# ELT WFRTC
# Software patterns and library solutions for low latency multithreading

*Calle Rosenquist, ESO*

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

# Real time core

- CPU-based

- Single server (currently 128 core AMD EPYC Zen 3)

- Minimal Linux

- Control software

  - Single process

  - Statically linked

  - C++

  - Real-time threads running on isolated cores without scheduler ("isolcpus")
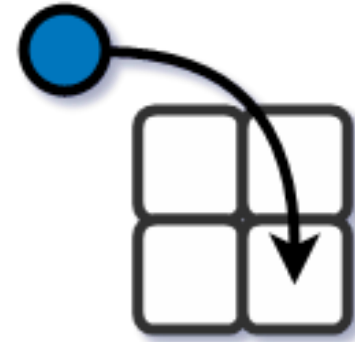
ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*
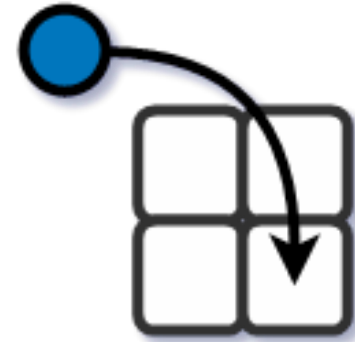
# NUMA
## non-uniform memory access

# "Pinning" real time threads

- In Linux each thread have two masks

  - "CPUs allowed" - where to run

  - "Mems allowed" - where to allocate memory

- By *pinning* we allow only one core and its local memory

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Memory allocations

- Critical heap memory allocated from specific NUMA nodes

- Avoid allocators that rely on thread policy

  - E.g. default `operator new`, `malloc()`

  - Not easily predictable

  - Memory may come from allocator "free list"

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

# Support library: NUMA++

*Creating pinned threads*

```cpp
#include <numapp/memory.hpp>
#include <numapp/thread.hpp>

template <class F, class... Args>
auto MakePinnedThread(int cpu, std::string_view name, F&& f, Args... args)
    -> std::thread {
    using namespace numapp;
    auto node = GetNodeOfCpu(cpu);
    if (!node.has_value()) {
        throw std::invalid_argument("cpu invalid");
    }

    NumaPolicies policies;
    policies.SetCpuAffinity(CpuAffinity::MakeBindCpu(cpu));
    policies.SetMemPolicy(MemPolicy::MakeBindNode(*node));
    return MakeThread(
        name, policies, std::forward<F>(f), std::forward<Args>(args)...);
}
```

In short

- Lookup node from CPU core

- Create policies

  - Pin thread to core

  - Memory from local node

- Create thread with provided policies

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Support library: NUMA++

*Allocating memory from a specific node*

```cpp
#include <numapp/memory.hpp>
#include <numapp/mempolicy.hpp>

void* BindAlloc(std::size_t size, int node, std::error_code& ec) {
    using namespace numapp;

    MemPolicy policy = MemPolicy::MakeBindNode(node);
    void* ptr = Allocate(size, policy, MemPolicyFlag::Strict, ec);
    if (ec) {
        return nullptr;
    }
    ec = MemLock(ptr, size, LockFlag::PreFault);
    if (ec) {
        return nullptr;
    }
    return ptr;
}
```

- Create policy to use single node

- Allocate

  - Uses `mmap()`

- Lock and pre-fault to trigger physical memory allocation

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

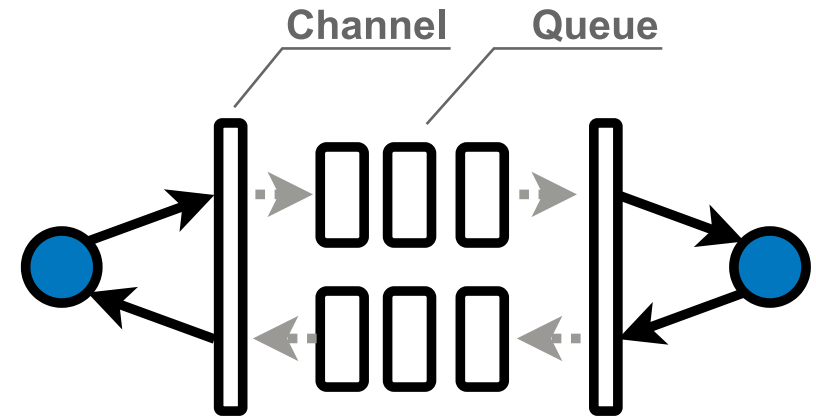# Inter-thread queues

# Inter-thread queues

- Deliver data from single producer to single consumer without locking or blocking

- Fix capacity at construction to avoid reallocation

- Use *views* (e.g., pointer + size) to avoid unnecessary memory copies

Data Classification: PUBLIC

# Channels

*WFRTC composition of two directional queues*

- Second "free-list" queue

- Two queues forms a *Channel*

- boost::lockfree::spsc_queue



```
1   template <class Queue>
2   class Channel {
3   public:
4       using ValueType = typename Queue::value_type;
5       using QueueType = Queue;
6
7       explicit constexpr Channel(QueueType& read, QueueType& write) noexcept;
8
9       void Push(ValueType const& value);
10      bool TryPush(ValueType const& value) noexcept;
11
12      void Pop(ValueType& value);
13      bool TryPop(ValueType& value) noexcept;
14  };
```

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Benchmark

- Measures push/pop operations under maximum contention

- Unpinned threads on unisolated cores

```
Running wfhrtcSpscQueueBenchmark
Run on (128 X 2450 MHz CPU s)
CPU Caches:
  L1 Data 32 KiB (x128)
  L1 Instruction 32 KiB (x128)
  L2 Unified 512 KiB (x128)
  L3 Unified 32768 KiB (x16)
Load Average: 0.15, 0.03, 0.11
-------------------------------------------------------------------

Benchmark                  Time           CPU    Iterations UserCounters...
-------------------------------------------------------------------

PopFixture/SpscPush      10.9 ns      10.9 ns      64391865 FullPercent=5.03837
PushFixture/SpscPop      34.9 ns      34.8 ns      19543691 EmptyPercent=2.70725
```
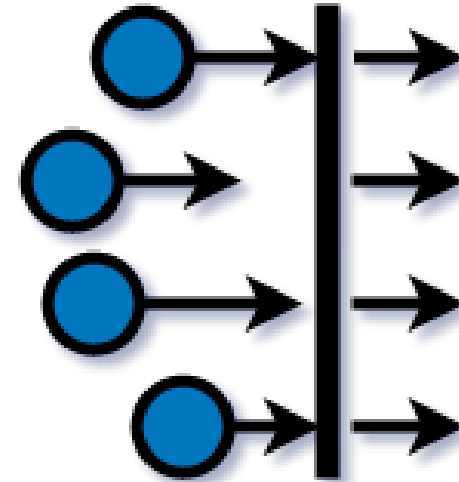
*AMD EPYC 7763 (128 Core Zen 3 Milan)*

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Scatter-gather

# Scatter-gather

WFRTC use this to

- Offload work to multiple threads
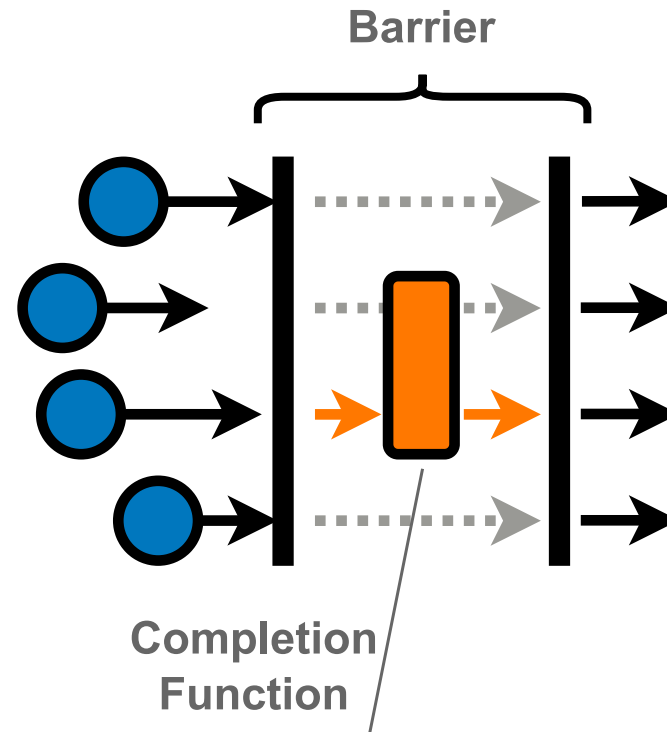
- Gather results when all threads have completed

Achieved using the *barrier* primitive

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

# Barrier

*Introducing the barrier*

- Threads *arrive and wait*

- When last thread arrives a *Completion Function* is invoked by one of the threads

  - *Completion Function* can **safely modify shared state**

- When *Completion Function* returns barrier is lifted and threads continue

- C++20 `std::barrier`



Barrier

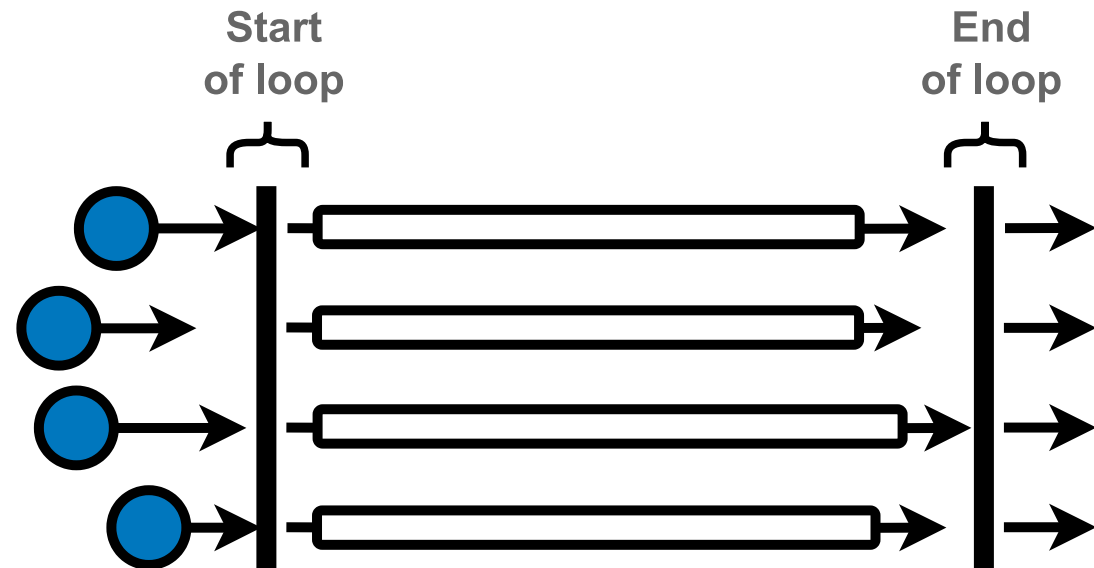Completion Function

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Barriers in WFRTC

*Computation loop*

**Start of loop barrier**

- Wait for start condition, e.g.
  - Inputs from *channel*
  - State change *signal*
- Distribute data or signal to threads

**End of loop barrier**

- Gather results
- Push result to output *channel*

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Support library: ion

*ion::Barrier*

- Semantics like `std::barrier`
- Busy-waits

```
1   template <class CompletionFunction = *unspecified*>
2   class Barrier {
3   public:
4       Barrier(std::ptrdiff_t expected, CompletionFunction func);
5
6       void ArriveAndWait();
7
8       auto Arrive(std::ptrdiff_t n = 1) -> ArrivalToken;
9       void Wait(ArrivalToken&& token) const;
10  };
```

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist,*
*2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Support library: ion

*Benchmark*

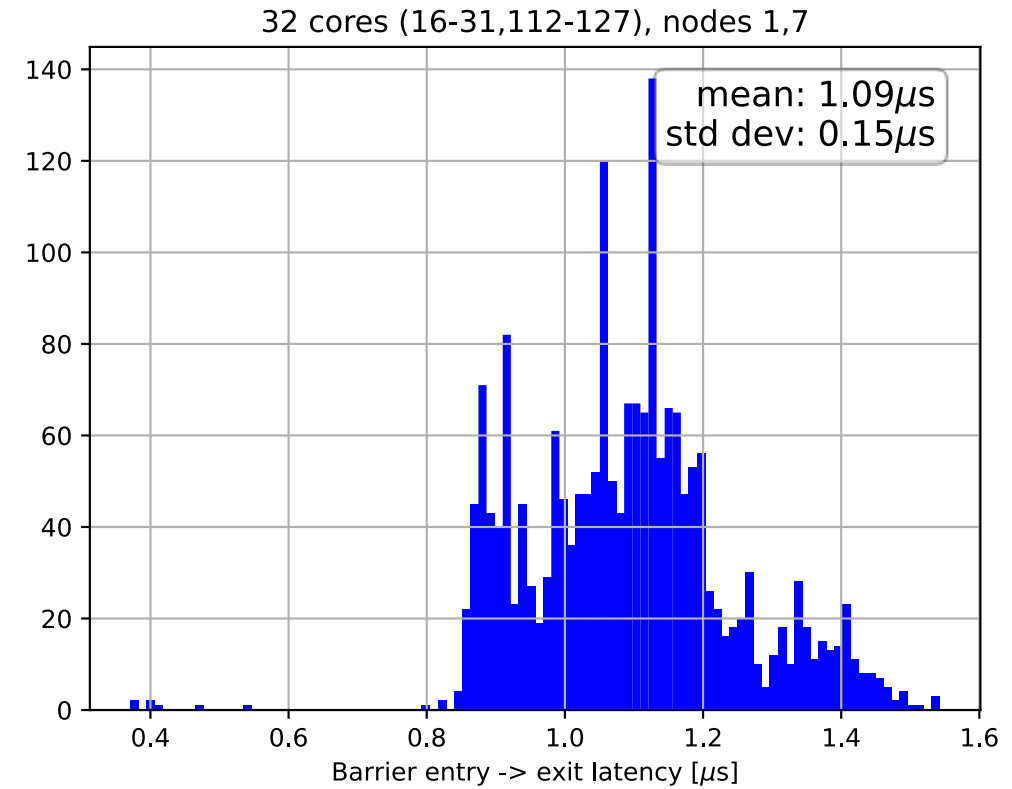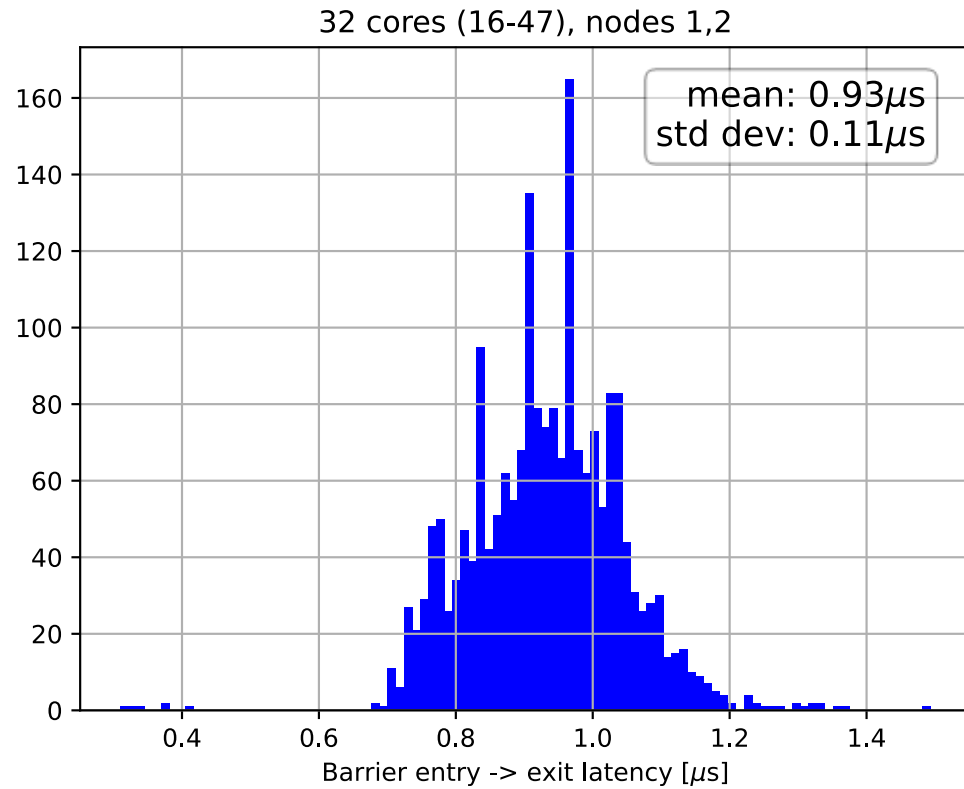One measure of synchronization latency

- Time between last thread arriving and first thread leaving barrier

- Empty completion function

- Last 2'000 samples from 10'000 iterations



*AMD EPYC 7763 (128 Core Zen 3 Milan)*

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Support library: ion

*As threads and NUMA distance increase so does latency*



32 cores (16-47), nodes 1,2 — mean: 0.93$\mu$s, std dev: 0.11$\mu$s

32 cores (16-31,112-127), nodes 1,7 — mean: 1.09$\mu$s, std dev: 0.15$\mu$s

Barrier entry -> exit latency [$\mu$s]

*AMD EPYC 7763 (128 Core Zen 3 Milan)*

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Signals

Data Classification: PUBLIC

# Signals

*Notify, observe, wait*

WFRTC use signals to

- Request state change (state enumeration)

- Indicate current state (state enumeration)

- Trigger telemetry sending (sample id)

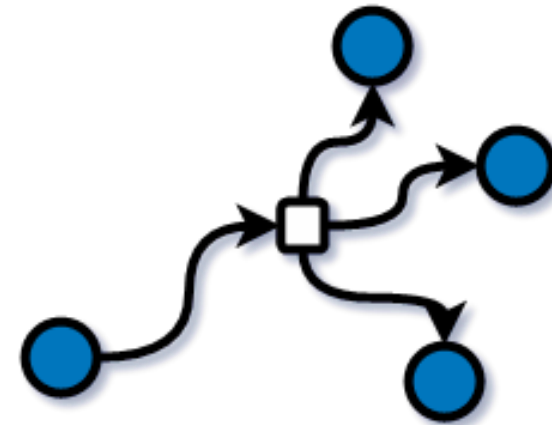ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Support library: ion

*Thin abstraction on top of `std::atomic<T>`*

`ion::SignalSource<T>`

- Lightweight thread safe, lock-free atomic value of type T

- No synchronization or order constraints

`ion::SignalToken<T>`

- Keeps association with source and remember last value

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Support library: ion

```cpp
template <class T>
class SignalSource {
public:
    using ValueType = T;
    explicit SignalSource(ValueType initial = ValueType{});
    auto Load() const -> ValueType;
    void Store(ValueType value);

    auto CompareExchangeWeak(ValueType& expected, ValueType desired) -> bool;
    auto CompareExchangeStrong(ValueType& expected, ValueType desired) -> bool;
};

template <class T, class UnaryOperation>
auto CompareExchangeTransform(SignalSource<T>& signal,
                              UnaryOperation op) noexcept -> T;
```

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

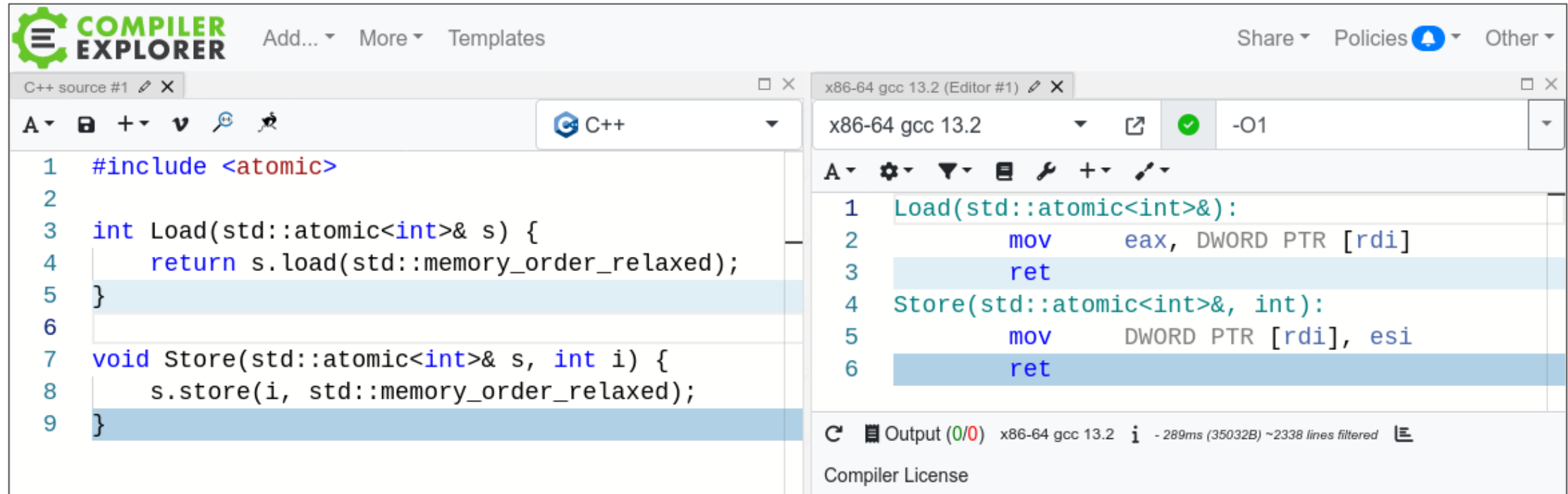# Support library: ion

```cpp
template <class T>
class SignalToken {
public:
    using ValueType = T;
    SignalToken(SignalSource<T> const* signal,
                std::optional<ValueType> last = std::nullopt);

    auto Load() const -> ValueType;
    auto GetLast() const -> ValueType;
    void SetLast(ValueType last);
    auto Update() -> ValueType;
};

template <class T>
auto Wait(SignalToken<T>& token) -> T;

template <class T, class Predicate>
auto Wait(SignalToken<T>& token, Predicate&& stop_waiting) -> T;

template <class T1, class T2>
auto WaitAny(SignalToken<T1>& token1, SignalToken<T2>& token2)
    -> WaitAnyResult<T1, T2>;
```
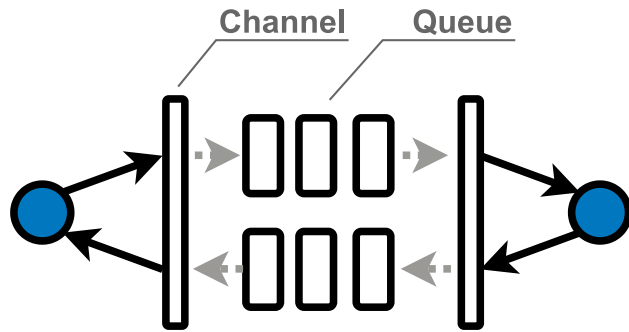
Data Classification: PUBLIC

# Support library: ion

Signals under the hood



*https://godbolt.org/z/6d7jMfefW*

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Summary

**1:1**

**1:N:1**

**N:M**

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Thank you!

**Calle Rosenquist**

**crosenqu@eso.org**

**Boost.Lockfree**
*www.boost.org/doc/libs/1_83_0/doc/html/lockfree.html*

**NUMA++ & ion**
*gitlab.eso.org/rtctk/roadrunner*

@ESOAstronomy

@esoastronomy

@ESO

european-southern-observatory

@ESOobservatory

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Bonus Slides

Data Classification: PUBLIC

# Barrier Example 1/2

```cpp
1   #include <chrono>
2   #include <iostream>
3   #include <thread>
4   #include <vector>
5
6   #include <ion/barrier.hpp>
7
8   /** Completion function */
9   void Complete() noexcept {
10      std::cout << "Iteration beginning\n";
11  }
12
13  using Barrier = ion::Barrier<void (*)() noexcept>;
14
15  void Worker(Barrier& barrier, std::size_t num_iterations) {
16      for (auto count = 0u; count < num_iterations; ++count) {
17          // Synchronize start of each iteration
18          barrier.ArriveAndWait();
19
20          // Fake some work
21          using namespace std::chrono_literals;
22          std::this_thread::sleep_for(100ms);
23      }
24  }
```

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Barrier Example 2/2

```cpp
26  int main() {
27      auto const num_threads = 10;
28      auto const num_iterations = 100;
29
30      Barrier barrier(num_threads, &Complete);
31
32      std::vector<std::thread> threads;
33      threads.reserve(num_threads);
34      for (auto thread_idx = 0; thread_idx < num_threads; ++thread_idx) {
35          threads.emplace_back(&Worker, std::ref(barrier), num_iterations);
36      }
37
38      // Let it complete and then join
39      for (auto& thread : threads) {
40          thread.join();
41      }
42  }
```

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

Data Classification: PUBLIC

# Signal Example

```cpp
1   #include <cstdint>
2   #include <iostream>
3   #include <thread>
4
5   #include <ion/signal.hpp>
6
7   int main() {
8       auto source = ion::SignalSource<int>();
9
10      // To avoid race-condition between storing new signal value and what token
11      // initializes to we create the token in the main thread.
12      std::thread thread([token = ion::SignalToken<int>(&source)]() mutable {
13          auto current = Wait(token);
14          std::cout << "Got signal: " << current << std::endl;
15      });
16
17      // Update signal
18      source.Store(42);
19
20      thread.join();
21  }
```

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist, 2023-11-08, RTC4AO*

# Support library: ion

*Benchmark*

```
Running ionBenchmark
Run on (128 X 2450 MHz CPU s)
CPU Caches:
  L1 Data 32 KiB (x128)
  L1 Instruction 32 KiB (x128)
  L2 Unified 512 KiB (x128)
  L3 Unified 32768 KiB (x16)
Load Average: 0.00, 0.06, 0.12
------------------------------------------------------------------
Benchmark                                       Time           CPU   Iterations
------------------------------------------------------------------
BenchmarkSignalLoad/real_time/threads:1       0.410 ns      0.409 ns   1000000000
BenchmarkSignalLoad/real_time/threads:3       0.137 ns      0.409 ns   3000000000
BenchmarkSignalLoad/real_time/threads:5       0.082 ns      0.409 ns   5000000000
BenchmarkSignalLoad/real_time/threads:7       0.059 ns      0.409 ns   7000000000
BenchmarkSignalLoad/real_time/threads:9       0.046 ns      0.409 ns   9000000000
BenchmarkSignalLoad/real_time/threads:11      0.037 ns      0.409 ns   11000000000
BenchmarkSignalLoad/real_time/threads:13      0.032 ns      0.409 ns   13000000000
BenchmarkSignalLoad/real_time/threads:15      0.027 ns      0.409 ns   15000000000
BenchmarkSignalLoad/real_time/threads:16      0.026 ns      0.409 ns   16000000000
BenchmarkSignalStoreLoad/real_time/threads:1  0.410 ns      0.409 ns   1000000000
BenchmarkSignalStoreLoad/real_time/threads:3  0.204 ns      0.609 ns   3000000000
BenchmarkSignalStoreLoad/real_time/threads:5  0.125 ns      0.624 ns   5000000000
BenchmarkSignalStoreLoad/real_time/threads:7  0.094 ns      0.658 ns   6899152162
BenchmarkSignalStoreLoad/real_time/threads:9  0.075 ns      0.674 ns   9000000000
BenchmarkSignalStoreLoad/real_time/threads:11 0.063 ns      0.688 ns   10598253435
BenchmarkSignalStoreLoad/real_time/threads:13 0.056 ns      0.727 ns   12852271835
BenchmarkSignalStoreLoad/real_time/threads:15 0.047 ns      0.707 ns   14715672750
BenchmarkSignalStoreLoad/real_time/threads:16 0.044 ns      0.710 ns   15882213616
```

*Divided by # threads*

Under varying number of unpinned threads:

- `SignalSource::Load()`

- `SignalSource::Store()` [thread 1]
  `SignalSource::Load()` [threads 2..N-1]

*AMD EPYC 7763 (128 Core Zen 3 Milan)*

ELT WFRTC Software patterns and library solutions for low latency multithreading, *Calle Rosenquist,*
*2023-11-08, RTC4AO*

Data Classification: PUBLIC