Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○○○

Order!
○○○

CS-453 (project)
Memory ordering

Sébastien Rouault

Distributed Computing Laboratory

September 24, 2019

Order?
●○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○○○

Order!
○○○

## Order?
### A single thread

```
// Single thread

int a = 0;
int b = 0;
print(a, b); // a = 0, b = 0

a = 1;
print(a, b); // a = ·, b = ·

b = 1;
print(a, b); // a = ·, b = ·
```

Order?
●○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○○○

Order!
○○○

# Order?

### A single thread

```
// Single thread

int a = 0;
int b = 0;
print(a, b); // a = 0, b = 0

a = 1;
print(a, b); // a = 1, b = 0

b = 1;
print(a, b); // a = ·, b = ·
```

# Order?

### A single thread

```
// Single thread

int a = 0;
int b = 0;
print(a, b); // a = 0, b = 0

a = 1;
print(a, b); // a = 1, b = 0

b = 1;
print(a, b); // a = 1, b = 1
```

# Order?

### Two threads

```
// Global var.

int a = 0;
int b = 0;

// Thread A

a = 1; // write
b = 1; // write
```

```
// Thread B

auto v = b; // read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1   □
    // a = 1, v = 0
    // a = 0, v = 1
    // a = 0, v = 0
}
```

3 / 15

## Order?

### Two threads

```
// Global var.

int a = 0;
int b = 0;

// Thread A

a = 1; // write
b = 1; // write
```

```
// Thread B

auto v = b; // read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1    ☑
    // a = 1, v = 0
    // a = 0, v = 1
    // a = 0, v = 0
}
```

3 / 15

Order?      But why complicated? ☺      C11/C++11's solutions ☺      Order!

○●      ○○○○      ○○○○○      ○○○

## Order?

### Two threads

```
// Global var.

int a = 0;
int b = 0;

// Thread A


a = 1; // write
b = 1; // write
```

```
// Thread B

auto v = b; // read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1    ☑
    // a = 1, v = 0    ☐
    // a = 0, v = 1
    // a = 0, v = 0
}
```

3 / 15

## Order?
### Two threads

```
// Global var.

int a = 0;
int b = 0;

// Thread A

a = 1; // write
b = 1; // write
```

```
// Thread B

auto v = b; // read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1    ☑
    // a = 1, v = 0    ☐
    // a = 0, v = 1
    // a = 0, v = 0
}
```

## Order?
### Two threads

```
// Global var.

int a = 0;
int b = 0;

// Thread A

a = 1; // write
b = 1; // write
```

```
// Thread B

auto v = b; // read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1   ☑
    // a = 1, v = 0   ☐
    // a = 0, v = 1   ☐
    // a = 0, v = 0
}
```

3 / 15

Order?      But why complicated? ☺      C11/C++11's solutions ☺      Order!

○●      ○○○○      ○○○○○      ○○○

## Order?
### Two threads

```
// Global var.

int a = 0;
int b = 0;

// Thread A

a = 1; // write
b = 1; // write
```

```
// Thread B

auto v = b; // read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1   ☑
    // a = 1, v = 0   ☐
    // a = 0, v = 1   ☑
    // a = 0, v = 0
}
```

3 / 15

# Order?

### Two threads

```
// Global var.

int a = 0;
int b = 0;

// Thread A

a = 1; // write
b = 1; // write
```

```
// Thread B

auto v = b; // read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1   ☑
    // a = 1, v = 0   ☐
    // a = 0, v = 1   ☑
    // a = 0, v = 0   ☐
}
```

## Order?

### Two threads

```
// Global var.

int a = 0;
int b = 0;

// Thread A

a = 1; // write
b = 1; // write
```

```
// Thread B

auto v = b; // read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1   ☑
    // a = 1, v = 0   ☐
    // a = 0, v = 1   ☑
    // a = 0, v = 0   ☐
}
```

Order?
00

But why complicated? ☺
●000

C11/C++11's solutions ☺
00000

Order!
000

# But why complicated? ☹

Order?
○○

But why complicated? ☺
○●○○

C11/C++11's solutions ☺
○○○○○

Order!
○○○

# But why complicated? ☹

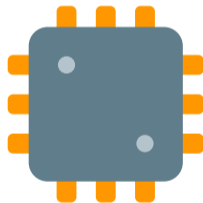## Compiler/hardware reordering

```
a = 1;
b = 1;
```

⇕

```
b = 1;
a = 1;
```

Memory consistency model?

Unrelated R/W
(& R/R, W/W)
could be carried
out-of-order.

*(More of that in
other courses, e.g., CS-471.)*

5 / 15

Order?
oo

But why complicated? ☺
ooo●

C11/C++11's solutions ☺
ooooo

Order!
ooo

## But why complicated? ☹

It even gets a bit worse. . .

```
// Global var.

int a = 0;
int b = 0;

// Thread A

a = 1; // write
b = 1; // write
```

```
// Thread B

auto v = b; // read  U.B.!!
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1   ☑
    // a = 1, v = 0   ☐
    // a = 0, v = 1   ☑
    // a = 0, v = 0   ☐
}
```

Order?
00

But why complicated? ☺
000●

C11/C++11's solutions ☺
00000

Order!
000

# But why complicated? ☹
### Main takeaway

C11/C++11  do  **not**  ensure "by default"

that  reads/writes

are  carried/observed

in  program order

by  different threads

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
●○○○○

Order!
○○○

# C11/C++11's solutions ☺

Order?
oo

But why complicated? ☺
oooo

C11/C++11's solutions ☺
o●ooo

Order!
ooo

## C11/C++11's solutions ☺

### Atomic variables

```
#include <atomic>

std::atomic<T> foo = T{};
```

With T being:

- Trivially copyable
- Copy and move constructible
- Copy and move assignable

Order?
00

But why complicated? ☺
0000

C11/C++11's solutions ☺
00●00

Order!
000

# C11/C++11's solutions ☺

Thread fences

Specifies   constraints
on the   ordering
of   memory accesses

```
#include <atomic>

std::atomic_thread_fence(std::memory_order_ /*...*/ );

std::atomic<T> foo = T{};
foo.load(std::memory_order_ /*...*/ );
foo.store(T{}, std::memory_order_ /*...*/ );
```
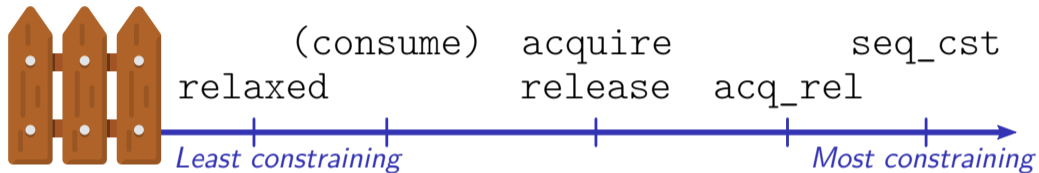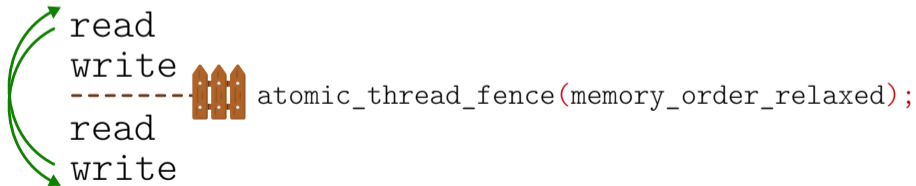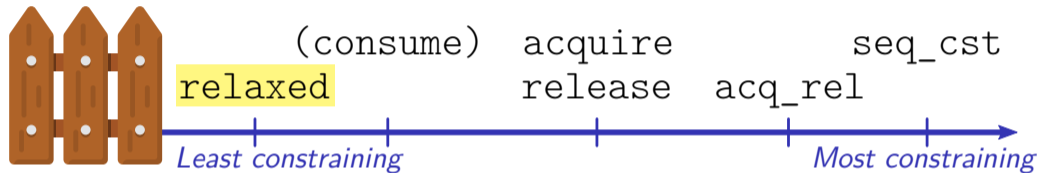
10 / 15

Order?
00

But why complicated? ☺
0000

C11/C++11's solutions ☺
000●0

Order!
000

# C11/C++11's solutions ☺

### Thread fences



```
                    (consume)  acquire            seq_cst
relaxed                        release  acq_rel
```

*Least constraining*                                    *Most constraining*

# C11/C++11's solutions ☺

### Thread fences



```
                (consume)   acquire              seq_cst
relaxed                     release   acq_rel

Least constraining                          Most constraining
```

```
read
write
------  ███  atomic_thread_fence(memory_order_relaxed);
read
write
```

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○●○

Order!
○○○

# C11/C++11's solutions ☺

### Thread fences



(consume)  acquire                    seq_cst
relaxed                  release    acq_rel

*Least constraining*                              *Most constraining*

read
write
------- �️ atomic_thread_fence(memory_order_acquire);
read
write

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○●○

Order!
○○○

# C11/C++11's solutions ☺

### Thread fences



(consume)   acquire            seq_cst
relaxed     release     acq_rel

*Least constraining*                    *Most constraining*

```
read
write
-------      atomic_thread_fence(memory_order_release);
read
write
```

11 / 15

Order?
oo

But why complicated? ☺
oooo

C11/C++11's solutions ☺
ooo●o

Order!
ooo

# C11/C++11's solutions ☺

### Thread fences



```
              (consume)   acquire        seq_cst
relaxed                   release   acq_rel
```

*Least constraining*                              *Most constraining*

```
read
write
------- 🪵 atomic_thread_fence(memory_order_acq_rel);
read
write
```

11 / 15

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○●○

Order!
○○○

# C11/C++11's solutions ☺

### Thread fences



(consume)    acquire        seq_cst
relaxed      release   acq_rel

*Least constraining*                              *Most constraining*

```
read
write
------   atomic_thread_fence(memory_order_seq_cst);
read
write                              (total order)
```

11 / 15

Order?
oo

But why complicated? ☺
oooo

C11/C++11's solutions ☺
ooooo●

Order!
ooo

# C11/C++11's solutions ☺

Thread experiment

```
// Global      // Threads {0,1}      // Threads {2,3}
a = 0;         print(a,c);           print(a,c);
b = 0;         print(b,d);           print(b,d);
c = 0;         b = 1;                d = 1;
d = 0;         a = 1;                c = 1;
```

Order?
00

But why complicated? ☺
0000

C11/C++11's solutions ☺
0000●

Order!
000

# C11/C++11's solutions ☺

### Thread experiment

```
// Global     // Threads {0,1}          // Threads {2,3}
a = 0;        print(a,c);  relaxed      print(a,c);  relaxed
b = 0;        print(b,d);               print(b,d);
c = 0;        b = 1;       relaxed      d = 1;       relaxed
d = 0;        a = 1;                    c = 1;
```

Values of (a,b,c,d) that could be read by threads 0 and 2

| Thread 0: | 0100 | 1100 | 1100 | 1001 |
|-----------|------|------|------|------|
| Thread 2: | 0001 | 1111 | 0011 | 0110 |

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○○●

Order!
○○○

# C11/C++11's solutions ☺

### Thread experiment

```
// Global    // Threads {0,1}      // Threads {2,3}
a = 0;       print(a,c);  acquire  print(a,c);  acquire
b = 0;       print(b,d);           print(b,d);
c = 0;       b = 1;       relaxed  d = 1;       relaxed
d = 0;       a = 1;                c = 1;
```

Values of (a,b,c,d) that could be read by threads 0 and 2

| Thread 0: | 0100 | 1100 | 1100 | 1001 |
|---|---|---|---|---|
| Thread 2: | 0001 | 1111 | 0011 | 0110 |

Order?
oo

But why complicated? ☺
oooo

C11/C++11's solutions ☺
oooo●

Order!
ooo

# C11/C++11's solutions ☺

Thread experiment

```
// Global      // Threads {0,1}          // Threads {2,3}
a = 0;         print(a,c); ▯▯▯ acquire    print(a,c); ▯▯▯ acquire
b = 0;         print(b,d);                print(b,d);
c = 0;         b = 1;      ▯▯▯ release     d = 1;      ▯▯▯ release
d = 0;         a = 1;                     c = 1;
```

Values of (a,b,c,d) that could be read by threads 0 and 2

| | | | |
|---|---|---|---|
| Thread 0: | 0100 | 1100 | 1100 | ~~1001~~ |
| Thread 2: | 0001 | 1111 | 0011 | ~~0110~~ |

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○○●

Order!
○○○

# C11/C++11's solutions ☺

Thread experiment

```
// Global       // Threads {0,1}      // Threads {2,3}
a = 0;          print(a,c);           print(a,c);
b = 0;          print(b,d);           print(b,d);
c = 0;          b = 1;                d = 1;
d = 0;          a = 1;                c = 1;
```

seq_cst (for Threads {0,1} print block)
seq_cst (for Threads {0,1} assignment block)
seq_cst (for Threads {2,3} print block)
seq_cst (for Threads {2,3} assignment block)

Values of (a,b,c,d) that could be read by threads 0 and 2

| Thread 0: | 0100 | 1100 | ~~1100~~ | ~~1001~~ |
| Thread 2: | 0001 | 1111 | ~~0011~~ | ~~0110~~ |

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○○○

Order!
●○○

## Order!

### Original code

```
// Global var.



int a = 0;
int b = 0;


// Thread A


a = 1; // write
b = 1; // write
```

```
// Thread B


auto v = b; // read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1   ☑
    // a = 1, v = 0   ☐
    // a = 0, v = 1   ☑
    // a = 0, v = 0   ☐
}
```

13 / 15

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○○○

Order!
○●○

# Order!

### Corrected code

```cpp
// Global var.

#include <atomic>

int a = 0;
std::atomic<int> b = 0;

// Thread A

a = 1; // write
b = 1; // atomic write
```

```cpp
// Thread B

auto v = b; // atomic read
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1   ☑
    // a = 1, v = 0   ☐
    // a = 0, v = 1   ☐
    // a = 0, v = 0   ☐
}
```

14 / 15

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○○○

Order!
○●○

# Order!

Corrected code

```
// Global var.

#include <atomic>

int a = 0;
std::atomic<int> b = 0;

// Thread A

a = 1; // write
b.store(1, release);
```

```
// Thread B

auto v = b.load( acquire );
if (v == 1) {
    print(a, v); // read
    // a = 1, v = 1    ☑
    // a = 1, v = 0    ☐
    // a = 0, v = 1    ☐
    // a = 0, v = 0    ☐
}
```

14 / 15

Order?
○○

But why complicated? ☺
○○○○

C11/C++11's solutions ☺
○○○○○

Order!
○○●

# Order!

### I want to know more

### Here you go

- https://preshing.com/...
  - ... 20120612/an-introduction-to-lock-free-programming
  - ... 20120913/acquire-and-release-semantics
- https://en.cppreference.com/w/{c,cpp}/...
  - ... atomic{,/memory_order}
  - ... language/memory_model
- Memory Barriers: a Hardware View for Software Hackers, Paul E. McKenney

### Next time

- *Read–Modify–Write* atomic primitives (e.g. compare & swap)
- Workshop: "Writing my own lock"