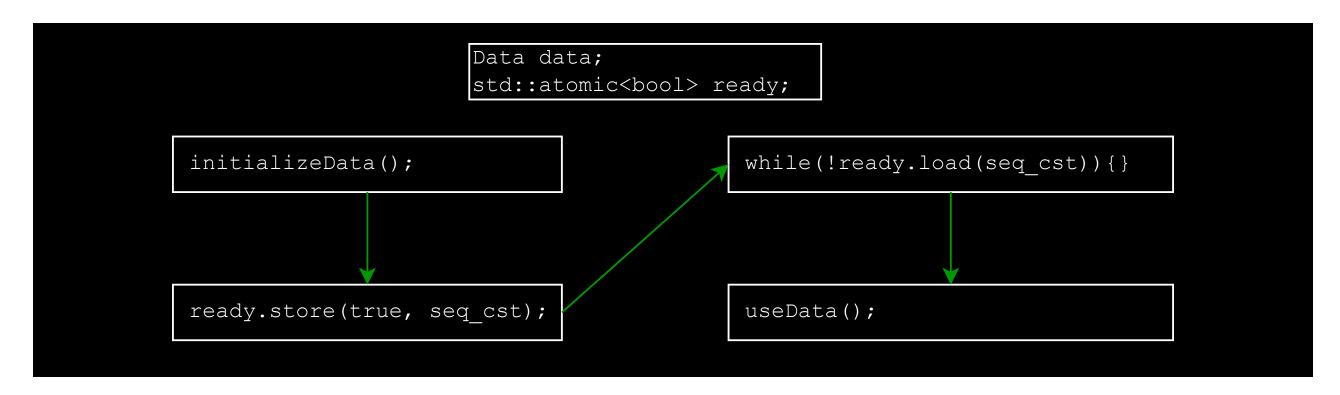
- Sequentially consistent:
 - Establishes happens-before across threads
 - Single total order for all seq_cst operations
 - Stricter than necessary about reordering
- Relaxed:
 - No happens-before across threads
 - One total modification order per atomic



- Sequentially consistent:
 - Establishes happens-before across threads
 - Single total order for all seq_cst operations
 - Stricter than necessary about reordering
- Relaxed:
 - No happens-before across threads
 - One total modification order per atomic
- Avoid half-written writes

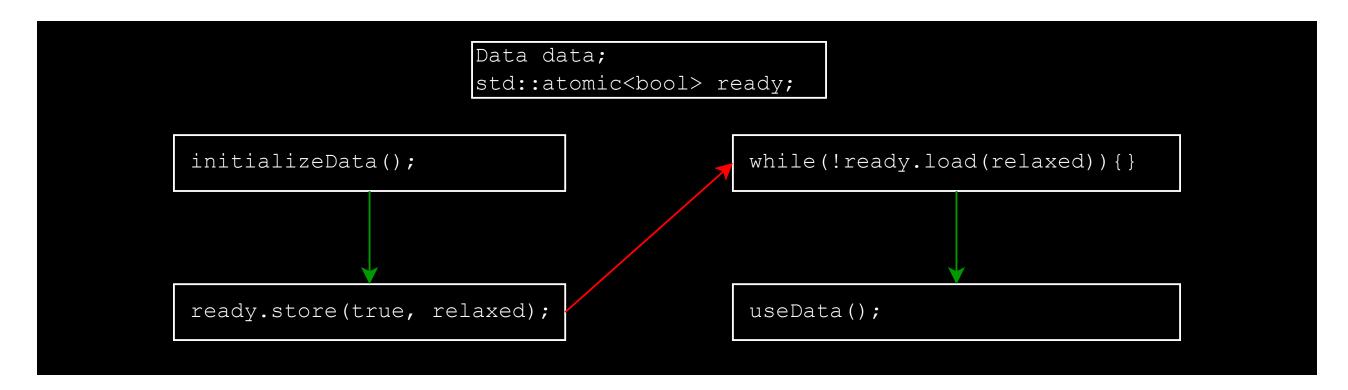


SEQUENTIALLY CONSISTENT





RELAXED

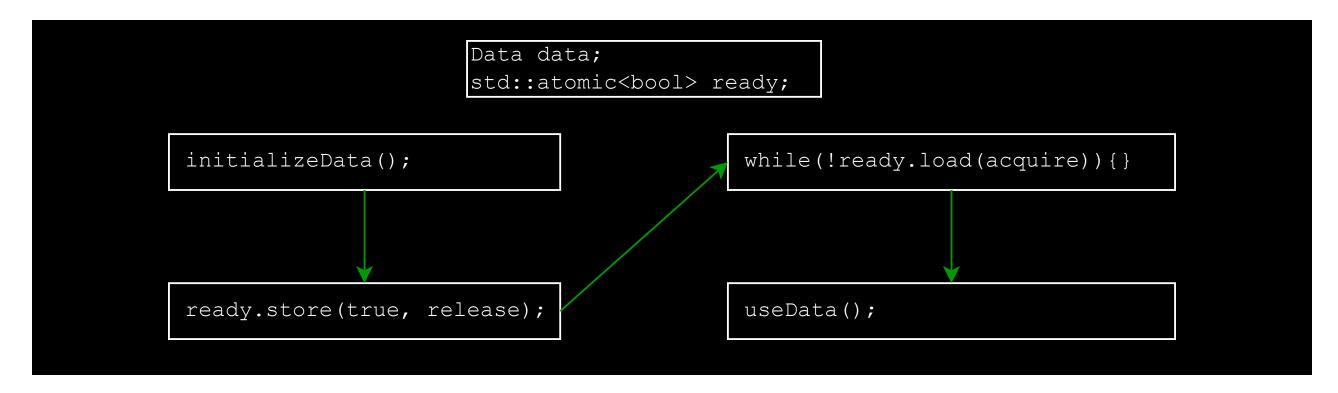




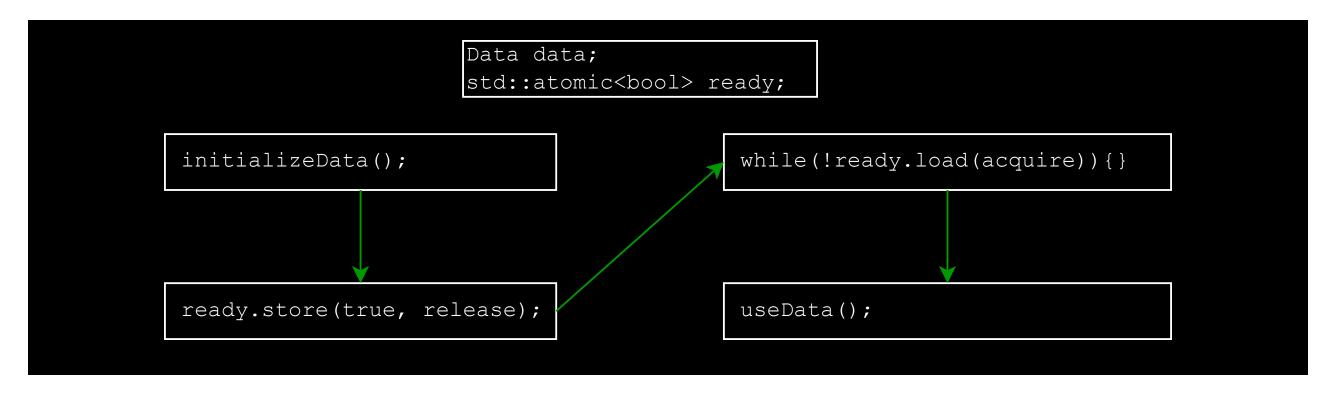
ACQUIRE / RELEASE

If an acquire-load observes the value of a release-store, the release-store store happens before the acquire-load



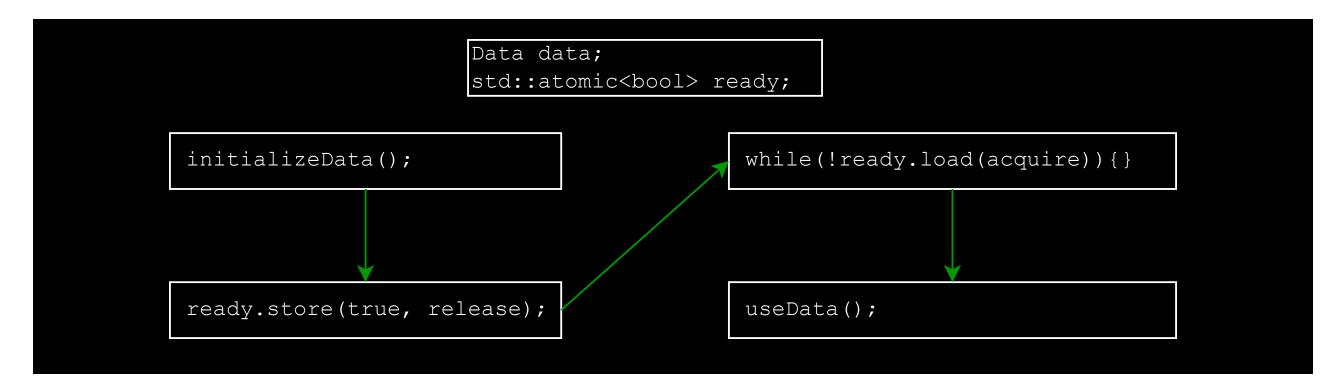






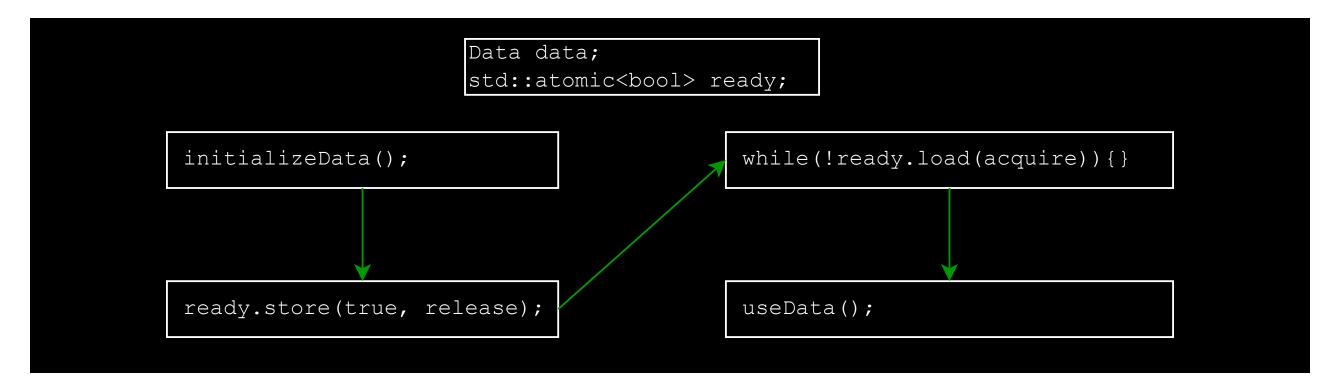
• Store happens-before load





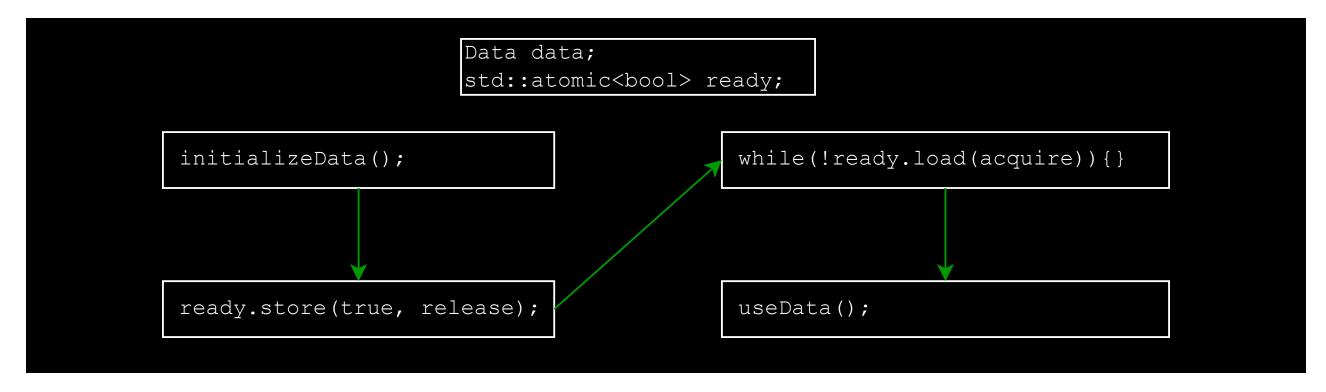
- Store happens-before load
- Can't reorder memory ops after the release-store





- Store happens-before load
- Can't reorder memory ops after the release-store
- Can't reorder memory ops before the acquire-load





- Store happens-before load
- Can't reorder memory ops after the release-store
- Can't reorder memory ops before the acquire-load
- But opposite is ok



MULTI THREADED EXAMPLE (X86)

```
Data data;
std::atomic<bool> ready{false};

void producer()
{
    initializeData();
    ready.store(true, std::memory_order_release)
}

void consumer()
{
    while(!ready.load(std::memory_order_acquire useData();
}
```



MULTI THREADED EXAMPLE (X86)

```
Data data;
std::atomic<bool> ready{false};

void producer()
{
    initializeData();
    ready.store(true, std::memory_order_release)
}

void consumer()
{
    while(!ready.load(std::memory_order_acquiruseData();
}
```

```
1 producer():
    sub
            rsp, 8
   call
           initializeData()
            BYTE PTR ready[rip], 1
   mov
            rsp, 8
    add
    ret
7 consumer():
8 .L5:
            eax, BYTE PTR ready[rip
   movzx
            al, al
   test
            . L5
   jе
            useData()
    jmp
```



MULTI THREADED EXAMPLE (X86)

```
Data data;
std::atomic<bool> ready{false};

void producer()
{
    initializeData();
    ready.store(true, std::memory_order_release)
}

void consumer()
{
    while(!ready.load(std::memory_order_acquiruseData();
}
```

```
1 producer():
           rsp, 8
    sub
   call
           initializeData()
           BYTE PTR ready[rip], 1
   mov
           rsp, 8
    add
   ret
7 consumer():
8 .L5:
           eax, BYTE PTR ready[rip
   movzx
           al, al
   test
           . L5
   jmp
           useData()
```



Also, memory_order::

- acq_rel (read-modify-write ops)
- consume (don't use)



How does the compiler know what is safe and not?



CPU MEMORY MODEL





Out of order instruction execution?



- Out of order instruction execution?
- Out of order micro-op execution?



- Out of order instruction execution?
- Out of order micro-op execution?
- Pipeline?



- Out of order instruction execution?
- Out of order micro-op execution?
- Pipeline?
- Forwarding networks?



- Out of order instruction execution?
- Out of order micro-op execution?
- Pipeline?
- Forwarding networks?
- Store buffers?



- Out of order instruction execution?
- Out of order micro-op execution?
- Pipeline?
- Forwarding networks?
- Store buffers?
- Store/load coalescing?



- Out of order instruction execution?
- Out of order micro-op execution?
- Pipeline?
- Forwarding networks?
- Store buffers?
- Store/load coalescing?
- Memory system?



- Out of order instruction execution?
- Out of order micro-op execution?
- Pipeline?
- Forwarding networks?
- Store buffers?
- Store/load coalescing?
- Memory system?
- Caches?





• Different way of thinking than C++ memory model!



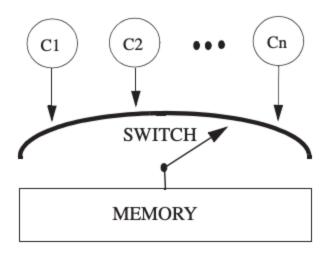
- Different way of thinking than C++ memory model!
- Imagine a single global order of all memory operations



- Different way of thinking than C++ memory model!
- Imagine a single global order of all memory operations
- What are the allowed orders



- Different way of thinking than C++ memory model!
- Imagine a single global order of all memory operations
- What are the allowed orders



Sarin, Hill, Wood: A Primer on Memory Consistency and Cache Coherence



Producer thread

mov reg_0 42

store *data reg_0

mov reg_1 1

store *ready reg_1



Producer thread

mov reg_0 42

store *data reg_0

mov reg_1 1

store *ready reg_1

Consumer thread

■wait

load reg_0 *ready

jz reg_0 .wait

load reg_1 *data

(use reg_1)



Producer thread

store *data reg_0

store *ready reg_1

Consumer thread

load reg_0 *ready

load reg_1 *data



Producer thread

Consumer thread

store *data reg_0

load reg_0 *ready

store *ready reg_1

load reg_1 *data

What is a valid order here?



```
Producer thread

Store *data reg_0

Store *ready reg_1

Consumer thread

load reg_0 *ready

load reg_1 *data
```

What is a valid order here?

• In the (imaginary) global, total memory order,



```
Producer thread

Store *data reg_0

Consumer thread

load reg_0 *ready
```

store *ready reg_1 load reg_1 *data

What is a valid order here?

- In the (imaginary) global, total memory order,
- imagine each thread takes turns (as-if rule).



```
Producer thread

Store *data reg_0

Store *ready reg_1

Consumer thread

load reg_0 *ready

load reg_1 *data
```

What is a valid order here?

- In the (imaginary) global, total memory order,
- imagine each thread takes turns (as-if rule).
- How much of program order is preserved?



MEMORY ORDER

Producer thread

store *data reg_0

store *ready reg_1

Memory order

store *data reg_0

store *ready reg_1

load reg_0 *ready

load reg_1 *data

Consumer thread

load reg_0 *ready

load reg_1 *data



MEMORY ORDER

Producer thread

store *data reg_0

store *ready reg_1

Memory order

store *data reg_0

load reg_0 *ready

load reg_0 *ready

store *ready reg_1

load reg_0 *ready

load reg_1 *data

Consumer thread

load reg_0 *ready

load reg_1 *data



Producer thread

store *data reg_0

store *ready reg_1

Memory order

load reg_0 *ready

load reg_0 *ready

store *ready reg_1

load reg_0 *ready

load reg_1 *data

store *data reg_0

Consumer thread

load reg_0 *ready



Producer thread

store *data reg_0

store *ready reg_1

Memory order

store *data reg_0

load reg_0 *ready

load reg_0 *ready

store *ready reg_1

load reg_0 *ready

load reg_1 *data

Sequential consistency

Consumer thread

load reg_0 *ready



Producer thread

store *data reg_0
store *ready reg_1

Memory order

store *data reg_0

load reg_0 *ready

load reg_0 *ready

store *ready reg_1

load reg_0 *ready

load reg_1 *data

- Sequential consistency
- All of program order is preserved

Consumer thread

load reg_0 *ready



Producer thread

store *data reg_0

store *ready reg_1

Memory order

store *data reg_0

load reg_0 *ready

load reg_0 *ready

store *ready reg_1

load reg_0 *ready

load reg_1 *data

- Sequential consistency
- All of program order is preserved
- Easy to reason about

Consumer thread

load reg_0 *ready



```
Producer thread
```

store *data reg_0
store *ready reg_1

```
Memory order
```

```
store *data reg_0
load reg_0 *ready
```

load reg_0 *ready

```
store *ready reg_1
```

load reg_0 *ready

```
load reg_1 *data
```

- Sequential consistency
- All of program order is preserved
- Easy to reason about
- Less possibilities for optimization ("no-one" does this)

Consumer thread

```
load reg_0 *ready
```



SEQUENTIAL CONSISTENCY

Preserves

- Load → Load
- Store → Store
- Load → Store
- Store → Load



Preserves

- Load → Load
- Store → Store
- Load → Store
- Store → Load



Motivated by store buffers



- Motivated by store buffers
- Stores are buffered on the core into a store buffer



- Motivated by store buffers
- Stores are buffered on the core into a store buffer
- A later load might take its value before the buffer is drained to cache



Thread 1 Thread 2

store x 1 store y 1

load y load x



Thread 1 Thread 2

store x 1 store y 1

load y load x

Can x and y be 0?



Thread 1 Memory order Thread 2

store x 1 load y store y 1

load y store y 1 load x

load x

store x 1



```
Thread 1 Memory order Thread 2

store x 1 load y store y 1

load y load x

store x 1

x == 0, y == 0
```

Fine on TSO, not on SC!



Thread 1

store x 1

load x



Thread 1

store x 1

load x

Can store x be ordered after load x?



Thread 1 Memory order

store x 1 load x

load x store x 1



Thread 1 Memory order

store x 1 load x

load x store x 1



Thread 1 Memory order store x 1 load x load x

Yes! But.

 Memory order: The order that memory operations go to a hypothetical serialized memory



Thread 1 Memory order store x 1 load x load x

- Memory order: The order that memory operations go to a hypothetical serialized memory
- Value of a load follows program order, otherwise memory order



Thread 1 Memory order store x 1 load x load x

- Memory order: The order that memory operations go to a hypothetical serialized memory
- Value of a load follows program order, otherwise memory order
- load x picks 1 from the core's store buffer



Thread 1 Memory order store x 1 load x load x

- Memory order: The order that memory operations go to a hypothetical serialized memory
- Value of a load follows program order, otherwise memory order
- load x picks 1 from the core's store buffer
- load x never even goes to cache



Thread 1 Memory order store x 1 load x load x

- Memory order: The order that memory operations go to a hypothetical serialized memory
- Value of a load follows program order, otherwise memory order
- load x picks 1 from the core's store buffer
- load x never even goes to cache
- Just a model



MULTI THREADED EXAMPLE (X86/TSO)

```
Data data;
std::atomic<bool> ready{false};

void producer()
{
    initializeData();
    ready.store(true, std::memory_order_release)
}

void consumer()
{
    while(!ready.load(std::memory_order_acquiruseData();
}
```

```
1 producer():
    sub
           rsp, 8
   call
           initializeData()
           BYTE PTR ready[rip], 1
   mov
            rsp, 8
    add
    ret
7 consumer():
8 .L5:
           eax, BYTE PTR ready[rip
   movzx
            al, al
   test
            . L5
   jе
           useData()
    jmp
```



MULTI THREADED EXAMPLE (X86/TSO)

```
Data data;
std::atomic<bool> ready{false};

void producer()
{
    initializeData();
    ready.store(true, std::memory_order_release)
}

void consumer()
{
    while(!ready.load(std::memory_order_acquiruseData();
}
```

```
1 producer():
    sub
           rsp, 8
   call
           initializeData()
           BYTE PTR ready[rip], 1
   mov
           rsp, 8
    add
   ret
7 consumer():
8 .L5:
           eax, BYTE PTR ready[rip
   movzx
           al, al
   test
           . L5
   jmp
           useData()
```



• Do we really need to keep order of all stores



• Do we really need to keep order of all stores and all loads



• Do we really need to keep order of all stores and all loads and all loads before stores?



- Do we really need to keep order of all stores and all loads and all loads before stores?
- What if want non-FIFO store buffers?



- Do we really need to keep order of all stores and all loads and all loads before stores?
- What if want non-FIFO store buffers?
- What if we want coalescing of stores/loads?



- Do we really need to keep order of all stores and all loads and all loads before stores?
- What if want non-FIFO store buffers?
- What if we want coalescing of stores/loads?
- What if we want fancy speculation and prediction?



- Do we really need to keep order of all stores and all loads and all loads before stores?
- What if want non-FIFO store buffers?
- What if we want coalescing of stores/loads?
- What if we want fancy speculation and prediction?
- What if we want other optimizations in the memory system / cache?



Producer thread

store *data0 reg_0

store *data1 reg_1

store *ready reg_2

Consumer thread

load reg_0 *ready

load reg_1 *data0



Producer thread

store *data0 reg_0

store *data1 reg_1

store *ready reg_2

Memory order

store *data1 reg_1

store *data0 reg_0

store *ready reg_2

load reg_0 *ready

load reg_2 *data1

load reg_1 *data0

Consumer thread

load reg_0 *ready

load reg_1 *data0



• As a programmer (with an interest in performance), I want to



- As a programmer (with an interest in performance), I want to
 - Give the architecture a lot of freedom to optimize



RELAXED MEMORY MODELS

- As a programmer (with an interest in performance), I want to
 - Give the architecture a lot of freedom to optimize
 - Tell it which memory operations I care about



RELAXED MEMORY MODELS

- As a programmer (with an interest in performance), I want to
 - Give the architecture a lot of freedom to optimize
 - Tell it which memory operations I care about
 - At least a few safe defaults



RELAXED MEMORY MODELS

- As a programmer (with an interest in performance), I want to
 - Give the architecture a lot of freedom to optimize
 - Tell it which memory operations I care about
 - At least a few safe defaults
 - A memory model to reason about



RELAXED MEMORY MODELS: RISC-V SAFE DEFAULTS: SAME ADDRESS STORE → STORE

Producer thread

Memory order?

store *data0 reg_0

store *data0 reg_1

store *data0 reg_1 store *data0 reg_0



RELAXED MEMORY MODELS: RISC-V SAFE DEFAULTS: DATA DEPENDENCY

Producer thread

load reg_0 *data0

store *data1 reg_0

Memory order?

store *data1 reg_0

load reg_0 *data0



RELAXED MEMORY MODELS: RISC-V SAFE DEFAULTS: ADDRESS DEPENDENCY

Producer thread

load reg_0 pointer

store *reg_0 2

Memory order?

store *reg_0 2

load reg_0 pointer



FENCES

Producer thread

store *data0 reg_0

store *data1 reg_1

store *ready reg_2

Memory order

store *data1 reg_1

store *data0 reg_0

store *ready reg_2

load reg_0 *ready

load reg_2 *data1

load reg_1 *data0

Consumer thread

load reg_0 *ready

load reg_1 *data0

load reg_2 *data1



FENCES

Producer thread

store *data0 reg_0

store *data1 reg_1

fence w, w

store *ready reg_2

Memory order

store *data1 reg_1

store *data0 reg_0

store *ready reg_2

load reg_0 *ready

load reg_2 *data1

load reg_1 *data0

Consumer thread

load reg_0 *ready

fence r, r

load reg_1 *data0

load reg_2 *data1





RISC-V allows all sorts of reordering



- RISC-V allows all sorts of reordering
- Except 13 specific rules, e.g.



- RISC-V allows all sorts of reordering
- Except 13 specific rules, e.g.
 - Don't reorder past a store to overlapping address



- RISC-V allows all sorts of reordering
- Except 13 specific rules, e.g.
 - Don't reorder past a store to overlapping address
 - Don't reorder across data/control/address dependency



- RISC-V allows all sorts of reordering
- Except 13 specific rules, e.g.
 - Don't reorder past a store to overlapping address
 - Don't reorder across data/control/address dependency
 - Don't reorder across fences (depending on type of fence)



- RISC-V allows all sorts of reordering
- Except 13 specific rules, e.g.
 - Don't reorder past a store to overlapping address
 - Don't reorder across data/control/address dependency
 - Don't reorder across fences (depending on type of fence)
 - Don't reorder across AMOs with acquire/release semantics



- RISC-V allows all sorts of reordering
- Except 13 specific rules, e.g.
 - Don't reorder past a store to overlapping address
 - Don't reorder across data/control/address dependency
 - Don't reorder across fences (depending on type of fence)
 - Don't reorder across AMOs with acquire/release semantics
 - And more, see "Preserved Program Order" in RISC-V ISA manual







1

Architecture memory model

1



1

Architecture memory model

1

Microarchitecture (pipeline, store buffers)



1

Architecture memory model

1

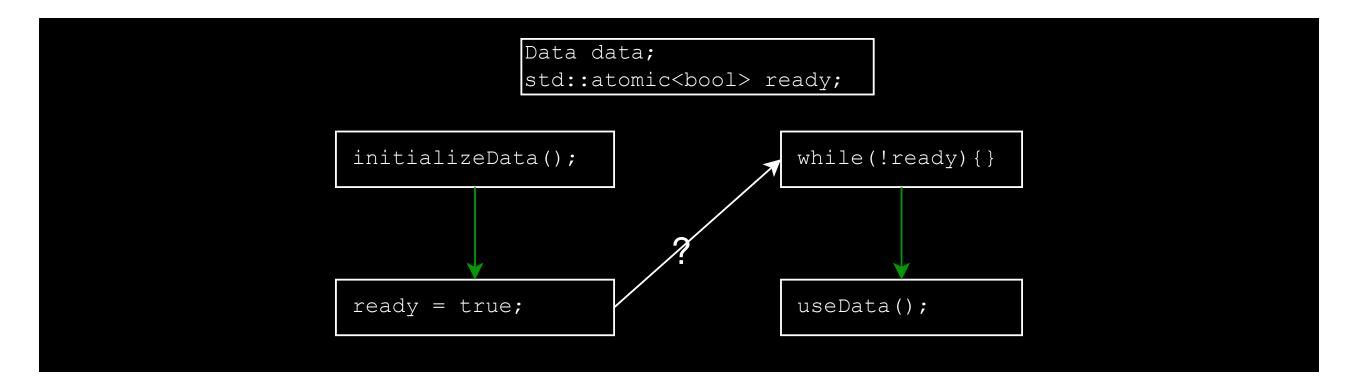
Microarchitecture (pipeline, store buffers) and cache coherence



REMEMBER

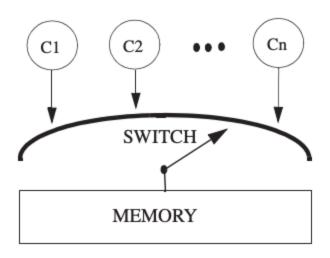


C++ MEMORY MODEL



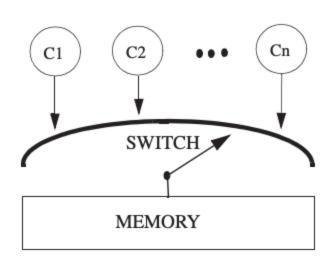


CPU MEMORY MODEL





CPU MEMORY MODEL



Producer thread

store *data reg_0

store *ready reg_1

Memory order

store *data reg_0

store *ready reg_1

load reg_0 *ready

load reg_1 *data

Consumer thread

load reg_0 *ready

load reg_1 *data



BONUS! DON'T TRUST THE ASSEMBLY



THE TWO MEMORY MODELS ANDERS SCHAU KNATTEN

C++ UNDER THE SEA 2025

Squarehead / CppQuiz.org