

# PARALLEL PROGRAMMING IN C++23

*Jakub Marecek and team*

version March 13, 2023

<https://github.com/jmarecek/book-parallel-cpp>

March 13, 2023

This work may be distributed and/or modified under the conditions of the CC-BY-NC-SA license.

Original copy of this book can be obtained at <https://github.com/jmarecek/book-parallel-cpp>.

This copy of the book is dated March 13, 2023.

# Preface

These lecture notes are intended for PDV course at the Czech Technical University, but could perhaps be of more widespread interest.

## About the authors

Jakub Marecek first taught a course on “Advanced C++ Programming” in 1999 in Brno, the Czech Republic. Since then, he has worked at the University of Nottingham, ARM Ltd., the University of Edinburgh, IBM Research, and the University of California, Los Angeles.

A wider team has contributed:

- Simon Toth has been a major source of inspiration for the project with his book *The Standard Algorithms in C++*. The build system for the lecture notes is taken from the book almost in verbatim, in compliance with the Creative Commons licence. The lecture notes also feature several examples from his book.
- Illia Kryvoviaz has fixed numerous issues with the examples and added the GodBolt links.
- Jan Mrkos has provided a number of suggestions that improved the SYCL chapter in particular.

## Feedback

If you like these lecture notes, please do send us feedback: positive, negative. Anything goes!

## Why cc-by-sa-nc

This book is licensed CC-BY-SA-NC, which allows derivative works, but does not permit commercial use including selling printed copies.



# Contents

<b>Preface</b>	<b>iii</b>
<b>1 Introduction</b>	<b>3</b>
<b>2 The Concepts</b>	<b>5</b>
2.0.1 Data race . . . . .	5
2.0.2 Deadlock . . . . .	6
2.0.3 Concurrent programming . . . . .	6
2.1 Structuring code: Processes, Threads, Tasks, Coroutines . . .	7
2.2 Hardware support . . . . .	8
2.2.1 Memory hierarchies . . . . .	8
2.2.2 GPGPUs . . . . .	9
2.2.3 Memory order . . . . .	10
2.2.4 Compare and swap . . . . .	11
2.3 Synchronization primitives . . . . .	13
2.3.1 Raw synchronization primitives . . . . .	13
2.3.2 Further synchronization features . . . . .	14
<b>3 The Syntax in C++23</b>	<b>15</b>
3.1 Threads, Tasks, Coroutines . . . . .	15
3.1.1 C++11 thread . . . . .	15
3.1.2 C++20 jthread . . . . .	17
3.1.3 Coroutines . . . . .	21
3.2 Synchronisation Primitives . . . . .	32
3.2.1 Atomic Variables . . . . .	32
3.2.2 Mutexes and Locks . . . . .	34
3.2.3 Barrier . . . . .	36
3.3 Algorithms in the Standard Template Library . . . . .	38
3.3.1 For Each . . . . .	38
3.3.2 Reduce . . . . .	39
3.3.3 Merge . . . . .	40

<b>4</b>	<b>The Syntax in OpenMP</b>	<b>41</b>
4.1	Threads, Tasks, Coroutines	42
4.1.1	OpenMP Task Region	42
4.1.2	Threads and their Sizing	45
4.1.3	Sections	45
4.1.4	Tasks	47
4.1.5	Kernels	48
4.2	Synchronisation Primitives	51
4.2.1	Atomic Variables	51
4.2.2	Reductions	52
4.2.3	Mutexes	52
4.2.4	Critical Sections	53
4.2.5	Barrier	54
4.2.6	Fences and Flushes	56
4.3	Algorithms	57
4.3.1	For Each	57
4.4	Exercises	58
<b>5</b>	<b>The Syntax in SYCL</b>	<b>59</b>
5.1	Motivation	59
5.2	More concepts	63
5.2.1	Device selector	63
5.2.2	Queues	63
5.2.3	Work items, Work groups, and Kernels	64
5.2.4	Asynchronous errors	65
5.2.5	Unified shared memory	66
5.2.6	Buffers and accessors	67
5.2.7	Barrier	68
5.3	More complex examples	68
5.4	Building code	71
<b>6</b>	<b>Introducing parallel data structures</b>	<b>73</b>

# Chapter 1

## Introduction

In a modern desktop computer equipped with a GPGPU, a single-threaded application cannot easily make use of more than 0.05 % of the available performance. One hence has to learn the dark art of parallel programming to make a full use of the available performance.

Indeed, consider an AMD Threadripper 3990X with NVIDIA GeForce RTX 4090 Ti. There, a single threaded application can utilise circa 49 GFLOPS, which is less than 0.05 % of the overall performance, when one runs a single-threaded application once. A multi-threaded application can utilise perhaps 3732 GFLOPS. A multi-threaded application not making use of the GPGPU can hence make use of less than 4 % of the overall performance. In contrast, a multi-threaded application can make use of almost 100 TFLOPS. (To unlock this performance, one needs to use recent versions of OpenMP, SYCL, or vendor- or application-specific interfaces such as CUDA, OpenGL, or DirectX Compute. We will focus on OpenMP and SYCL in subsequent Chapters.) Similar breakdown is available in a modern mobile phone (with a substantial mobile GPGPU, such ARM Mali-G710 MP7 in Google Tensor G2) and modern supercomputers.

This is related to the changing “shape” of computers: Computers no longer get faster, just “wider” and “heterogeneous”. Data-parallel computing is most scalable approach to the “wider” computers.

In these lecture notes, you will learn about parallel computing and data-parallel computing, in particular. In the first three chapters, this we will learn to design multi-threaded applications. First with C++23, and subsequently with OpenMP (“just in case”. In the following chapters, we will extend this to GPGPU-enabled systems and SYCL, and illustrate in a number of application domains.





## Chapter 2

# The Concepts

Parallelism means two or more tasks can be executed simultaneously. This is an option, which the compiler and operating system and processor can exercise, but does not come with any guarantees. Often, this means no shared variables or other resources, and need not require any synchronization primitives.

Concurrency means that two or more tasks start, run, and complete in overlapping time periods, while sharing some resources. If two tasks concurrently set shared variable  $x$  to 1 and 2, it is not clear what value it would have, subsequently. More broadly, concurrent access to a mutable shared memory can result in issues without the use of synchronization primitives (“data race”) and with the use of synchronization primitives (“deadlock”).

### 2.0.1 Data race

When we need to ensure mutual exclusion in access to two or more shared variables, e.g., read value of one of the variables and add it to another variable, we may need to use some synchronization primitives (e.g., mutexes). Without the use of synchronization primitives, we are facing the risk of a data race.

For example, consider the situation in a bank, where there are two clients: Alice and Bob. Transaction T1: Bob has \$100 in his account, but will be paying a \$50 bill. At the same time, in Transaction T2, Alice will be paying \$100 to Bob. Depending on the ordering of the reading and writing operations, one may obtain several outcomes:

- Transaction T1 will read \$100 valued of Bob’s account. Transaction T2 will read \$100 value. Transaction T1 will write \$200. Transaction T2 will write \$50 value.
- Transaction T1 will read \$100 valued of Bob’s account. Transaction T1 will write \$50. Transaction T2 will read \$50 value. Transaction T2 will write \$150 value.

- Transaction T1 will read \$100 valued of Bob's account. Transaction T2 will read \$100 value. Transaction T1 will write \$50. Transaction T2 will write \$200 value.
- Transaction T2 will read \$100 value. Transaction T2 will write \$200 value. Transaction T1 will read \$200 valued of Bob's account. Transaction T1 will write \$150.

Either Bob or the bank could be \$100 short.

## 2.0.2 Deadlock

When we need to ensure mutual exclusion in access to two or more shared variables, e.g., two temporary results associated with two mutexes, one may naively lock the first mutex first, and subsequently lock the other mutex. This, however, can lead to a deadlock. One needs to lock both mutexes at the same time.

Naively, one could run:

Locking multiple mutexes at once.

```
1 void thread_operation(){
2     std::scoped_lock l(mutex1,mutex2);
3     ...
4     complicated_task();
5 }
```

[Open in Compiler Explorer](#)

In theory, a deadlock (Czech: “problém uváznutí”) can occur when:

- each lock is owned by one thread
- each thread has locked at least one lock and needs to lock at least one more lock
- it is impossible to remove the lock ownership
- there is a cyclic dependency among the lock-using threads.

## 2.0.3 Concurrent programming

There are two essential models for concurrent programming: shared memory and message passing. In sharing memory, we have broadly four options:

- Confinement: Do not share memory between threads. This is often impossible.
- Immutability: Do not share any mutable data between threads.

- Thread-safe code: Use data types with additional guarantees for storing any mutable data shared between threads, or even better, use implementations of algorithms that are already parallelized and handle the concurrency issues for you. For example in C++, one can use the standard template library with a suitable execution policy. In particular, the header `execution` defines objects `std::execution::seq`, `std::execution::par`, `std::execution::par_unseq`, which can be passed as the first argument of any standard algorithm, e.g.,  
`std::vector<int> my_data; std::sort(std::execution::par, my_data.begin(), my_data.end())`  
 See Section 3.3.1 for more examples.
- Synchronization: Use synchronization primitives to prevent accessing the variable at the same time. This option is explored in this chapter in more detail.

Eventually, we will see that message passing can be implemented using the synchronization primitives and may be the least challenging to use correctly.

## 2.1 Structuring code: Processes, Threads, Tasks, Coroutines

Processes, threads, tasks, and coroutines execute instructions.

A *process* provides all of the prerequisites for executing instructions: loads an executable program, sets up a virtual address space, the environment (e.g. environment variables and a security context), the process control block (PCB, often stored in registers of the processor and on a per-process stack in kernel memory), opens handles to system objects (e.g., files, sockets), and often much more. In some sense, one can imagine “a virtual machine”.

Within a particular process, there is at least one *thread*. All threads of a particular process share the same virtual address space and handles to system objects. Each thread, independently, operates its own context (registers, stack, exception handlers). Unless declared otherwise, threads of a particular process share memory and are allocated “time slices” by the operating system. This can be seen as a “virtual processor” within a “a virtual machine” of a process, often with no guarantees on the time slicing.

All modern operating systems (OS) are multitasking, i.e., running multiple processes with the operating system forcibly interrupting the run one process to execute another process after a certain amount of time (“pre-emptive scheduling”). Switching between the processes involves swapping the process control block (PCB). In Intel architectures, this is known as the task state segment (TSS), and there is hardware support for the switch. AMD64 does not support task switches in hardware. Consequently, neither

Windows nor Linux kernels utilize the hardware support for the switch. Context switching thus has non-trivial impact on performance.

Most modern processors are multi-core and support multithreading in some form. This means that one each process can execute multiple “hardware threads” and there is some support for switching between those. In Intel architectures, hyper-threading means each hardware core can execute multiple threads, e.g., two, to take advantage of idle time (e.g., loading data, network communications).

Most modern operating systems *do not guarantee* fairness among the threads.

Within a particular thread, one may utilize multiple *coroutines*, which can be seen as subroutines that can run in multiple steps, but sometimes can serve as a light-weight alternative to hardware threads. Coroutines can be called, can return when completed, but also can suspend themselves, yielding control and partial results, and be resumed by another co-routine. Typical uses involve generators and factories and various other concepts within “lazy evaluation”, as well as event-driven architectures within cooperative multi-tasking.

That is: two coroutines within one thread never run in parallel, but one can have the runs of two or more coroutines interleaved. We can suspend a co-routine in one thread and resume it within another thread.

As it turns out, the “context switch” with user-level threads has a similar cost to a function call or suspending a coroutine (`co_yield`). Indeed, coroutines are typically implemented with user-level threads, which leads to cheaper context-switch compared with hardware threads. Within the user-level threads, one can distinguish stackful and stackless versions, where coroutine state is saved on the heap (as in C++).

A *task* is a rather abstract unit of work, e.g., a function, which can be executed by any thread, but often allocated to one of a many threads within a pool.

## 2.2 Hardware support

### 2.2.1 Memory hierarchies

Most modern computers utilize complicated cache hierarchies. One may have 32 KB of L1 cache and 256-512 KB of L2 cache on each core, with cca. 1 ns and 3-4 ns latencies, respectively. (This will get reported as 32 MB of L2 cache on a Threadripper 3990X.) One may have many megabytes of L3 cache shared between the cores in some non-uniform fashion, with non-uniform latency of 10-20 ns. Finally, one may access the memory in cca. 60 ns. Compare these times with the 0.27 ns cycle time of a processor clocked at 3.7 GHz.

Crucially, the times above are the best-case performance. To understand the worst-case performance, one needs to understand cache snooping (with a write-invalidate protocol). Most modern cores monitor a stream of memory loads and stores across the whole shared-memory machine. If a transaction concerns store in a block cached on the core (processor), the core sets some “invalid” flag to the block in its cache. If the core wishes to read a block with the “invalid” flag, it will ask all other cores for a valid value. Obtaining the correct value takes time.

As we will see in the distributed-memory part of the course, one can also maintain the coherence otherwise, which is known as a directory-based coherency mechanism. Either way, cache coherence defines the requirements on concurrent reads and writes to the same address. To some extent, it masks the presence of cache hierarchies in any modern computing system.

## 2.2.2 GPGPUs

As we have mentioned in the introduction, the bulk of the computational performance of most modern computing system (from a mobile phone to a supercomputer) is in its GPGPUs. Let us briefly consider a particular example of the NVIDIA Ampere architecture of GeForce RTX 3080 or NVIDIA A100:

- There are seven Graphics Processing Clusters (GPCs), sharing up to 40 MB of L2 cache and up to 40 GB of high-speed HBM2 memory
- Within each GPC, there are 12 Streaming Multiprocessors (SMs),
- Within each SM, there are 128 cores working with single-precision floating-point (FP32) precision and two double-precision (FP64) units. There is also 128 KB of L1/Shared Memory, shared across the 128 cores.
- Each SM is partitioned into four processing blocks (or partitions), each with a few kilobytes of L0 instruction cache and one warp scheduler.

Altogether, A100 has  $7 \cdot 12 \cdot 128 = 10752$  cores, but their use is rather constrained. Each warp can have at most 32 threads and runs them in lock-step on one processing blocks. Each streaming multiprocessors can run at most 64 warps of 32 threadblocks (i.e., 2048 threads per SM). Further constraints are due to the register use: each thread can use at most 255 registers, but there are only 65536 32-bit registers for the SM (yielding a limit of 257 threads per SM, at full register utilization). Further constraints are due to the use of memory hierarchy (esp. the 128 KB of shared memory, shared across the 128 cores).

Similar to the CPU, the GPU hence has a memory hierarchy:

- L1 cache with 33 cycle latency and shared memory with even lower latency, based on microbenchmarking [?]
- L2 cache with up to 2080 GB/s read bandwidth (200 cycle latency), based on microbenchmarking [?]
- on-board HBM2 memory with 1555 GB/sec bandwidth (290 cycle latency)
- intra-board NVLink with 50 Gb/sec per signal pair bandwidth
- access to RAM via PCI Express Gen 4 (PCIe Gen 4) at 31.5 GB/sec
- optionally, intra-node communication at 200 Gbit/sec using InfiniBand.

The interaction of the GPGPU memory hierarchy and CPU memory hierarchy is non-trivial, but summarized by the suggestion to reduce the number and volume of transfers between the host and the GPGPU, even at the price of increasing the volume of computation substantially. (Compare the numbers above to M.2 PCIe Gen4 SSDs with 7 GB/sec bandwidth.)

### 2.2.3 Memory order

Thus, there are several options for implementing synchronization primitives, known as “memory orders” or memory consistency. They define allowed behaviours of loads and stores to different addresses. All guarantee atomicity and modification-order consistency.

In `memory_order_relaxed`, no further guarantees are provided and specifically no order is imposed on concurrent memory accesses. This is also how weakly-ordered architectures (e.g. ARM, NVIDIA) operate, by default: if two threads access shared memory, the load in one thread does not have to read a value written by another thread very recently, unless they use synchronization primitives (e.g. a fence).

With `memory_order_release` and `memory_order_acquire` specifiers, we force weakly-ordered architectures to behave closer to strongly-ordered architectures (e.g., Intel). If one thread writes into shared memory atomically with `memory_order_release` and another thread reads the memory atomically with `memory_order_acquire`, the load in the second thread is guaranteed to read the value written by another thread.

With `memory_order_seq_cst` (sequential consistency), we additionally require a single total ordering of all modifications (with this specifier). A load with this specifier gets its value either from the last store with this specifier or from some store without this specifier that did not precede the most recent `memory_order_seq_cst` store. This is the default option and the strongest form of consistency, but it may also seem weak, in that one

considers the ordering of the instructions the processor chose, rather than the ordering the programmer wrote down.

(In earlier version of C++ standard, there were further memory models defined for the sake of DEC Alpha architecture. At least `memory_order_consume` is deprecated as of C++17.)

See [1] for an informal overview, [2] for a formal overview, and [https://www.youtube.com/watch?v=A\\_vAG6LIHwQ&ab\\_channel=ACCUConference](https://www.youtube.com/watch?v=A_vAG6LIHwQ&ab_channel=ACCUConference) for a nice lecture.

## 2.2.4 Compare and swap

Synchronization primitives are typically implemented using some hardware instructions, typically “compare-and-swap”. In locking, these make it possible to test whether the lock is free, and if so, acquire the lock within a single operation that the hardware guarantees to execute atomically.

The atomic compare and swap“ (CAS) instruction dereferences a pointer to an atomic variable and compares its value against a given value. If there is a match, it replaces the atomic variable with a given new value. That is:

- we declare an atomic variable (and a pointer to it)
- (\*) we save the value of an atomic variable to a local, private variable (by dereferencing the pointer)
- based on the saved value in a local, private variable, we compute the new value, which we would like to store in the atomic variable
- the CAS instruction is used. If the current value matches the value saved in the local, private variable, we will overwrite the value with the newly computed value. If the current value no longer matches the value saved in the local, private variable, we wait (some random and growing from a small starting value) and repeat from (\*).

In C++, the `atomic` header defines two variants of “compare and swap” and a specialization thereof for pointers:

Definitions of two variants of “compare and swap”.

```
1 // based on
  ↪ https://github.com/gcc-mirror/gcc/blob/master/libstdc%2B%2B-v3/include/std/atomic
2
3 template<typename _Tp> struct atomic {
4     using value_type = _Tp;
5
6     [...]
7
8     bool compare_exchange_strong(_Tp& __e, _Tp __i, memory_order
  ↪ __s,
9         memory_order __f) noexcept {
10         return __atomic_impl::__compare_exchange(_M_i, __e, __i,
  ↪ false, __s, __f);
11     }
12
13     bool compare_exchange_weak(_Tp& __e, _Tp __i, memory_order __s,
14         memory_order __f) noexcept {
15         return __atomic_impl::__compare_exchange(_M_i, __e, __i, true,
  ↪ __s, __f);
16     }
17
18     [...]
19 };
20
21 // Partial specialization for pointer types.
22 template<typename _Tp> struct atomic<_Tp*> {
23     using value_type = _Tp*;
24     using difference_type = ptrdiff_t;
25
26     [...]
27
28     bool compare_exchange_strong(__pointer_type& __p1,
  ↪ __pointer_type __p2, memory_order __m1, memory_order __m2)
  ↪ noexcept {
29         return _M_b.compare_exchange_strong(__p1, __p2, __m1, __m2);
30     }
31
32     [...]
33
34     };
```

[Open in Compiler Explorer](#)

The former is called with the desired value `i`, the new value `i`, and the memory order to consider if there is a match and if there is no match. Typically, if there is a match and we want to replace the value, we may use `std::memory_order_release`. If there is no match, we are just reading the value and `std::memory_order_acquire` would suffice. In the latter variant, we pass two pointers.



The difference between the “weak” and “strong” variant is in that the weak variant may return false even if there is a match, in certain cases, but can be much faster in certain architectures. This notably entails ARM architectures (RISC-V and MIPS), where the weak variant will be implemented using the so called load-link/store-conditional pair of instructions (load exclusive register / ldxr and store exclusive register / stxr in ARM version 8). These are much faster than the comparable instructions issuing a barrier (ldaxr/stlxx in ARM version 8). All four ARM instructions utilize only two registers, compared to three registers for CAS proper in Intel architectures (Compare and exchange / cmpxchg since 80486 and cmpxchg8b and cmpxchg16b since Intel Core 2). On recent Intel and AMD processors, cmpxchg is only marginally slower than a non-cached load.

In C++, the *only* synchronization primitive that is guaranteed to be hardware implemented is a particular atomic boolean type, which is known as `std::atomic_flag`. Unlike all specializations of `std::atomic`, it is guaranteed to be lock-free. Prior to C++20, it has been very restricted, because there was no way to check the value of `std::atomic_flag` without setting it. C++20 adds method `test()`.

## 2.3 Synchronization primitives

Synchronization primitives make it possible to synchronize or restrict access of multiple threads to some resources (e.g., global variables, file handles, sockets). You can use them as an interface, without knowing their implementation.

### 2.3.1 Raw synchronization primitives

Lock, Mutex, Semaphore, Atomic, Memory Fence, Condition Variable are synchronization primitives, which make it possible to synchronize or restrict access of multiple threads to some resources.

Lock is a very general term for a synchronization primitive. Mutexes are usually used by one thread only, while semaphores are shared between multiple threads. The *binary semaphore* is the most simple type of a lock, which provides exclusive access for both reading and writing. *counting semaphore* limits the use of a single resource by at most a given number of threads.

A spinlock, the thread simply waits (“spins”) until the lock becomes available. This is efficient if threads are blocked for a short time, because it avoids the overhead of operating system process re-scheduling. It is inefficient if the lock is held for a long time, or if the progress of the thread that is holding the lock depends on preemption of the locked thread. An intentionally simplistic implementation is presented below.

A silly implementation of a spin lock.

```
1 #include <thread>
2 #include <queue>
3 #include <iostream>
4 #include <atomic>
5 #include "atomic4.h"
6
7 int main() {
8     std::queue<int> numbers;
9     SpinLock numbers_lock;
10    int n = 0;
11    for (n = 0; n < 10; ++n) {
12        numbers.push(n);
13    }
14    for (n = 0; n < 2; n++) {
15        std::jthread t{[n, &numbers, &numbers_lock] () {
16            numbers_lock.lock();
17            int val = numbers.front();
18            numbers.pop();
19            numbers_lock.unlock();
20        }};
21    }
22 }
```

[Open in Compiler Explorer](#)

### 2.3.2 Further synchronization features

Fences help order non-atomic and atomic memory accesses, without any associated operations. On Intel architectures (including x86-64), `atomic_thread_fence` do not issue any instructions, except for

`std::atomic_thread_fence(std::memory_order::seq_cst)`.

Barrier provides a thread-coordination mechanism that blocks a group of threads until all threads in that group have reached the barrier. Such a barrier can be used repeatedly to wait until a number of threads have finished their operations.

Latch and is a downward counter, whose initial value is initialized and then threads may block on the latch until the counter is zero. One thread may decrement a latch multiple times, but no thread can increment the latch. Thus, it serves as a single-use barrier.

We will also see synchronized output streams. The synchronized buffer is flushed only when the destructor of the synchronized buffer is called, but provides for guarantees of atomicity for the access. (That is, `endl` and `std::flush` no longer flush!)

## Chapter 3

# The Syntax in C++23

### 3.1 Threads, Tasks, Coroutines

#### 3.1.1 C++11 thread

C++11 had a very basic support for threads, in terms of `std::thread` of header `thread`. The thread starts running once the constructor is called. The object is not CopyConstructible nor CopyAssignable. The challenge in C++11 threads is that one needs to call `join` or `detach` prior to the destructor being called. If neither was called, the program was `std::aborted`. Prior to calling either, one needs to check whether the thread is `joinable()`. At the same time, it was almost impossible to handle exceptions while being able to call `join` correctly. The use of the C++11 `thread` is thus considered harmful and we will present only two short examples.

An example of the use of a C++11 thread.

```
1 #include <iostream>
2 #include <chrono>
3 #include <thread>
4
5 using namespace std::this_thread;
6 using namespace std::chrono_literals;
7
8 void A() {
9     std::cout << "a";
10    sleep_for(5s);
11    std::cout << "A";
12 }
13
14 int main() {
15     std::thread t(A);
16     t.join();
17 }
```

[Open in Compiler Explorer](#)

An example of the use of a C++11 thread.

```
1 #include <iostream>
2 #include <chrono>
3 #include <thread>
4
5 using namespace std::this_thread;
6 using namespace std::chrono_literals;
7
8 void A() {
9     std::cout << "a";
10    sleep_for(5s);
11    std::cout << "A";
12 }
13
14 void B() {
15     std::cout << "b";
16     sleep_for(1s);
17     std::cout << "B";
18 }
19
20 void C() {
21     std::cout << "c";
22     std::thread t(A);
23     t.detach();
24     std::thread u(B);
25     u.join();
26     std::cout << "C";
27 }
28
29 int main() {
30     C();
31     std::thread t(B);
32     t.join();
33     A();
34 }
```

[Open in Compiler Explorer](#)

### 3.1.2 C++20 `jthread`

C++20 adds a new class `jthread` (“joining threads”), which does not require a call to `join` or `detach`. Instead, the destructor waits for completion of the code (“joins”) automatically.

This is an example of the “resource acquisition is initialization” idiom, which is generally one of the best practices in C++. In RAII, the resource allocation is tied to an object’s lifetime and is hence a class invariant. In a constructor, one allocates the resources. In a destructor, one releases the resources. There is no risk of a resource leak.

An example of the use of jthread.

```
1 #include <iostream>
2 #include <thread>
3 #include <vector>
4
5 void Hello();
6
7 int main(int argc, char* argv[]) {
8     std::vector<std::jthread> threads;
9     for (int cnt=0; cnt < 10; cnt++) {
10         threads.push_back(std::jthread(Hello));
11     }
12     return 0;
13 }
14
15 void Hello() {
16     using namespace std::chrono_literals;
17     std::this_thread::sleep_for(2s);
18 }
```

[Open in Compiler Explorer](#)

Notice that the above example would very like result in abnormal program termination, if we changed jthread to thread. (Why?) When we use standard output, it is prudent to wrap it in a syncstream:

An example of the use of jthread.

```
1 #include <iostream>
2 #include <syncstream>
3 #include <thread>
4 #include <vector>
5
6 void Foobar(int cnt);
7
8 int main(int argc, char* argv[]) {
9     std::vector<std::jthread> threads;
10    for (int cnt=0; cnt < 10; cnt++) {
11        threads.push_back(std::jthread(Foobar, cnt));
12    }
13    std::osyncstream(std::cout) << "Main thread" << std::endl;
14    return 0;
15 }
16
17 void Foobar(int cnt) {
18    std::osyncstream(std::cout) << "Thread " << cnt << std::endl;
19 }
```

[Open in Compiler Explorer](#)

Rather commonly, one uses the lambda function to define the thread.  
(This is the `[]()`.)

An example of the use of `jthread`.

```
1 #include <iostream>
2 #include <syncstream>
3 #include <thread>
4 using namespace std::chrono_literals;
5
6 int main() {
7
8     auto t1 = std::jthread([](){
9         std::osyncstream(std::cout) << "Another thread" <<
10         ↪ std::endl;
11         std::this_thread::sleep_for(1s);
12     });
13
14     std::this_thread::sleep_for(2s);
15     std::osyncstream(std::cout) << "Main thread" << std::endl;
16 }
```

[Open in Compiler Explorer](#)

Passing arguments to threads is, nevertheless, very useful. Notably, when we pass the first argument of type `std::stop_token` token, we request the thread to stop its execution by calling `request_stop()` on the `jthread` object:

An example of the use of `jthread`.

```
1 #include <iostream>
2 #include <syncstream>
3 #include <thread>
4 using namespace std::chrono_literals;
5
6 int main() {
7     auto t1 = std::jthread([](std::stop_token token){
8         while (!token.stop_requested()) {
9             std::osyncstream(std::cout) << "A thread";
10            std::this_thread::sleep_for(1s);
11        }
12        std::osyncstream(std::cout) << "Stop requested";
13    });
14
15    std::this_thread::sleep_for(2s);
16
17    std::osyncstream(std::cout) << "Main thread";
18    t1.request_stop();
19 }
```

[Open in Compiler Explorer](#)

More complicated procedures for the stopping of the thread possible. One can define `std::stop_callback` object inside the thread, whose constructor takes the stop token (`std::stop_token`) and a function. The function (in the example below another lambda) gets executed, when the thread is requested to stop via the `std::stop_token`:



An example of the use of `jthread`.

```
1 #include <iostream>
2 #include <syncstream>
3 #include <atomic>
4 #include <thread>
5 using namespace std::chrono_literals;
6
7 int main() {
8
9     auto t = std::jthread([](std::stop_token token) {
10         std::osyncstream(std::cout) << "Thread " <<
11         ↪ std::this_thread::get_id() << std::endl;
12         std::atomic<bool> flag = false;
13         std::stop_callback callback(token, [&flag]{
14             std::osyncstream(std::cout) << "Stop requested" <<
15             ↪ std::endl;
16             flag = true;
17         });
18         while (!flag) {
19             std::this_thread::sleep_for(1s);
20         }
21     });
22     std::osyncstream(std::cout) << "Main thread" << std::endl;
23     std::this_thread::sleep_for(3s);
24     t.request_stop(); // runs all callbacks!
25 }
```

Open in Compiler Explorer [↗](#)

(For a substantially more complex example, see [https://en.cppreference.com/w/cpp/thread/stop\\_callback](https://en.cppreference.com/w/cpp/thread/stop_callback).)

### 3.1.3 Coroutines

In C++23 and subsequent versions, we hope to see some standard syntax for defining coroutines (cf. P2502), such as:

An example of the use of coroutines, which currently does not compile in GCC 12.2.

```
1 #include <coroutine>
2 #include <generator>
3 #include <iostream>
4 #include <syncstream>
5
6 std::generator<int> work() {
7     for (int i = 0; i < 10; i++) {
8         co_yield i;
9     }
10 }
11
12 int main() {
13     for (int i : work()) {
14         std::osyncstream(std::cout) << i << '\n';
15     }
16 }
```

[Open in Compiler Explorer](#)

Already, there are three new keywords:

- **co\_await** **awaiter** suspends computation and block the co-routine until the computation is resumed by another co-routine calling “resume” method of the present coroutine. In the process, it tests whether it is possible to suspend the computation using an awaiter (such as `std::suspend_always{}`;) and, if so, saves all local variables to a heap-allocated handle.
- **co\_yield** yields a value and suspends computation as above, and
- **co\_return** returns a value. (There is no notion of an optional return type in-built.)

Unfortunately, defining the coroutine in C++20 take some more effort. In particular, it requires:

- defining the behaviour of the coroutine, which is known as a **promise** (different from `std::promise`), and requires one returns the type used to access the state of the coroutine on the heap, which is known as the handle,
- defining how to store the state of the coroutine on the heap, using template class `std::coroutine_handle` parametrized by the promise

Clearly, one needs to declare one, define the other, and then return to declare the first one. We will see how to do this later. Optionally, we

can also define an awaiter, which controls suspension and resumption behaviour.

Another difficulty in using coroutines is the fact that the coroutine may live longer than the scope it has been called from. It is hence *not* advisable to pass by reference, except perhaps `std::ref` or `std::cref`. One can either pass by value or pass, e.g., `std::unique_ptr`:

An example of the use of coroutines, which currently does not compile in GCC 12.2.

```
1 // inspired by
  ↳ https://www.incredibuild.com/blog/cpp-coroutines-lets-play-with-them
2
3 #include <coroutine>
4 #include <generator>
5 #include <iostream>
6 #include <syncstream>
7 #include <memory>
8 #include <string>
9
10 std::generator<char> split-by-value(std::string s) {
11     for (char ch : ps) {
12         co_yield ch;
13     }
14 }
15
16 std::generator<char>
  ↳ split-by-uniqueptr(std::unique_ptr<std::string> ps) {
17     for (char ch : *ps) {
18         co_yield ch;
19     }
20 }
21
22 int main() {
23     for (char ch : split-by-value("test")) {
24         std::osyncstream(std::cout) << ch << '\n';
25     }
26     for (char ch :
  ↳ split-by-uniqueptr(std::make_unique<std::string>("west"))
  ↳ {
27         std::osyncstream(std::cout) << ch << '\n';
28     }
29 }
```

[Open in Compiler Explorer](#)

Below, we will see two examples of the use of coroutines that compile with GCC 11 and 12 (enable `-std=c++2b -fcoroutines`). These focus on defining the promise class to be used with template class `std::coroutine_handle<promise>`. The promise class defines the

behaviour of the coroutine by implementing methods:

- `coroutine get_return_object()` is called to initialize the coroutine and create the coroutine handle, which can be the rather formulaic `coroutine_handle::from_promise(*this)`;
- `std::suspend_always initial_suspend()`, suggests whether the coroutine starts right after initialization (`std::suspend_never()`) or upon resumption (`std::suspend_always()`). (Both awaiters are described below.)
- `std::suspend_always final_suspend() noexcept`, which can be rather formulaic `std::suspend_always()`
- `void return_void()` or `void return_value(const auto& value)`, which is called upon reaching the end of the coroutine and upon reaching `co_return`. The latter (`return_value`) often just stores the result locally.
- `void unhandled_exception()`, which can be rather formulaic `std::terminate()`, or can save the exception via `std::current_exception()`.

The promise class is instantiated for each instance of the coroutine, and its methods are called as follows:

A schema of the coroutine, in terms of its calls of the methods of the promise class.

```
1 {
2   co_await promise.initial_suspend();
3   try {
4     ...
5   }
6   catch (...) {
7     promise.unhandled_exception();
8   }
9   // finally
10  co_await promise.final_suspend();
11 }
```

[Open in Compiler Explorer](#)

Once we have a promise class, we can specialize template class `std::coroutine_handle`, which can be seen as the equivalent of a pointer and its method “destroy” as equivalent to a “free”, and use the handle specialized to our own promise class to define a promise class:

An example of the use of coroutines.

```
1 // https://en.cppreference.com/w/cpp/language/coroutines
2
3 #include <coroutine>
4 #include <iostream>
5 #include <syncstream>
6
7 struct promise;
8
9 struct coroutine : std::coroutine_handle<promise> {
10     using promise_type = struct promise;
11 };
12
13 struct promise {
14     coroutine get_return_object() { return
15         ↪ {coroutine::from_promise(*this)}; }
16     std::suspend_always initial_suspend() noexcept { return
17         ↪ {}; }
18     std::suspend_always final_suspend() noexcept { return {};
19         ↪ }
20     void return_void() {}
21     void unhandled_exception() {}
22 };
23
24 int main() {
25     coroutine h = [](int i) -> coroutine {
26         std::osyncstream(std::cout) << i;
27         co_return;
28     }(0);
29     h.resume();
30     h.destroy();
31 }
```

[Open in Compiler Explorer](#)

Sometimes, we also store the promise type in a `promise_type` type member, and disable (= `delete`) copy and move constructors.

## Awaiters

Finally, let us consider awaiters, which can be called when a coroutine is suspended or resumed. Key methods of an awaiter include:

- `await_ready()` is called immediately before suspension of a coroutine. If it returns `true`, the coroutine will not be suspended.
- `await_suspend(handler)` is called immediately after the suspension of the coroutine. The `handler` of type `std::coroutine_handle` can be used to pass the state of the coroutine (e.g., to another

thread).

- `await_resume()` is called when the coroutine is resumed after a successful suspension. If it returns a value, this will be returned by the `co_await` routine.

The two awaiters we have seen so far (`std::suspend_never()` and `std::suspend_always()`) just returned boolean constants in `await_ready()`:

An example of two standard awaiters.

```
1 //
  ⇨ https://github.com/gcc-mirror/gcc/blob/master/libstdc%2B%2B-v3/include/std/coroutine
2
3 // 17.12.5 Trivial awaitables
4 /// [coroutine.trivial.awaitables]
5 struct suspend_always
6 {
7     constexpr bool await_ready() const noexcept { return
  ⇨ false; }
8     constexpr void await_suspend(coroutine_handle<>) const
  ⇨ noexcept {}
9     constexpr void await_resume() const noexcept {}
10 };
11
12 struct suspend_never
13 {
14     constexpr bool await_ready() const noexcept { return true;
  ⇨ }
15     constexpr void await_suspend(coroutine_handle<>) const
  ⇨ noexcept {}
16     constexpr void await_resume() const noexcept {}
17 };
```

[Open in Compiler Explorer](#)

By defining `await_transform()` in the promise type, the compiler will use `co_await promise.await_transform(<expr>)` instead of any call of `co_await <expr>` in the coroutine.

An example defining a Generator (1/2).

```
1 // This code is based on the work of Simon Toth,
2 //
3   ↪ https://itnext.io/c-20-coroutines-complete-guide-7c3fc08db89d
4 // The caller-level type
5 struct Generator {
6     // The coroutine level type
7     struct promise_type {
8         using Handle = std::coroutine_handle<promise_type>;
9
10        Generator get_return_object() {
11            return Generator{Handle::from_promise(*this)};
12        }
13        std::suspend_always initial_suspend() { return {}; }
14        std::suspend_always final_suspend() noexcept { return
15            ↪ {}; }
16        std::suspend_always yield_value(int value) {
17            current_value = value;
18            return {};
19        }
20        void unhandled_exception() { }
21        int current_value;
22    };
23    explicit Generator(promise_type::Handle coro) :
24        ↪ coro_(coro) {}
25    // Make move-only
26    Generator(const Generator&) = delete;
27    Generator& operator=(const Generator&) = delete;
28    Generator(Generator&& t) noexcept : coro_(t.coro_) {
29        ↪ t.coro_ = {}; }
30    Generator& operator=(Generator&& t) noexcept {
31        if (this == &t) return *this;
32        if (coro_) coro_.destroy();
33        coro_ = t.coro_;
34        t.coro_ = {};
35        return *this;
36    }
37    int get_next() {
38        coro_.resume();
39        return coro_.promise().current_value;
40    }
41 private:
42     promise_type::Handle coro_;
43 };
```

Open in Compiler Explorer [↗](#)

An example defining a Generator (2/2).

```
1 // This code is based on the work of Simon Toth,
2 //
   ↳ https://itnext.io/c-20-coroutines-complete-guide-7c3fc08db89d
3
4 #include <coroutine>
5 #include <iostream>
6 #include "coroutines4.h"
7
8 Generator myCoroutine() {
9     int x = 0;
10    while (true) {
11        co_yield x++;
12    }
13 }
14
15 int main() {
16     auto c = myCoroutine();
17     int x = 0;
18     while ((x = c.get_next()) < 10) {
19         std::cout << x << "\n";
20     }
21 }
```

[Open in Compiler Explorer](#)

Eventually, in C++23 and C++26, the support for message-passing architectures based on <https://github.com/lewissbaker/cppcoro> should be available. For a small sample, see a custom implementation [here](#):



An example of the use of coroutines in message-passing architectures (1/3).

```
1 // This code is based on the work of Simon Toth,
2 //
3   ↪ https://itnext.io/c-20-coroutines-complete-guide-7c3fc08db89d
4 class Event {
5 public:
6     Event() = default;
7
8     Event(const Event&) = delete;
9     Event(Event&&) = delete;
10    Event& operator=(const Event&) = delete;
11    Event& operator=(Event&&) = delete;
12
13    class Awaiter;
14    Awaiter operator co_await() const noexcept;
15
16    void notify() noexcept;
17
18 private:
19     friend class Awaiter;
20     mutable std::atomic<void*> suspendedWaiter{nullptr};
21     mutable std::atomic<bool> notified{false};
22
23 };
24
25 class Event::Awaiter {
26 public:
27     Awaiter(const Event& eve): event(eve) {}
28
29     bool await_ready() const;
30     bool await_suspend(std::coroutine_handle<> corHandle)
31         ↪ noexcept;
32     void await_resume() noexcept {}
33
34 private:
35     friend class Event;
36
37     const Event& event;
38     std::coroutine_handle<> coroutineHandle;
39 };
40
41 struct Task {
42     struct promise_type {
43         Task get_return_object() { return {}; }
44         std::suspend_never initial_suspend() { return {}; }
45         std::suspend_never final_suspend() noexcept { return
46             ↪ {}; }
47         void return_void() {}
48         void unhandled_exception() {}
49     };
50 }
```

### An example of the use of coroutines in message-passing architectures (2/3)

```
1 // This code is based on the work of Simon Toth,
2 //
   ↪ https://itnext.io/c-20-coroutines-complete-guide-7c3fc08db89d
3
4 bool Event::Awaiter::await_ready() const {
5     // allow at most one waiter
6     if (event.suspendedWaiter.load() != nullptr){
7         throw std::runtime_error("More than one awaiter is not
           ↪ valid");
8     }
9
10    return event.notified; // `false' suspends the coroutine
11 }
12
13 bool Event::Awaiter::await_suspend(std::coroutine_handle<>
   ↪ corHandle) noexcept {
14
15    coroutineHandle = corHandle;
16
17    if (event.notified) return false;
18
19    // store the waiter for later notification
20    event.suspendedWaiter.store(this);
21
22    return true;
23 }
24
25 void Event::notify() noexcept {
26    notified = true;
27
28    // try to load the waiter
29    auto* waiter =
       ↪ static_cast<Awaiter*>(suspendedWaiter.load());
30
31    // check if a waiter is available
32    if (waiter != nullptr) {
33        // resume the coroutine => await_resume
34        waiter->coroutineHandle.resume();
35    }
36 }
37
38 Event::Awaiter Event::operator co_await() const noexcept {
39     return Awaiter{ *this };
40 }
41
42 Task receiver(Event& event) {
43     auto start = std::chrono::high_resolution_clock::now();
44     co_await event;
45     std::cout << "Got the notification! " << std::endl;
46     auto end = std::chrono::high_resolution_clock::now();
47     std::chrono::duration<double> elapsed = end - start;
48     std::cout << "Waited " << elapsed.count() << " seconds."
       ↪ << std::endl;    30
49 }
```

### An example of the use of coroutines in message-passing architectures (3/3)

```
1 // This code is based on the work of Simon Toth,
2 //
   ↳ https://itnext.io/c-20-coroutines-complete-guide-7c3fc08db89d
3
4 #include <coroutine>
5 #include <chrono>
6 #include <iostream>
7 #include <functional>
8 #include <string>
9 #include <stdexcept>
10 #include <atomic>
11 #include <thread>
12
13 #include "coroutines5a.h"
14 #include "coroutines5b.h"
15
16 using namespace std::chrono_literals;
17
18 int main(){
19
20     std::cout << std::endl;
21
22     std::cout << "Notification before waiting" << std::endl;
23     Event event1{};
24     auto senderThread1 = std::thread([&event1]{
25         ↳ event1.notify(); });
26     auto receiverThread1 = std::thread(receiver,
27         ↳ std::ref(event1));
28
29     receiverThread1.join();
30     senderThread1.join();
31
32     std::cout << std::endl;
33
34     std::cout << "Notification after 2 seconds waiting" <<
35         ↳ std::endl;
36     Event event2{};
37     auto receiverThread2 = std::thread(receiver,
38         ↳ std::ref(event2));
39     auto senderThread2 = std::thread([&event2]{
40         std::this_thread::sleep_for(2s);
41         event2.notify();
42     });
43
44     receiverThread2.join();
45     senderThread2.join();
46
47     std::cout << std::endl;
48 }
```

Open in Compiler Explorer [↗](#)

For further nice examples, see also Boost Asio, e.g., [https://www.boost.org/doc/libs/1\\_78\\_0/doc/html/boost\\_asio/example/cpp20/channels/throttling\\_proxy.cpp](https://www.boost.org/doc/libs/1_78_0/doc/html/boost_asio/example/cpp20/channels/throttling_proxy.cpp), as discussed, e.g., at [https://www.youtube.com/watch?app=desktop&v=ZNttI\\_WswMU&ab\\_channel=ACCUConference](https://www.youtube.com/watch?app=desktop&v=ZNttI_WswMU&ab_channel=ACCUConference). For more details, see Simon Toth's Complete guide at <https://itnext.io/c-20-coroutines-complete-guide-7c3fc08db89d>. For the technical specification, see <https://github.com/GorNishanov/coroutines-ts>.

## 3.2 Synchronisation Primitives

### 3.2.1 Atomic Variables

Since C++11, there is an excellent support for atomic variables in header `<atomic>`. The primary template can be instantiated with types that are TriviallyCopyable, CopyConstructible, and CopyAssignable. First, let us consider two simple examples:

An example of the use of atomic variables.

```
1 #include <iostream>
2 #include <atomic>
3 #include <thread>
4
5 std::atomic<int> i(0);
6
7 int main() {
8     auto t1 = std::jthread([](){
9         int j;
10        do { j = i; }
11        while (j == 0);
12        std::cout << j << std::endl;
13    });
14    auto t2 = std::jthread([](){
15        i = 1;
16    });
17    return 0;
18 }
```

[Open in Compiler Explorer](#)

An example of the use of atomic variables.

```
1 #include <iostream>
2 #include <atomic>
3 #include <thread>
4
5 std::atomic<int> i(0);
6
7 int main() {
8     auto t1 = std::jthread([](){
9         int j;
10        do { j = i.load(std::memory_order_relaxed); }
11        while (j == 0);
12        std::cout << j << std::endl;
13    });
14    auto t2 = std::jthread([](){
15        i.store(1, std::memory_order_relaxed);
16    });
17    return 0;
18 }
```

[Open in Compiler Explorer](#)

Finally, let us consider two, more complete examples of a stack:

A stack implemented using atomic variables.

```
1 // based on
  ↪ https://en.cppreference.com/w/cpp/atomic/atomic/compare_exchange
  ↪
2 #include <atomic>
3 #include <stack>
4
5 template<typename T>
6 struct Node {
7     T data;
8     Node* next;
9     Node(const T& data) : data(data), next(nullptr) {}
10 };
11
12 template<typename T> class stack {
13     std::atomic<Node<T>*> head;
14 public:
15     void push(const T& data) {
16         Node<T>* new_node = new Node<T>(data);
17         new_node->next = head.load(std::memory_order_relaxed);
18         while(!head.compare_exchange_weak(new_node->next,
19             ↪ new_node, std::memory_order_release,
20             ↪ std::memory_order_relaxed));
21     }
22 };
23
24 int main() {
25     std::stack<int> s; s.push(1); s.push(2); s.push(3);
26 }
```

Open in Compiler Explorer [↗](#)

We will elaborate upon this later in Chapter 6.

### 3.2.2 Mutexes and Locks

Standard Template Library in header `<mutex>` provides multiple mutexes (of type `BasicLockable` that implement `lock` and `unlock` methods): `mutex`, `recursive_mutex`, `timed_mutex`, `recursive_timed_mutex`, and `unique_lock`.

A good practice for the use of mutexes is to lock them via the RAII idiom. Since C++11, this is available as `std::unique_lock` and `std::lock_guard`, and since C++17 `scoped_lock` in header `<mutex>`. Crucially, using `scoped_lock` provides the ability to lock multiple mutexes at once, avoiding deadlock. One may hence advise to use one or more mutex with a `scoped_lock` on top.

An example of the use of a unique lock.

```
1 #include <mutex>
2 #include <thread>
3 #include <chrono>
4 #include <iostream>
5
6 int main() {
7     using namespace std::chrono_literals;
8     struct Shared {
9         int value;
10        std::mutex mux;
11    };
12    Shared shared{0, {}};
13    auto t1 = std::jthread([&shared]{
14        std::this_thread::sleep_for(1s);
15        for (int i = 0; i < 10; i++) {
16            std::this_thread::yield();
17            {
18                std::unique_lock lock(shared.mux);
19                shared.value += 10;
20            } // mutex unlocks!
21            std::this_thread::sleep_for(1s);
22        }
23    });
24 }
```

[Open in Compiler Explorer](#)

A slightly more verbose example considers:

An example of the use of a unique lock.

```
1 #include <mutex>
2 #include <thread>
3 #include <chrono>
4 #include <iostream>
5
6 int main() {
7     using namespace std::chrono_literals;
8     struct Shared {
9         int value;
10        std::mutex mux;
11    };
12    Shared shared{0, {}};
13    auto t = std::jthread([&shared]{
14        std::this_thread::sleep_for(1s);
15        for (int i = 0; i < 10; i++) {
16            {
17                std::unique_lock lock(shared.mux);
18                shared.value += 1;
19            }
20            std::this_thread::sleep_for(1s);
21        }
22    });
23    auto observer = std::jthread([&shared]{
24        while (true) {
25            {
26                std::unique_lock lock(shared.mux);
27                std::cout << shared.value << std::endl;
28                if (shared.value == 10)
29                    break;
30            }
31            std::this_thread::sleep_for(1s);
32        }
33    });
34 }
```

[Open in Compiler Explorer](#)

### 3.2.3 Barrier

Since C++20, there is support for barriers in header `<barrier>`. The constructor takes an integer value, which is the number of threads that the barrier is expected to block. The key methods are `arrive_and_wait()` and `arrive_and_drop()`. The former functions as one would expect. The latter decrements the initial expected count for all uses by one, as if one thread could never reach the barrier subsequently. This can be very useful in error management:



An example of the use of a barrier.

```
1 #include <barrier>
2 #include <syncstream>
3 #include <iostream>
4 #include <vector>
5 #include <thread>
6 #include <algorithm>
7 #include <random>
8
9 int main() {
10     std::barrier b(5);
11     std::vector<std::jthread> ts;
12     std::generate_n(std::back_inserter(ts), 5, [&b]{
13         return std::jthread([&b]{
14             std::mt19937
15                 ↪ gen(std::hash<std::thread::id>{}(std::this_thread::get_id()));
16             std::bernoulli_distribution d(0.3);
17             int cnt = 1;
18             while (true) {
19                 std::osyncstream(std::cout) <<
20                 ↪ std::this_thread::get_id() << "/" << cnt
21                 ↪ << std::endl;
22                 std::this_thread::yield();
23                 if (d(gen)) {
24                     b.arrive_and_drop();
25                     return;
26                 } else {
27                     b.arrive_and_wait();
28                 }
29                 cnt++;
30             }
31         });
32     });
33 }
```

[Open in Compiler Explorer](#)

More complicated uses of barriers may use the template parameter `CompletionFunction` and have a callable executed whenever the barrier is hit (reaches zero):

An example of the use of a barrier.

```
1 #include <barrier>
2 #include <syncstream>
3 #include <iostream>
4 #include <vector>
5 #include <thread>
6 #include <algorithm>
7 #include <random>
8
9 int main() {
10
11     std::barrier b(4,[id = 1]() mutable noexcept {
12         std::osyncstream(std::cout) << id << " OK" <<
13         ↪ std::endl;
14         id++;
15     });
16     std::vector<std::jthread> runners;
17     std::generate_n(std::back_inserter(runners), 4, [&b]{
18         return std::jthread([&b]{
19             std::osyncstream(std::cout) <<
20             ↪ std::this_thread::get_id() << "/1" <<
21             ↪ std::endl;
22             std::this_thread::yield();
23             b.arrive_and_wait();
24             std::osyncstream(std::cout) <<
25             ↪ std::this_thread::get_id() << "/2" <<
26             ↪ std::endl;
27             std::this_thread::yield();
28             b.arrive_and_wait();
29         });
30     });
31     runners.clear();
32     std::osyncstream(std::cout) << std::endl;
33 }
```

[Open in Compiler Explorer](#)

## 3.3 Algorithms in the Standard Template Library

Since C++17, there is an excellent Parallel Standard Template Library in header `<algorithm>`.

### 3.3.1 For Each

The most useful algorithm from the Standard Template Library (STL) in terms of parallel programming is surely `for_each`:

An example of the use of `for_each`.

```
1 struct Custom {
2     void expensive_operation() {
3         // ...
4     }
5 };
6
7 std::vector<Custom> data(10);
8
9 std::for_each(std::execution::par_unseq,
10             data.begin(), data.end(),
11             [](Custom& el) {
12                 el.expensive_operation();
13             });
```

[Open in Compiler Explorer](#)

As in the serial version of STL, the callable within `for_each` is permitted to change the state of elements, if the underlying range is mutable, but cannot invalidate iterators.

### 3.3.2 Reduce

Similarly useful is the reduce operation (also known as fold, accumulate, aggregate, compress, or inject). Indeed, a whole approach to parallel programming (Map Reduce) is named after the algorithm. There, one applies an associative operation to each piece of data to obtain a partial result, and then (perhaps using a binary-tree reduction, which also may enforce an order which may or may not be desirable) obtains the final result by applying the same associative operation to the partial results. The binary-tree reduction makes it possible to utilize  $O(\log(n))$  rounds of computation on  $n$  processors.

An example of the use of reduce.

```
1 std::vector<int> data{1, 2, 3, 4, 5};
2
3 auto sum = std::reduce(data.begin(), data.end(), 0);
4 // sum == 15
5
6 sum = std::reduce(std::execution::par_unseq,
7   data.begin(), data.end(), 0);
8 // sum == 15
9
10 auto product = std::reduce(data.begin(), data.end(), 1,
11   std::multiplies<>{});
12 // product == 120
13
14 product = std::reduce(std::execution::par_unseq,
15   data.begin(), data.end(), 1, std::multiplies<>{});
16 // product == 120
```

[Open in Compiler Explorer](#)

### 3.3.3 Merge

Finally, in implementing parallel sorting algorithms, we will utilize the parallel merge operation:

An example of the use of a merge.

```
1 std::vector<int> data1{1, 2, 3, 4, 5, 6};
2 std::vector<int> data2{3, 4, 5, 6, 7, 8};
3
4 std::vector<int> out(data1.size()+data2.size(), 0);
5 std::merge(std::execution::par_unseq,
6   data1.begin(), data1.end(),
7   data2.begin(), data2.end(),
8   out.begin());
9 // out == {1, 2, 3, 3, 4, 4, 5, 5, 6, 6, 7, 8}
```

[Open in Compiler Explorer](#)

## Chapter 4

# The Syntax in OpenMP

OpenMP is a specification for parallel programming of shared-memory systems in Fortran, C, and C++. The current version of the specification can be downloaded from <https://www.openmp.org/wp-content/uploads/OpenMP-API-Specification-5-2.pdf>, with many examples at <https://github.com/OpenMP/Examples>.

The specification is built on top of the fork-join model of parallel execution (sériově-paralelní model uspořádání vláken), but generally does not provide any guarantees as to how a particular directive or function is implemented. This also means that running on a different hardware may result in different order of execution of floating-point operations and vastly different performance. Prime implementations include

- libgomp (GOMP) for GCC,
- libomp for clang, and
- liomp5 (IOMP) for ICC/Clang.

There are also two “subprime” implementation available in Microsoft Visual Studio, which has its own OpenMP runtime (`/openmp`) and an experimental support (`/openmp:llvm`) for the clang/llvm OpenMP runtime, in both cases at version 2.0. If you wish to use Microsoft Visual Studio, notice that Intel oneAPI release 2023.0 officially supports Microsoft Visual Studio 2022.

OpenMP specifies preprocessor directives (pragmas) and a library of functions exported via `omp.h`. Ideally, one and the same program written with OpenMP should be possible to run as serial code, or with any number of threads. Many OpenMP programs use only the preprocessor directives, which makes it possible to compile them as serial code even without OpenMP-aware compiler.

There is a long history of the evolution of OpenMP since 1997, from a 50-page long specification of OpenMP 1.0, through 150-page version 3.0 (introducing “tasks”) and 250-page version 4.0 (adding GPGPU support), to much longer versions recently. Traditional implementations of OpenMP have been rather closely built on top of Pthreads, which results in the lack of fine-grained scheduling, memory management, network management, signaling, etc. The lack of fine-grained scheduling notably meant the lack of user-level threads (co-routines) and the lack of queries as to the number of hardware threads utilized by other processes, which often results in high overhead when the number of threads (across all processes!) increases above the number of hardware threads supported. Since OpenMP 5.0, the distinction between threads and tasks has been erased and thread teams are also cast into tasks. There are now also OpenMP implementations over lightweight threads, notably (BOLT is OpenMP over Lightweight Threads, <https://www.bolt-omp.org/>).

## 4.1 Threads, Tasks, Coroutines

### 4.1.1 OpenMP Task Region

Initially, OpenMP application starts with a single thread (initial/master thread). This can spawn parallel regions, typically with multiple threads in a thread pool (“thread team”) in the fork-join manner. This means that in any thread, one can either wait for all the “sibling” threads to finish (“join them”) or spawn a further, nested parallel region.

A key construct in OpenMP is thus `#pragma omp parallel`, which delineates a parallel region and opens a “team of OpenMP threads”, which could be seen as a thread pool of threads or user-level threads. By default, the number of threads is set based on the available hardware threads, but this can be affected by the environment variables (`OMP_NUM_THREADS`) and modifiers of the pragma (`num_threads(2)`) and function calls (`omp_set_num_threads()`).

An example of the use of OpenMP parallel regions.

```
1 #include <iostream>
2 #include <syncstream>
3 #include <omp.h>
4
5 int main() {
6
7     #pragma omp parallel default(none)
8     {
9         int iam = omp_get_thread_num();
10        int nt = omp_get_num_threads();
11        std::osyncstream(std::cout) << iam << "/" << nt <<
        ↪ std::endl;
12    }
13
14    return 0;
15 }
```

[Open in Compiler Explorer](#)

The construct can take a number of modifiers, including:

- **num\_threads**: number of threads to use in the team
- **private**(list of variables): those variables will be private to each thread
- **firstprivate**(list of variables): those variables will be private to each thread, but initially their value will be copied from the master thread using the default copy constructor.
- **lastprivate**(list of variables): those variables will be private to each thread, but at the end, their value will be copied to the master thread using the default copy constructor.
- **shared**(list of variables): these variables will be shared between the master thread and all threads in the new team. It is the programmer's responsibility to keep the variables constructed as long as the parallel region is running
- **default**: values of **private**, **firstprivate**, **shared**, **none** suggest what should be the default behaviour for variables not listed above. The default is **shared**, which is suboptimal from both performance and thread-safety perspective. It is wise to issue **default(none)**.
- **reduction**(reduction-identifier : list) suggests that variables in list should be treated as shared, when they are used by the function **reduction-identifier**, which could also take the special values of **+**, **-**, **\***, **&**, **|**, **'**, **||**, **max**, **min**. The list can include array el-

ements and, when `reduction-identifier` is a static function of a class, accessible data objects of the object.

- `proc_bind`: values of `master` and `close` and `spread` suggest how far from the master thread should be executed the new threads (same core, close in non-uniform architectures, as far as possible in non-uniform architectures).

Whether the nested parallel regions create their own thread teams or use the existing thread teams depends on settings that we can affect through `omp_set_nested`, or environment variables `OMP_NESTED` (which can be true or false) and `OMP_MAX_ACTIVE_LEVELS`, which controls the maximum number of nested active parallel regions. See, for example:

An example of the use of OpenMP parallel regions.

```
1  #include <iostream>
2  #include <syncstream>
3  #include <omp.h>
4
5  int main() {
6
7      omp_set_nested(1);
8      int iam, nt;
9
10     #pragma omp parallel num_threads(2) default(none)
        ↳ private(iam,nt)
11     {
12         iam = omp_get_thread_num();
13         nt = omp_get_num_threads();
14         std::osyncstream(std::cout) << "L1: " << iam << "/" << nt
        ↳ << std::endl;
15
16     #pragma omp parallel num_threads(2) default(none)
        ↳ private(iam,nt)
17     {
18         iam = omp_get_thread_num();
19         nt = omp_get_num_threads();
20         std::osyncstream(std::cout) << "L2 " << iam << "/" <<
        ↳ nt << std::endl;
21     }
22 }
23     return 0;
24 }
```

[Open in Compiler Explorer](#)



### 4.1.2 Threads and their Sizing

As has been mentioned above, ideally one and the same program written with OpenMP should be possible to run as serial code, or with any number of threads. This requires sizing the work in each thread depending on the number of threads, such as in:

An example of the use of OpenMP parallel regions.

```
1 // Based on parallel.1.c in standard OpenMP Examples 5.1
2 //
   ↳ https://www.openmp.org/wp-content/uploads/openmp-examples-5.1.pdf
3
4 #include <omp.h>
5
6 void subdomain(float *x, int  istart, int  ipoints) {
7     int i;
8     for (i = 0; i < ipoints; i++) x[istart+i] = 123;
9 }
10
11 void sub(float *x, int  npoints) {
12     int iam, nt, ipoints, istart;
13
14 #pragma omp parallel default(shared)
   ↳ private(iam,nt,ipoints,istart)
15     {
16         iam = omp_get_thread_num();
17         nt = omp_get_num_threads();
18         ipoints = npoints / nt;    /* size of partition */
19         istart = iam * ipoints;    /* starting array index */
20         if (iam == nt-1)          /* last thread may do more */
21             ipoints = npoints - istart;
22         subdomain(x, istart, ipoints);
23     }
24 }
25
26 int main() {
27     float array[10000];
28     sub(array, 10000);
29     return 0;
30 }
```

Open in Compiler Explorer [↗](#)

### 4.1.3 Sections

An alternative, non-iterative structuring of the code is possible with sections. Each section is a block of code executed by one thread of the current thread team (corresponding to the innermost enclosing parallel region). One can use `private`, `firstprivate`, `lastprivate`,

and **reduction** modifiers, similar to the parallel construct. There is an implied barrier of the sections region, unless eliminated by a **nowait** clause. See, for example:

An example of the use of OpenMP sections.

```
1  #include <iostream>
2  #include <syncstream>
3  #include "omp.h"
4
5  const int thread_count = 2;
6
7  void work(const int& i) {
8      int iam = omp_get_thread_num();
9      std::osyncstream(std::cout) << "Hello from work(" << i << ")
   ↪   by t = " << iam << std::endl;
10 }
11
12 int main(int argc, char* argv[]) {
13     #pragma omp parallel num_threads(thread_count)
14     {
15         #pragma omp sections
16         {
17             #pragma omp section
18             { work(1); }
19             #pragma omp section
20             { work(2); }
21         }
22     }
23     return 0;
24 }
```

Open in Compiler Explorer [↗](#)

An example of the use of OpenMP sections.

```
1 #include <iostream>
2 #include <syncstream>
3 #include "omp.h"
4
5 const int thread_count = 2;
6
7 void work(const int& i) {
8     int iam = omp_get_thread_num();
9     std::osyncstream(std::cout) << "Hello from work(" << i << ")
    ↪ by t = " << iam << std::endl;
10 }
11
12 int main(int argc, char* argv[]) {
13     #pragma omp parallel num_threads(thread_count)
14     {
15         #pragma omp sections
16         {
17             #pragma omp section
18             { work(2); work(3); }
19             #pragma omp section
20             { work(4); }
21         }
22     }
23     return 0;
24 }
```

Open in Compiler Explorer [↗](#)

#### 4.1.4 Tasks

The closest to a coroutine in OpenMP is the concept of a task. While it does not come with a promise of an implementation with a user-level thread library (cf. Argobots, Converse threads, Qthreads, MassiveThreads, Nanos++, Maestro, GnuPth, StackThreads/MP, Prothreads, Capriccio, StateThreads, TiNy-threads, etc), it often is implemented thus.

An example of the use of OpenMP tasks.

```
1 #include <iostream>
2 #include <syncstream>
3 #include "omp.h"
4
5 const int thread_count = 4;
6
7 void work() {
8     int iam = omp_get_thread_num();
9     int nt = omp_get_num_threads();
10    std::osyncstream(std::cout) << "Thread " << iam << " of "
    ↪ << nt << std::endl;
11 }
12
13 int main(int argc, char* argv[]) {
14     #pragma omp parallel num_threads(thread_count)
15     {
16         #pragma omp single
17         {
18             work();
19         }
20         #pragma omp task
21         work();
22     }
23 }
```

[Open in Compiler Explorer](#)

#### 4.1.5 Kernels

There is a nascent support for the use of GPGPUs via `target` construct. While the block following the target construct can be arbitrary, in principle, most GPUs are not able to support arbitrary code. Specifically, there limitations on the use of synchronization primitives and coherence among L1 caches. Ideally, one should combine the offloading of the code to the GPGU (via the target), across multiple partitions (via teams construct) etc. This can get quite non-trivial:

An example of the use of OpenMP targets and teams.

```
1 // Copyright EPCC, The University of Edinburgh,  
  ↳ www.epcc.ed.ac.uk  
2 // Creative Commons AttributionNonCommercial-ShareAlike 4.0  
  ↳ International License.  
3  
4 #pragma omp target teams distribute parallel for\  
5 map(to:B,C), map(tofrom:sum) reduction(+:sum)  
6 for (int i=0; i<N; i++){  
7     sum += B[i] + C[i];  
8 }
```

Open in Compiler Explorer [↗](#)

Another example of the use of OpenMP targets and teams.

```
1 // Based on:  
2 //  
  ↳ https://www.openmp.org/wp-content/uploads/2021-10-20-Webinar-OpenMP-Offload-Programming-Intro.pdf  
3  
4 void saxpy(float a, float* x, float* y, int sz, int  
  ↳ num_blocks) {  
5 #pragma omp target teams distribute parallel for simd \  
6     num_teams(num_blocks) map(to:x[0:sz]) map(tofrom:y[0:sz])  
7     for (int i = 0; i < sz; i++) {  
8         y[i] = a * x[i] + y[i];  
9     }  
10 }
```

Open in Compiler Explorer [↗](#)

We could also use `nowait` as in other OpenMP constructs, but that does not match the “lockstep” execution on the GPGPUs. Since version 5.0, one can also utilize the `teams` construct on its own, optionally with a target (such a GPGPU).

The advantage is that from the “top-level perspective”, the use of the example above remains standard:

Compiling the example above.

```
1 #include <vector>
2 #include "omp.h"
3 #include "target2.h"
4
5 int main() {
6     float a = 3.1415;
7     int sz = 1048576;
8     int num_blocks = 4096;
9     std::vector<float> vx(sz);
10    float* x = vx.data();
11    std::vector<float> vy(sz);
12    float* y = vy.data();
13    saxpy(a, x, y, sz, num_blocks);
14    return 0;
15 }
```

[Open in Compiler Explorer](#)

Compiling OpenMP with offloading is non-trivial. As has been mentioned in the introduction to this chapter, there are multiple implementations of OpenMP. Most likely, you have used `gcc -fopenmp`, so far. Depending on what GPGPUs you wish to target, you may need to switch:

- NVIDIA CUDA Compiler Driver NVCC may be the easiest to use with NVIDIA hardware. For A100 of the RCI cluster, use `cc80: nvcc {mp=gpu -gpu=cc80`
- clang/LLVM compilers use `clang++ -fopenmp -fopenmp-targets=<target triple>` where the triple is documented at [https://llvm.org/doxygen/Triple\\_8h\\_source.html](https://llvm.org/doxygen/Triple_8h_source.html) such as `clang++ -fopenmp -fopenmp-targets=nvptx64 -`
- AMD ROC is built on top of Clang, so starting with `clang -fopenmp -offload-arch` is a good idea. In more detail, `-target x86_64-pc-linux-gnu -fopenmp -fopenmp-`
- Intel Compiler Collection is, since 2021, also built on top of Clang, except still uses the Intel OpenMP library, so you start with `icx -fiopenmp -fopenmp-targets=<target triple>`. Intel has support for `spir64` and `spirv64` architectures, as well as the above.
- NVIDIA on Power with IBM Open XL C/C++ for Linux on Power: `xlc {qsmp {qoffload {tgtarch=sm_70`

We refer to: <https://www.archer.ac.uk/training/course-material/2019/06/AdvOpenMP-manch/L10-OpenMPTargetOffload.pdf> from which we took the first example and to <https://www.openmp.org/wp-content/>

[uploads/2021-10-20-Webinar-OpenMP-Offload-Programming-Introduction.pdf](#), from which we adapted the second example, for further details.

## 4.2 Synchronisation Primitives

### 4.2.1 Atomic Variables

OpenMP has a rich support for atomic variables. At the simplest one may consider:

An example of the use of atomic variables.

```
1 #include <iostream>
2 #include <syncstream>
3 #include "omp.h"
4
5 int main() {
6
7     int i = 0;
8
9     #pragma omp parallel
10    {
11
12        #pragma omp section
13        {
14            int j;
15            #pragma omp atomic
16            do { j = i; } while (j == 0);
17            std::osyncstream(std::cout) << j << std::endl;
18        }
19        #pragma omp section
20        {
21            #pragma omp atomic
22            i = 1;
23        }
24    }
25    return 0;
26 }
```

[Open in Compiler Explorer](#)

More broadly, one can specify:

- operations for which the atomicity is enforced, out of **read**, **write**, **update**, **capture**, out of which update is the default. Capture makes it possible to use operators such as **+=**, e.g., `{v = x; x binop= expr;}`.
- memory order, out of **seq\_cst**, **acq\_rel**, **release**, **acquire**, **relaxed** as discussed in Chapter 2. The default is relaxed-consistency shared memory model.

E.g. `#pragma omp atomic write` or `#pragma omp atomic seq_cst`. Note that you need to use `{}` to create a new block, wherein the operations would be atomic, if you wish to consider some sequence of read and write operations.

### 4.2.2 Reductions

As we have briefly mentioned on page 44, one may consider a reduction, which produces a single value from an associative operations such as addition, multiplication, taking of the minimum, maximum, or custom associative functions. The goal is for each thread to run the reduction with a private copy and then to produce the final result with the same reduction, perhaps in a hierarchical fashion. A simplistic blueprint is provided in the following example:

An example of the use of reductions.

```
1 #include "omp.h"
2
3 int work() { return 0; }
4
5 int main() {
6     int sum = 0;
7     #pragma omp parallel for reduction(+:sum)
8     for (int i = 0; i < 42; i++) sum += work();
9 }
```

[Open in Compiler Explorer](#)

### 4.2.3 Mutexes

OpenMP has only a limited support for mutexes, as it does not support any RAII variant, which makes them hard to use correctly. One can, however, construct one own's RAII variant:



An example of the use of OpenMP mutexes.

```
1  #include <omp.h>
2
3  int count;
4  omp_nest_lock_t countMutex;
5
6  struct RAIIMutexInit {
7      RAIIMutexInit() { omp_init_nest_lock (&countMutex); }
8      ~RAIIMutexInit() { omp_destroy_nest_lock(&countMutex); }
9  } countMutexInit; // notice the scope!
10
11 struct RAIIMutexHold {
12     RAIIMutexHold() { omp_set_nest_lock (&countMutex); }
13     ~RAIIMutexHold() { omp_unset_nest_lock (&countMutex); }
14 };
15
16 void work() {
17     CountMutexHold releaseAtEndOfScope;
18     count++;
19 }
20
21 int main() {
22     work();
23     return 0;
24 }
```

Open in Compiler Explorer [↗](#)

#### 4.2.4 Critical Sections

Instead of mutexes, one can safely use critical sections in OpenMP. A critical section is a block of code which can be executed by at most one thread at one time. This can be used to protect non-trivial non-associative update operations, for which we cannot use reductions:

An example of the use of a critical section.

```
1 #include <iostream>
2 #include <syncstream>
3 #include "omp.h"
4
5 const int thread_count = 2;
6
7 int main() {
8     int a = 42, b = 1;
9     #pragma omp parallel num_threads(thread_count) shared(a)
10     ↪ private(b)
11     {
12         b = omp_get_thread_num()+2;
13         #pragma omp critical
14         {
15             a = a / b;
16             int t = omp_get_thread_num();
17             std::osyncstream(std::cout) << "b = " << b << " in t = "
18             ↪ << t << "\n";
19             std::osyncstream(std::cout) << "b = " << a << " in t = "
20             ↪ << t << "\n";
21         }
22     }
23     std::osyncstream(std::cout) << "a = " << a << "\n";
24     std::osyncstream(std::cout) << "b = " << b << "\n";
25     return 0;
26 }
```

[Open in Compiler Explorer](#)

The use of critical sections does come with a substantial penalty in terms of performance, though.

#### 4.2.5 Barrier

OpenMP provides a straightforward, explicit barrier construct:

An example of the use of a barrier.

```
1 #include <iostream>
2 #include <syncstream>
3 #include "omp.h"
4
5 void work() {
6     std::osyncstream(std::cout) << "a";
7     #pragma omp barrier
8     std::osyncstream(std::cout) << "A";
9 }
10
11 int main() {
12     #pragma omp parallel num_threads(5)
13     work();
14 }
```

[Open in Compiler Explorer](#)

Especially with nested parallel regions, the behaviour can be quite non-trivial. All threads of the current team must complete execution of all tasks bound to the same parallel region prior to continuing past the barrier. What is the current team, however, e.g., whether it is created by the innermost enclosing parallel region, depends on the settings of the nesting:

An example of the use of a barrier, whose study is left as an exercise.

```
1 #include <iostream>
2 #include <syncstream>
3 #include "omp.h"
4
5 void work(int n) { std::osyncstream(std::cout) << n; }
6 void sub3(int n)
7 {
8     work(n);
9 #pragma omp barrier
10    work(n);
11 }
12 void sub2(int k)
13 {
14 #pragma omp parallel shared(k)
15     sub3(k);
16 }
17 void sub1(int n)
18 {
19     int i;
20 #pragma omp parallel private(i) shared(n)
21     {
22 #pragma omp for
23         for (i=0; i<n; i++)
24             sub2(i);
25     }
26 }
27 int main()
28 {
29     sub1(2);
30     sub2(2);
31     sub3(2);
32     return 0;
33 }
```

[Open in Compiler Explorer](#)

Barrier is also implied by the entry and exit in **parallel** regions. There is also an implicit barrier at the end of a **for**, **sections**, **scope**, and **workshare** constructs, unless one explicitly adds a **nowait** clause.

#### 4.2.6 Fences and Flushes

Close to a memory fence, the **flush** construct provides point at which a thread is guaranteed to have a consistent view of memory, similar to a fence.

The **flush** is also implied by the entry and exit in parallel regions, critical regions, operations with locks etc.

## 4.3 Algorithms

While OpenMP does not really implement any algorithms, some of the data-parallel constructs are similar to algorithms in the STL library. Notably:

### 4.3.1 For Each

OpenMP makes it possible to use “for each”:

An example of the use of for each.

```
1 // Please compile with -std=c++2b -fopenmp
2
3 #include <iostream>
4 #include <syncstream>
5 #include "omp.h"
6
7 int main(int argc, char* argv[]) {
8     std::cout << "Hello from the main thread\n";
9
10    #pragma omp parallel for
11    for (int i=0; i<10; i++)
12        std::osyncstream(std::cout) << "Item " << i << " is
        ↳ processed by thread" << omp_get_thread_num() <<
        ↳ std::endl;
13    return 0;
14 }
```

Open in Compiler Explorer [↗](#)

An example of the use of for each.

```
1 #include <iostream>
2 #include <syncstream>
3 #include "omp.h"
4
5 int main() {
6     #pragma omp parallel
7     {
8         std::osyncstream(std::cout) << "Thread" <<
            ↪ omp_get_thread_num() << std::endl;
9
10    #pragma omp for
11        for (int i=0; i<10; i++)
12            std::osyncstream(std::cout) << "For i = " << i << " in t
            ↪ = " << omp_get_thread_num() << std::endl;
13    }
14
15    std::osyncstream(std::cout) << "Main thread\n";
16    return 0;
17 }
```

Open in Compiler Explorer [↗](#)

## 4.4 Exercises

A useful exercise may be to consider a simple ray tracer <https://github.com/ssloy/tinyraytracer/blob/master/tinyraytracer.cpp> with basic OpenMP parallelization and to try to target a GPGPU.

## Chapter 5

# The Syntax in SYCL

SYCL is an open specification for the design of code targeting either CPUs or GPGPUs that utilizes templates and lambda functions. By “open” in “open specification” we mean it is not tied specifically to NVIDIA. The current version is SYCL 2020, released in 2021. If you want to learn more about SYCL, see <https://github.com/codeplaysoftware/syclacademy> for an extensive tutorial.

SYCL is a high-level specification, in that it builds on top of two lower-level standards, while aiming for succinct, modern C++ code, without any language extensions. First, SYCL is built on top of OpenCL and its device and execution models. Second, SYCL shares the intermediate representation of kernels, known as SPIR(-V), with OpenGL and Vulkan. SYCL aims to bring performance portability across various accelerators.

In the text, we omit the `sycl` namespace, where there is no risk of confusion, just as we have omitted the `std` earlier.

### 5.1 Motivation

In the previous chapter, we have seen the use of `target` to implement SAXPY. Now we clearly could have used some native, closed API, such as CUDA, if we had an NVIDIA GPGPU:

### SAXPY in CUDA C++.

```
1 #include <vector>
2 #include <iostream>
3 #include <cuda.h>
4 #define N 1048576
5
6 __global__ void saxpy_kernel(float a, float* x, float* y,
7                               ↪ float* z){
8     int i = blockIdx.x*blockDim.x + threadIdx.x;
9     if (i < n) y[i] = a*x[i] + y[i];
10 }
11
12 int main(){
13     std::vector<float> vx(N);
14     float* x = vx.data();
15     std::vector<float> vy(N);
16     float* y = vy.data();
17     std::vector<float> vz(N);
18     float* z = vz.data();
19     float *dx, *dy, *dz;
20     cudaMalloc(&dx, N*sizeof(float));
21     cudaMalloc(&dy, N*sizeof(float));
22     cudaMalloc(&dz, N*sizeof(float));
23     cudaMemcpy(dx, x, N*sizeof(float), cudaMemcpyHostToDevice);
24     cudaMemcpy(dy, y, N*sizeof(float), cudaMemcpyHostToDevice);
25     int nblocks = (n + 255) / 256;
26     saxpy_kernel<<<nblocks, 256>>>(3.1415, dx, dy, dz);
27     cudaMemcpy(z, dz, N*sizeof(float), cudaMemcpyDeviceToHost);
28     cudaFree(dx);
29     cudaFree(dy);
30     cudaFree(dz);
31     // we do not need free(x), free(y), free(z)
32 }
```

[Open in Compiler Explorer](#)

Admittedly, code in CUDA C++ still looks like C and is not easily portable to any other platform (AMD, Intel, ARM Mali). We could have also used Thrust, the CUDA C++ template library (<https://thrust.github.io/>):



SAXPY in Thrust, the CUDA C++ template library.

```
1 // Based on
  ↳ https://github.com/NVIDIA/thrust/blob/main/examples/saxpy.cu
2
3 #include <thrust/device_vector.h>
4 #include <thrust/transform.h>
5 #include <thrust/copy.h>
6 #include <thrust/fill.h>
7 #include <thrust/functional.h>
8 #include <iostream>
9
10 struct saxpy_functor {
11     const float a;
12     saxpy_functor(float _a) : a(_a) {}
13     __host__ __device__
14     float operator()(const float& x, const float& y) const {
15         return a * x + y;
16     }
17 };
18
19 int main(){
20     thrust::device_vector<float> dx(1048576);
21     thrust::fill(dx.begin(), dx.end(), 1.0);
22     thrust::device_vector<float> dy(1048576);
23     thrust::fill(dx.begin(), dx.end(), 2.0);
24     // Y <- A * X + Y
25     // thrust::transform(dx.begin(), dx.end(), dy.begin(),
  ↳ dy.begin(), 3.1415f * _1 + _2);
26     thrust::transform(dx.begin(), dx.end(), dy.begin(),
  ↳ dy.begin(), saxpy_functor(3.1415));
27     thrust::copy(dy.begin(), dy.end(),
  ↳ std::ostream_iterator<float>(std::cout, "\n"));
28 }
```

[Open in Compiler Explorer](#)

The code is beautiful, but we it is not easily portable to another platform (AMD, Intel, ARM Mali). There may be several good reasons not to be tied specifically to NVIDIA:

- In consumer PC market, the biggest vendor of PC GPU is Intel with 71 percent of the market by unit count as of Q4 2022. NVIDIA has a market share of 17 percent and AMD 12 percent. In discrete PC GPU market, NVIDIA occupies 80-90% of the market, according to various estimates. Note that this market is rapidly declining.
- In mobile GPGPUs, ARM Mali is estimated to have close to 40% market share.
- In gaming consoles, both Xbox and PlayStation use accelerators

by AMD.

- In high-performance computing, NVIDIA is dominant, but the top 1 system (Frontier) currently uses AMD Instinct™ 250X accelerators, and number of other systems use ARM architectures.

In this chapter, we hence consider SYCL, where the same example would resemble:

SAXPY in SYCL, the open specification.

```
1 // Based on https://github.com/jeffhammond/dpcpp-tutorial
2
3 #include <iostream>
4 #include "CL/sycl.hpp"
5
6 class saxpy3;
7
8 int main(int argc, char * argv[]) {
9     std::vector<float> vx(1048576, 1.0);
10    std::vector<float> vy(1048576, 2.01);
11    std::vector<float> vz(1048576, 0.0);
12    sycl::queue q(sycl::default_selector{});
13    try {
14        const float A(aval);
15        sycl::buffer<float,1> dx { vx.data(),
16        ↪ sycl::range<1>(vx.size()) };
17        sycl::buffer<float,1> dy { vy.data(),
18        ↪ sycl::range<1>(vy.size()) };
19        sycl::buffer<float,1> dz { vz.data(),
20        ↪ sycl::range<1>(vz.size()) };
21
22        q.submit([&](sycl::handler& h) {
23            sycl::accessor x(dx, h, sycl::read_only);
24            sycl::accessor y(dy, h, sycl::read_only);
25            sycl::accessor z(dz, h, sycl::read_write);
26            h.parallel_for<class saxpy3>( sycl::range<1>{length},
27            ↪ [=] (sycl::id<1> it) {
28                const size_t i = it[0];
29                z[i] += 3.1415 * x[i] + y[i];
30            });
31        });
32        q.wait();
33    }
34    catch (sycl::exception & e) {
35        std::cout << e.what() << std::endl;
36        return 1;
37    }
38    return 0;
39 }
```

Open in Compiler Explorer

The code is also clearly modern C++, and can target accelerators by a number of vendors (NVIDIA, AMD, Intel, ARM), as well as FPGA.

## 5.2 More concepts

To make as efficient use of GPGPUs as possible, SYCL utilizes a number of concepts that we have not seen earlier. This builds upon the discussion of the GPGPU architecture in Section 2.2.2, page 9.

### 5.2.1 Device selector

Selectors are used to pick device to run on. In header `<device_selector.h>`, there is an abstract class `device_selector` with numerous implementations such as `gpu_selector`, `host_selector`, `opencl_selector`, and `default_selector`.

One can also implement its own subclasses that specify to the runtime how to perform device selection. For example, it may query the amount of memory on the GPGPU and if it is sufficient, use GPGPU. If it were not sufficient, it could use CPU as a fallback.

In a `device_selector`, one overrides `int operator()(const sycl::device& dev) const` override and returns an integer for the priority. The higher integer, the higher priority.

### 5.2.2 Queues

A queue, `queue`, is an abstraction of a device, through which we orchestrate work on the device. In a constructor of a queue, we pass the device, which cannot be changed later, but one can create further queues for the same device. A key method is `queue.submit`, which passes a “command group function object” for asynchronous execution.

At the simplest, a command group” is a callable (a named type, a lambda function, or `std::function`), which receives a “command group handler” from the SYCL as an argument of `operator()` so as to access the API. A command group submission to the queue is atomic.

A command group objects may also combine the callable and a set of requirements (edges of a task-graph), but if declare the accessors for memory access correctly (see below), the task-graph will be constructed based on the accessors.

Asynchronous execution also means an undefined order of execution, unless we use `wait` or suggest the dependencies between the “actions”

in the form of a task graph. We can also declare the queue to be in-order, similar to sorted in OpenMP: `queue q{property::queue::in_order()};`.

### 5.2.3 Work items, Work groups, and Kernels

Within an action submitted to the queue, we execute kernels. Kernels are callables

- receiving an index to the run of the kernel as `auto idx` or `id<1> idx` or similar.
- returning nothing; with `void` return type
- which cannot allocate memory dynamically which cannot use certain other features (e.g., RTTI).

Within the `single_task` function method of the “command group handler” API, we pass a C++ function object as a parameter and have it executed once. Kernel can also be a class that overloads operator `()`, e.g. `void operator()(id<1> idx)`.

Most often, we want the kernel executed many times, in a data-parallel fashion. In the so-called nd-ranges (“kernel grid” on NVIDIA), we partition the index-set of data hierarchically first into global ranges, and then into local ranges. The local range corresponds to a work-group and each element corresponds to a work item (= single run of a kernel). Ideally, all work items within a work group are executed in lock step (i.e., the same hardware instruction in all work items at the same time). The work-group local memory can often be accessed very efficiently, via `local_accessor`, and can be used to coordinate multiple work items (= single runs) within a work group.

The threads of one work group (“thread block” on NVIDIA) are sent to one Streaming Multiprocessors (SM), but one SM can execute threads by multiple work group in its multiple processing blocks. At most one work group per processing block.

To summarize, each work item can access:

- private memory
- work-group local memory
- global memory accessible to all work items within an nd-range, but whose access can be very expensive, as it involves copying data across PCIe bus
- constant memory, which is a part of the global memory, but which can be very cheap to access.

### 5.2.4 Asynchronous errors

The SYCL implementation may throw “synchronous errors” (one at a time). In contrast, asynchronous errors are produced by a command group or a kernel (with many kernels running at any point). By default, asynchronous errors are not propagated to the host. One can, however, define an error handler and pass it to a queue `queue q(default_selector{}, exception_handler)`. The error handler receives an `exception_list`, wherein one can iterate over `std::exception_ptr`.

See <https://www.codingame.com/playgrounds/48226/introduction-to-sycl/error-handling> for a great tutorial with code that is editable, compilable, and runnable online. Let us simplify their main example somewhat:

#### Error handling in SYCL.

```
1 // based on error_handling.cpp of "Introduction to SYCL",
2 //
   ↳ https://www.codingame.com/playgrounds/48226/introduction-to-sycl
3
4 #include <iostream>
5 #include <sycl/sycl.hpp>
6
7 using namespace sycl;
8
9 class exception1;
10
11 int main(int, char**) {
12     auto exception_handler = [] (exception_list exceptions) {
13         for (std::exception_ptr const& e : exceptions) {
14             try {
15                 std::rethrow_exception(e);
16             } catch(exception const& e) {
17                 std::cout << "Caught asynchronous SYCL exception:\n"
18                     << e.what() << std::endl;
19             }
20         }
21     };
22
23     queue q(default_selector{}, exception_handler);
24     // actual use of the q
25
26     try {
27         q.wait_and_throw();
28     } catch(exception const& e) {
29         std::cout << "Caught synchronous SYCL exception:\n"
30             << e.what() << std::endl;
31     }
32     return 0;
33 }
```

[Open in Compiler Explorer](#)

### 5.2.5 Unified shared memory

At the cost of some latency, one can use a unified shared memory across both the host and the device, wherein one uses the same pointer on both the host and the device. This requires:

```
template <typename T> T* malloc_shared(size_t count,
```

and the corresponding

```
void free(void* ptr, sycl::queue& syclQueue)
```

See the example below:

An example of the use of unified shared memory.

```
1 auto shared = malloc_shared<double>(42, q);
2
3 q.submit([&](handler& cgh){
4     cgh.parallel_for(range{42}, [=](id<1> tid){
5         shared[tid] = 0.0;
6     });
7 });
```

[Open in Compiler Explorer](#)

### 5.2.6 Buffers and accessors

A buffer is a constrained view of a 1-, 2-, or 3-dimensional array. The constraints specify how it can be accessed on the host, the device or both. A buffer is constructed with a pre-allocated, trivially copyable C++ objects (e.g., STL container). Within the contract for the use of the buffer, one promises not to amend the memory used to initialise the buffer during the lifetime of the buffer. Buffer promises to update the memory in the host upon destruction, in RAII spirit.

In the case of one-dimensional arrays, one can call the constructor with an iterator: `template <typename InputIterator> buffer(InputIterator first, InputIterator last)`

Once in a kernel, an **accessor** specifies constraints on the use of a buffer therein. Two key choices are:

- access mode: `read`, `write`, and `read_write`, where write access mode also implicitly defines dependencies between tasks
- access target: `global_memory` suggests that the data resides in the global memory space of the device. Other options are device specific. `no_init` suggests that the initial data can be discarded (not moved to the device).

See the example below:

An example of the use of buffers and accessors.

```
1  buffer<double> A{range{42}};
2
3  q.submit([&](handler& cgh){
4      accessor aA{A, cgh};
5      cgh.parallel_for(range{42}, [=](id<1> & idx){
6          aA[idx] = 0.0;
7      })
8  });
9
10 host_accessor result{A};
11 for (int i = 0; i < N; i++) {
12     assert(result[i] == 0);
13 }
```

[Open in Compiler Explorer](#)

### 5.2.7 Barrier

Depending on the details of the use of a barrier, one may wish to use `sycl::queue::wait()` and `sycl::queue::wait_and_throw()`, or `item::barrier(access::fence_space)` within a kernel.

## 5.3 More complex examples

First, let us consider a complete, working example:



An example of vector addition in SYCL.

```
1  #include <iostream>
2  #include <vector>
3  #include <CL/sycl.hpp>
4
5  using namespace std;
6  using namespace cl::sycl;
7  class buffer3;
8
9  int main(int, char**) {
10
11     std::vector<float> a { 2.0, 3.0, 7.0, 4.0 };
12     std::vector<float> b { 4.0, 6.0, 1.0, 3.0 };
13     std::vector<float> c { 0.0, 0.0, 0.0, 0.0 };
14
15     default_selector device_selector;
16
17     queue q(device_selector);
18
19     std::cout << "Running on "
20               <<
21               ↪ queue.get_device().get_info<info::device::name>()
22               << "\n";
23     {
24         buffer bufA(a);
25         buffer bufB(b);
26         buffer bufC(c);
27
28         q.submit([&] (handler& cgh) {
29
30             auto accA = accessor(bufA, cgh, read_only);
31             auto accB = accessor(bufB, cgh, read_only);
32             auto accC = accessor(bufC, cgh, write_only);
33             cgh.single_task<buffer3>(bufC.get_range(), [=](id<1>
34               ↪ i) {
35                 accC[i] = accA[i] + accB[i];
36             });
37         });
38         q.wait_and_throw();
39     }
40     return 0;
41 }
```

[Open in Compiler Explorer](#)

Next, let us consider an example that uses work groups and local memory in an attempt to utilize more of the performance available in the GPGPU.

An example of the use of work groups and local memory.

```
1 // based on memory4.cpp of "Introduction to SYCL",
2 //
  ↪ https://www.codingame.com/playgrounds/48226/introduction-to-sycl
3
4 #include <array>
5 #include <cstdlib>
6 #include <iostream>
7 #include <random>
8 #include <cassert>
9
10 #include <CL/sycl.hpp>
11
12 class reduction_kernel;
13 namespace sycl = cl::sycl;
14
15 int main(int, char**) {
16     std::array<int32_t, 16> arr;
17
18     std::mt19937 mt_engine(std::random_device{}());
19     std::uniform_int_distribution<int32_t> idist(0, 10);
20
21     std::cout << "Data: ";
22     for (auto& el : arr) {
23         el = idist(mt_engine);
24         std::cout << el << " ";
25     }
26     std::cout << std::endl;
27
28     sycl::buffer<int32_t, 1> buf(arr.data(),
  ↪ sycl::range<1>(arr.size()));
29
30     sycl::device device =
  ↪ sycl::default_selector{}.select_device();
31
32     sycl::queue queue(device, [] (sycl::exception_list el) {
33         for (auto ex : el) { std::rethrow_exception(ex); }
34     });
35
36     size_t wgroup_size = 32;
37     auto part_size = wgroup_size * 2;
38
39     auto has_local_mem = device.is_host()
40         ||
  ↪ (device.get_info<sycl::info::device::local_mem_type>()
41         != sycl::info::local_mem_type::none);
42     auto local_mem_size =
  ↪ device.get_info<sycl::info::device::local_mem_size>();
43     if (!has_local_mem
44         || local_mem_size < (wgroup_size * sizeof(int32_t)))
45     {
46         throw "Device doesn't have enough local memory!";
47     }
48
49     #include "buffer4.h"
50
51     auto acc = buf.get_access<sycl::access::mode::read>();
52     std::cout << "Sum: " << acc[0] << std::endl;
53
54     return 0;
55 }
```

## 5.4 Building code

Compiling code with SYCL should require passing “-fsycl” to a compiler that supports the SYCL 2020 specification, but it often turns out to be substantially more complicated than with standard C++23 or OpenMP.

There are many implementations of the SYCL specification (incl. OpenSYCL, Intel DPC++, CodePlay ComputeCPP) with minor differences. One way of getting around the complexity is to install the OpenSYCL compiler: <https://github.com/OpenSYCL/OpenSYCL> (formerly known as hipSYCL).

Independent of the compiler, to build code with SYCL that targets a GPGPU, you need to link with the appropriate libraries, which depends on what is your target. If you target, e.g., common OpenCL variants:

Linking SYCL against OpenCL.

```
1 set(CMAKE_CXX_FLAGS "${CMAKE_CXX_FLAGS} -O3 -fsycl
   ↪ -std=c++17")
2 set(CMAKE_EXE_LINKER_FLAGS "${CMAKE_EXE_LINKER_FLAGS} -lOpenCL
   ↪ -lsycl")
```

Open in Compiler Explorer

If you target MKL:

Linking SYCL against MKL.

```
1 set(CMAKE_CXX_FLAGS "${CMAKE_CXX_FLAGS} -O3 -fsycl
   ↪ -std=c++17")
2 set(CMAKE_EXE_LINKER_FLAGS " -fsycl -lmkl_sycl
   ↪ -lmkl_intel_ilp64 -lmkl_sequential -lmkl_core")
```

Open in Compiler Explorer

If you target CUDA (or NVIDIA OpenCL), you may need to set up paths to CUDA, in addition to the use of:

```
-fsycl -fsycl-targets=nvptx64-nvidia-cuda
```

Similarly for AMD:

Linking SYCL against AMD.

```
1 -fsycl -fsycl-targets=amdgc-n-amd-amdhsa -Xsycl-target-backend  
   ↪ {offload-arch=gfx906
```

[Open in Compiler Explorer](#)

For details of the use of SYCL with CUDA, see <https://github.com/codeplaysoftware/SYCL-For-CUDA-Examples>

# Bibliography

- [1] M. Pöter and J. L. Träff, “Memory models for C/C++ programmers,” *CoRR*, vol. abs/1803.04432, 2018. [Online]. Available: <http://arxiv.org/abs/1803.04432>
- [2] J. Alglave, “A formal hierarchy of weak memory models,” *Formal Methods in System Design*, vol. 41, no. 2, pp. 178–210, 2012. [Online]. Available: <https://doi.org/10.1007/s10703-012-0161-5>

