

DS 7347

# High-Performance Computing (HPC) and Data Science

## Session 26

---

Robert Kalescky

Adjunct Professor of Data Science

HPC Research Scientist

July 26, 2022

Research and Data Sciences Services

Office of Information Technology

Center for Research Computing

Southern Methodist University



Session Questions

Parallelization Limits

C++ Parallel Algorithms

Parallel Python

Assignments and Project

## Session Questions

---



## Tuesday

Generally, what's the relationship between memory speed and size? Describe the hierarchy.

## Thursday

How would you parallelize painting a room and building a car? What is the maximum amount of work that can be done in parallel?

## Parallelization Limits

---



## Concurrency

- Structure to handle multiple tasks
- Tasks can run serially or in parallel

## Parallelization

- Multiple tasks running at the same time
- Parallelization is a special case of concurrency



**Moore's Law** Performance doubling every 18 months

**Amdahl's Law** There's an upper limit to adding more resources

**Gustafson's Law** There's an upper limit to adding more resources

**Brooks's Law** There's an upper limit to adding more human resources

**Law of Diminishing Returns** There's an ideal set of resources for a problem

# C++ Parallel Algorithms

---





- C++17 introduces high-level parallel versions of many algorithms found in the Standard Template Library
- Several new algorithms were included specifically to aid in using the new parallel algorithms, *e.g.* `std::reduce` and `std::transform_reduce`
- Each parallel algorithm has an execution policy that defines how the algorithm is parallelized



The C++17 standard defines three execution policies:

**`std::execution::seq`** Sequential execution, *i.e.* no parallelism

**`std::execution::par`** Parallel execution via one or more threads

**`std::execution::par_unseq`** Parallel execution via one or more threads  
with each thread possibly vectorized



When using an execution policy other than `std::execution::seq`:

- The compiler **may** execute the algorithm in parallel when possible and advantageous (some conforming C++17 implementations may ignore the parallelization flag and run via `std::execution::seq`)
- It is up to **you** to make sure that the algorithm and data are safe to be run in parallel

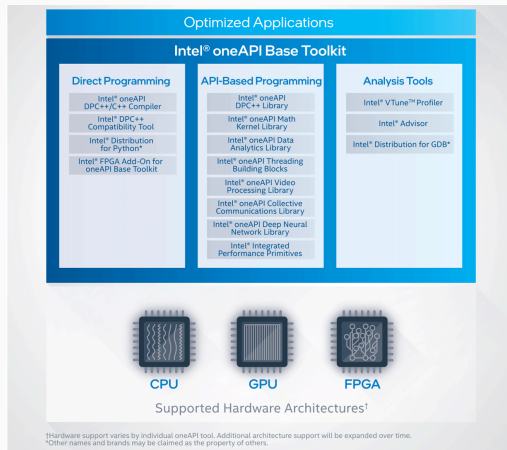


Figure 1: Components of the oneAPI 2021 base installation.



- Standards-based C++ compiler based on the open source LLVM compiler infrastructure
  - Full support up through C++17
  - Initial support for C++20
- SYCL compiler
- DPC++ compiler (SYCL with Intel extensions)
- Partial support for OpenMP 4.5 and 5.0 offloading

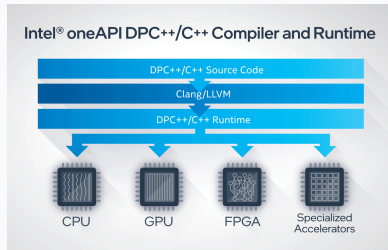


Figure 2: oneAPI DPC++ and C++ compiler and runtime.



- Productivity APIs for heterogeneous computing
- Optimized standards-based and familiar APIs:
  - C++ STL
  - Parallel STL (PSTL)
  - Boost.Compute
  - Standard SYCL
- Custom iterators for parallel algorithms
- Use device and host containers to target GPUs and FPGAs or run your code across multi-node CPUs



- Simplifies the work of adding parallelism to complex applications
- Runtime library automatically maps logical parallelism onto threads, making more efficient use of node resources
- Scalable data-parallel programming
  - Multiple threads to work on different parts of a collection
  - Scales well to larger numbers of processors by dividing the collection into smaller pieces
  - Program performance increases as processors and accelerators are added

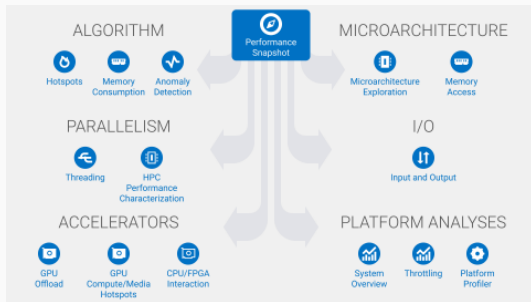


- Design and analysis tool for achieving high application performance
- Helps identify efficient threading, vectorization, and memory use, and GPU offload
- Supports C, C++, Fortran, DPC++, OpenMP, and Python.
- Intel Advisor features:
  - Offload Advisor
  - Automated Roofline Analysis
  - Vectorization Advisor
  - Threading Advisor
  - Flow Graph Analyzer





- Application performance profiler for CPU, GPU, and FPGA
- Supports DPC++, C, C++, Fortran, OpenCL, Python, assembly, or any combination
- Get coarse-grained system data for an extended period or detailed results mapped to source code.
- Low system overhead



**Figure 3:** Intel VTune Profiler overview.



- Helps find memory errors and nondeterministic threading errors and problems
- Dynamic memory and threading error debugger for C, C++, and Fortran applications



- Profiles and analyzes MPI applications
- Discover temporal dependencies and bottlenecks
- Check the correctness of your application
- Locate potential programming errors, buffer overlaps, and deadlocks
- Visualize and understand parallel application behavior
- Evaluate profiling statistics and load balancing
- Analyze performance of subroutines or code blocks
- Learn about communication patterns, parameters, and performance data
- Identify communication hot spots
- Decrease time to solution and increase application efficiency

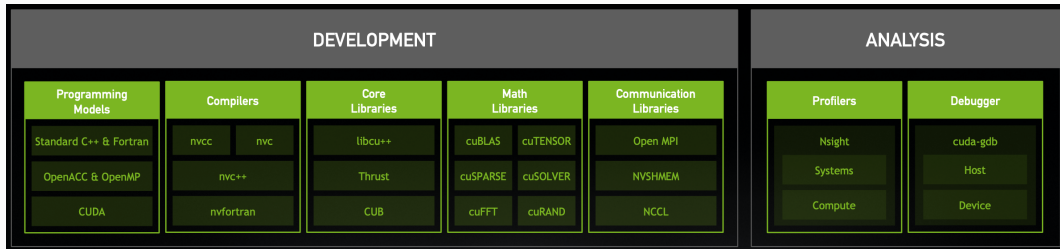


Figure 4: Components of the NVIDIA HPC SDK.



Standards-based compilers based on the open-source LLVM compiler infrastructure.

**nvc** C11 compiler for NVIDIA GPUs and AMD, Intel, OpenPOWER, and Arm CPUs

**nvc++** C++17 compiler for NVIDIA GPUs and AMD, Intel, OpenPOWER, and Arm CPUs

**nvfortran** Fortran 2003 compiler with many Fortran 2008 features implemented for NVIDIA GPUs and AMD, Intel, OpenPOWER, and Arm CPUs

**nvcc** CUDA C and CUDA C++ compiler driver for NVIDIA GPUs



**C++ -stdpar** C++ 17 Parallel Algorithms introduce parallel and vector concurrency through execution policies

**OpenACC** OpenACC directives for parallelization on CPUs and NVIDIA GPUs

**OpenMP** OpenMP directives for parallelization on CPUs and NVIDIA GPUs

**CUDA** C, C++, and Fortran variants of CUDA for parallelization on NVIDIA GPUs

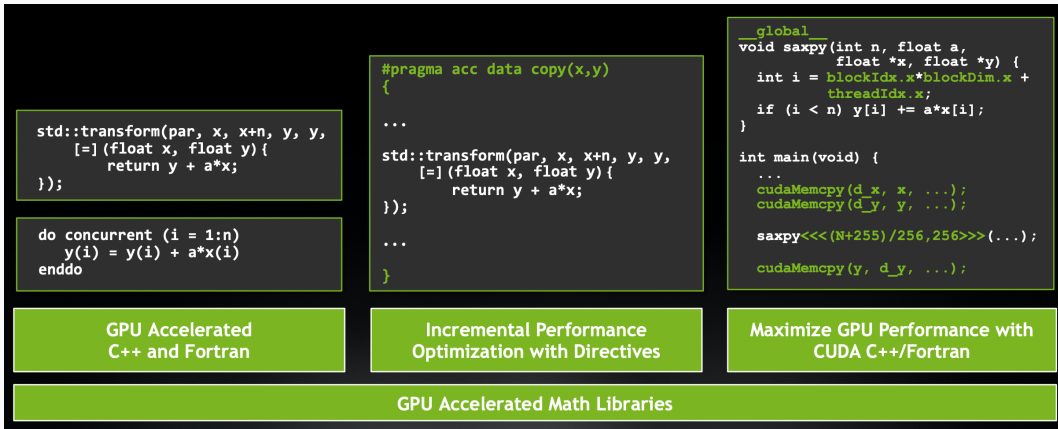


Figure 5: Parallel Programming Models

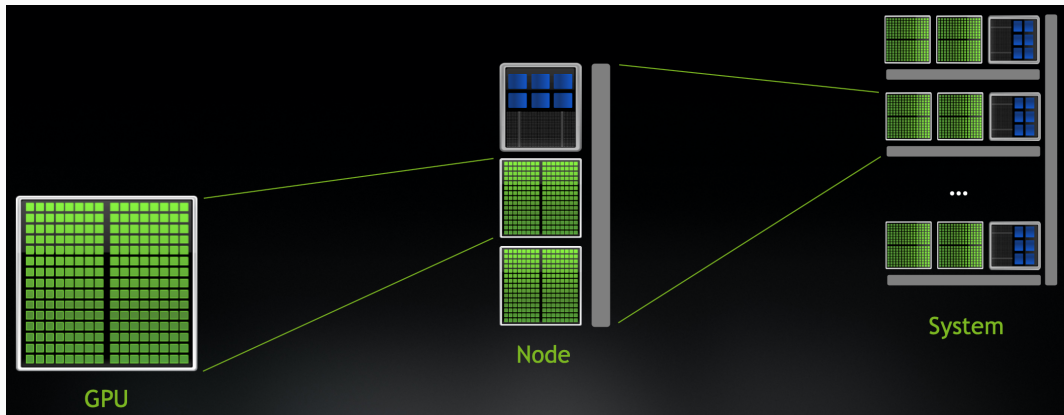


Figure 6: Parallel Programming Models





| C++17   | C++20  | C++23 and Beyond   |
|---|--|--|
| <p><b>Parallel Algorithms</b></p> <ul style="list-style-type: none"><li>➤ In <b>NVC++ 20.5</b></li><li>➤ Parallel and vector concurrency</li></ul> <p><b>Forward Progress Guarantees</b></p> <ul style="list-style-type: none"><li>➤ Extend the C++ execution model for accelerators.</li></ul> <p><b>Memory Model Clarifications</b></p> <ul style="list-style-type: none"><li>➤ Extend the C++ memory model for accelerators.</li></ul> | <p><b>Scalable Synchronization Library</b></p> <ul style="list-style-type: none"><li>➤ Express thread synchronization that is portable and scalable across CPUs and accelerators.</li><li>➤ In <b>libc++ in CUDA 10.2:</b><ul style="list-style-type: none"><li>➤ <code>std::atomic&lt;T&gt;</code></li></ul></li><li>➤ In <b>libc++ in CUDA 11.0:</b><ul style="list-style-type: none"><li>➤ <code>std::barrier</code></li><li>➤ <code>std::counting_semaphore</code></li><li>➤ <code>std::atomic&lt;T&gt;::wait/notify_*</code></li></ul></li><li>➤ In <b>libc++ in the future:</b><ul style="list-style-type: none"><li>➤ <code>std::atomic_ref&lt;T&gt;</code></li></ul></li></ul> | <p><b>Executors</b></p> <ul style="list-style-type: none"><li>➤ Simplify launching and managing parallel work across CPUs and accelerators.</li></ul> <p><b><code>std::mdspan/mdarray</code></b></p> <ul style="list-style-type: none"><li>➤ HPC-oriented multi-dimensional array abstractions.</li></ul> <p><b>Linear Algebra</b></p> <ul style="list-style-type: none"><li>➤ Standard C++ API to vendor BLAS libraries.</li></ul> <p><b>Extended Floating Point Types</b></p> <ul style="list-style-type: none"><li>➤ First-class support for formats new and old: <code>std::float16_t/float64_t</code></li></ul> |

Figure 7: Communications Libraries



```
std::sort(std::execution::par, c.begin(), c.end());  
  
std::unique(std::execution::par, c.begin(), c.end());
```

- Introduced in C++17
- Parallel and vector concurrency via execution policies

```
std::execution::par, std::execution::par_seq, std::execution::seq
```
- Several new algorithms in C++17 including
  - `std::for_each_n(POLICY, first, size, func)`
- Insert `std::execution::par` as first parameter when calling algorithms
- **NVC++ 20.5**: automatic GPU acceleration of C++17 parallel algorithms
  - Leverages CUDA Unified Memory

Figure 8: Communications Libraries



**CUDA-GDB** Debugging CUDA applications.

**Nsight Compute** Interactive kernel profiler for CUDA applications

**Nsight System** System-wide performance analysis tool designed to visualize application algorithms

**Compute Sanitizer** Functional correctness checking tools suite

**NVTX** API for annotating application events, code ranges, and resources for use with Nsight



- C++ STL Algorithms Documentation
- NVIDIA HPC SDK Documentation
- C++ Concurrency In Action, 2<sup>nd</sup> Edition
- CppCon 2018: Bryce Adelstein Lelbach “The C++ Execution Model”
- C++17 Parallel Algorithms on NVIDIA GPUs with PGI C++



- [The oneAPI Specification](#)
- [Data Parallel C++ \(Free eBook\)](#)
- [Intel oneAPI Base Toolkit](#)
- [Intel oneAPI HPC Toolkit](#)
- [Intel oneAPI Toolkit Samples](#)

# Parallel Python

---



- `threading`
- `multiprocessing`
- `asyncio`
- `queue`
- `concurrent.futures`
- Coroutines



- Coroutines (`asyncio`)
  - Minimal overhead
  - Subject to the GIL
- `threading`
  - Non-trivial overhead
  - Subject to the GIL
- `multiprocessing`
  - Substantial overhead





- Only one thread in a Python process can execute Python code at a time
- Calling certain libraries and functions release the GIL, *.e.g.* NumPy



```
1  from threading import Thread
2
3  def fib(n):
4      if n <= 2:
5          return 1
6      elif n == 0:
7          return 0
8      elif n < 0:
9          raise Exception('fib(n) is undefined for n < 0')
10     return fib(n - 1) + fib(n - 2)
11
12 if __name__ == '__main__':
13     import argparse
14
15     parser = argparse.ArgumentParser()
16     parser.add_argument('-n', type=int, default=1)
17     parser.add_argument('number', type=int, nargs='?', default=34)
18     args = parser.parse_args()
19
20     assert args.n >= 1, 'The number of threads has to be > 1'
21     for i in range(args.n):
22         t = Thread(target=fib, args=(args.number, ))
23         t.start()
```



```
1  import concurrent.futures as cf
2
3  def fib(n):
4      if n <= 2:
5          return 1
6      elif n == 0:
7          return 0
8      elif n < 0:
9          raise Exception('fib(n) is undefined for n < 0')
10     return fib(n - 1) + fib(n - 2)
11
12 if __name__ == '__main__':
13     import argparse
14
15     parser = argparse.ArgumentParser()
16     parser.add_argument('-n', type=int, default=1)
17     parser.add_argument('number', type=int, nargs='?', default=34)
18     args = parser.parse_args()
19
20     assert args.n >= 1, 'The number of threads has to be > 1'
21     with cf.ThreadPoolExecutor(max_workers=args.n) as pool:
22         results = pool.map(fib, [args.number] * args.n)
```



```
1  import concurrent.futures as cf
2
3  def fib(n):
4      if n <= 2:
5          return 1
6      elif n == 0:
7          return 0
8      elif n < 0:
9          raise Exception('fib(n) is undefined for n < 0')
10     return fib(n - 1) + fib(n - 2)
11
12 if __name__ == '__main__':
13     import argparse
14
15     parser = argparse.ArgumentParser()
16     parser.add_argument('-n', type=int, default=1)
17     parser.add_argument('number', type=int, nargs='?', default=34)
18     args = parser.parse_args()
19
20     assert args.n >= 1, 'The number of threads has to be > 1'
21     with cf.ProcessPoolExecutor(max_workers=args.n) as pool:
22         results = pool.map(fib, [args.number] * args.n)
```



```
1  import multiprocessing as mp
2
3  def fib(n):
4      if n <= 2:
5          return 1
6      elif n == 0:
7          return 0
8      elif n < 0:
9          raise Exception('fib(n) is undefined for n < 0')
10     return fib(n - 1) + fib(n - 2)
11
12 def worker(inq, outq):
13     while True:
14         data = inq.get()
15         if data is None:
16             return
17         fn, arg = data
18         outq.put(fn(arg))
```



```
20 if __name__ == '__main__':
21     import argparse
22
23     parser = argparse.ArgumentParser()
24     parser.add_argument('-n', type=int, default=1)
25     parser.add_argument('number', type=int, nargs='?', default=34)
26     args = parser.parse_args()
27
28     assert args.n >= 1, 'The number of threads has to be > 1'
29
30     tasks = mp.Queue()
31     results = mp.Queue()
32     for i in range(args.n):
33         tasks.put((fib, args.number))
34
35     for i in range(args.n):
36         mp.Process(target=worker, args=(tasks, results)).start()
37
38     for i in range(args.n):
39         print(results.get())
40
41     for i in range(args.n):
42         tasks.put(None)
```



- All Numba array operations
- Many NumPy array operations
- Many NumPy basic math functions
- Some NumPy BLAS functions



```
1  from numba import njit, prange
2  import numpy as np
3
4  @njit(parallel=True, fastmath=True)
5  def two_d_array_reduction_prod(n):
6      shp = (5000, 5000)
7      result1 = 2 * np.ones(shp, np.int_)
8      tmp = 2 * np.ones_like(result1)
9
10     for i in prange(n):
11         result1 *= tmp
12
13     return result1
14
15 if __name__ == '__main__':
16     import argparse
17
18     parser = argparse.ArgumentParser()
19     parser.add_argument('-n', type=int, default=1)
20     args = parser.parse_args()
21
22     two_d_array_reduction_prod(args.n)
```





- Fluent Python
- High Performance Python
- In-Memory Analytics with Apache Arrow
- Productive and Efficient Data Science with Python
- Data Science with Python and Dask
- Distributed Computing with Python
- Parallel and High Performance Computing
- C++ High Performance - Second Edition

## Assignments and Project

---



- Complete and commit all assignments and labs by Tuesday, July 26 to receive grades
- All assignments are listed in the README
- Available for office hours by request and 30 minutes before and after Thursday's (07/21/22) session



- Implement initial improvements for your three optimization targets to discuss at next Thursday's (07/28/22) session
- Available for office hours by request and 30 minutes before and after Thursday's (07/21/22) and Tuesday's (07/26/22) sessions
- Meet with each individual to discuss optimization targets