# Paralelní a distribuované výpočty (B4B36PDV)

## Jakub Mareček
jakub.marecek@fel.cvut.cz

Artificial Intelligence Center
Department of Computer Science
Faculty of Electrical Engineering
Czech Technical University in Prague

# What comes next?

# What comes next?

2015: http://www.openmp.org/wp-content/uploads/openmp-4.5.pdf

2020: https://www.openmp.org/wp-content/uploads/OpenMP-API-Specification-5-1.pdf

2022: https://www.openmp.org/wp-content/uploads/OpenMP-API-Specification-5-2.pdf

2023: TR 12 introduces a preview of OpenMP 6 with support for C++23: https://www.openmp.org/wp-content/uploads/openmp-TR12.pdf

# What is OpenMP?
A Specification

- 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 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.
- There is a long history of the evolution of OpenMP since 1997, from a 50-page long specification of OpenMP 1.0 to 795 pages in TR 12.

# What is OpenMP?
## A Specification

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 oficially supports Microsoft Visual Studio 2022.

# What is OpenMP?

- The specification is built on top of the fork-join model of parallel execution (sériově-paralelní model uspořádání vláken).



By Wikipedia user A1 - w:en:File:Fork_join.svg, CC BY 3.0, https://commons.wikimedia.org/w/index.php?curid=32004077

# What is OpenMP?
A Specification

- The specification is built on top of the fork-join model of parallel execution (sériově-paralelní model uspořádání vláken).
- 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/).

# What is OpenMP?

## A Specification



Page counts (not counting front matter, appendices or index) for versions of OpenMP

"The Ongoing Evolution of OpenMP," https://doi.org/10.1109/JPROC.2018.2853600

# What is OpenMP?

A Specification

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

# What comes next?

- Structuring code
  - OpenMP Task Region
  - Threads and their Sizing
  - Section
  - Task
  - Kernels, Teams, and Targets

- Synchronization primitives
  - Atomic variables
  - Mutexes and locks
  - Critical sections
  - Barrier
  - Fences and Flushes

- For each



Performance Comparison: OpenMP Loops, Tasks, and Dependences
2 X 24-core Intel Xeon Platinum 8160 CPU (48 cores total)

"The Ongoing Evolution of OpenMP," https://doi.org/10.1109/JPROC.2018.2853600

# Structuring code
## 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()).

# Structuring code

## Task Region

```cpp
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 }
```

# Structuring code

The **parallel** 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) .

# Structuring code

The **parallel** construct can take a number of modifiers, including:
- **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 elements 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).

# Structuring code

The **parallel** construct can take a number of modifiers, including:

- **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 elements 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).

```
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  }
```

The reduction 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.

# Structuring code
## Task Region

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_L
EVELS , which controls
the maximum number
of nested active
parallel regions.

```cpp
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 }
```

# Structuring code
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:

```c
 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();
            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  }
```

https://github.com/OpenMP/Examples/blob/main/sources/Example_parallel.1.c

# Structuring code

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.

```cpp
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); }
           }
22      }
23      return 0;
24  }
```
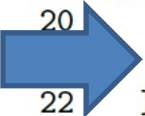
# Structuring code

## 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.

```cpp
 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  }
```

# Structuring code

## 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, Protothreads, Capriccio, StateThreads, TiNy-threads, etc), it often is implemented thus.

```cpp
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 #pragma omp task
20       work();
21     }
22   }
23 }
```

# Structuring code

- 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:

```
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  }
```

```
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 }
```

# Structuring code
## Teams and Targets

- 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.
- Ideally, one should combine the offloading of the code to the GPGU (via the target), across multiple partitions (via teams construct.
- We could also use nowait as in other OpenMP constructs, but that does not match the ``lockstep'' execution on the GPGPUs.

```
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  }
```

```
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 }
```

$$a \cdot \boxed{x_1} + \boxed{y_1} = \boxed{a \cdot x_1 + y_1}$$

$$a \cdot \boxed{x_2} + \boxed{y_2} = \boxed{a \cdot x_2 + y_2}$$

$$\vdots \qquad \vdots \qquad \vdots$$

$$a \cdot \boxed{x_d} + \boxed{y_d} = \boxed{a \cdot x_d + y_d}$$

https://hpc-wiki.info/mediawiki/hpc_images/a/ab/GPU_tutorial_saxpy_cuda_c.pdf

# Structuring code
## Teams and Targets

- 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.
- Ideally, one should combine the offloading of the code to the GPGU (via the target), across multiple partitions (via teams construct.
- We could also use nowait as in other OpenMP constructs, but that does not match the ``lockstep'' execution on the GPGPUs.

```
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   }
```

```
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  }
```

```
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  }
```

# Building the code
## Teams and Targets

Compiling OpenMP with offloading is non-trivial. 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 uses **clang**++ -**fopenmp** -**fopenmp-targets**=<**target triple**>, where the triple is documented at https://llvm.org/doxygen/Triple_8h_source.html
- AMD ROC is built on top of Clang, so starting with **clang** -**fopenmp** -**offload-arch**=**gfx908** is a good idea.
- 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**>.

# What comes next?

- Structuring code
  - OpenMP Task Region
  - Threads and their Sizing
  - Section
  - Task
  - Kernels, Teams, and Targets

- Synchronization primitives
  - Atomic variables
  - Mutexes and locks
  - Critical sections
  - Barrier
  - Fences and Flushes

- For each

# Synchronisation Primitives

## Atomic Variables

OpenMP has a rich support for atomic variables. 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 two weeks ago. The default is relaxed-consistency shared memory model.
- N.B. use {} after atomic!

```cpp
 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  }
```

# Synchronisation Primitives
## Mutexes

OpenMP has only a limited support for mutexes, as in does not support any "resource acquisition is initialization" (RAII) variant, which makes them hard to use correctly.

One can, however, construct one own's RAII variant.

```
 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  }
```

# Synchronisation Primitives

## 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.

```cpp
 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)
    ↪    private(b)
10    {
11      b = omp_get_thread_num()+2;
12  #pragma omp critical
13      {
14        a = a / b;
15        int t = omp_get_thread_num();
16        std::osyncstream(std::cout) << "b = " << b << " in t = "
          ↪    << t << "\n";
17        std::osyncstream(std::cout) << "b = " << a << " in t = "
          ↪    << t << "\n";
18      }
19    }
20    std::osyncstream(std::cout) << "a = " << a << "\n";
21    std::osyncstream(std::cout) << "b = " << b << "\n";
22    return 0;
23  }
```

# Synchronisation Primitives

Barrier

- OpenMP provides a straightforward, explicit barrier construct.
- 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.
- 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, single, scope**, and **workshare** constructs, unless one explicitly adds a nowait clause.

```cpp
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 }
```

# Synchronisation Primitives

## Barrier

- OpenMP provides a straightforward, explicit barrier construct.
- 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.
- 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, single, scope**, and **workshare** constructs, unless one explicitly adds a nowait clause.

```cpp
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 }
```

# Synchronisation Primitives

Fences and Flushes

As we have suggested above, the default memory ordering is relaxed.

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

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

# Algorithms

Last but not least, the most frequent "idiom"

While OpenMP does not really implement any algorithms, some of the data-parallel constructs are similar to algorithms in the STL library. Notably, OpenMP makes it possible to use "for each".

```cpp
 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;
```

# Algorithms

Last but not least, the most frequent "idiom"

While OpenMP does not really implement any algorithms, some of the data-parallel constructs are similar to algorithms in the STL library. Notably, OpenMP makes it possible to use "for each".

```cpp
 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  }
```

# Exercises

A simple exercise is to go through the examples at:
https://github.com/OpenMP/Examples
to understand the syntax in full.

A more demanding 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.

# Backup Slides

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 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).

# Backup Slides

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.)

# Backup Slides

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

# Backup Slides

```
 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  }
```

# Backup Slides

```
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, 2.0l);
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(),
          ↪   sycl::range<1>(vx.size()) };
16      sycl::buffer<float,1> dy { vy.data(),
          ↪   sycl::range<1>(vy.size()) };
17      sycl::buffer<float,1> dz { vz.data(),
          ↪   sycl::range<1>(vz.size()) };
18
19      q.submit([&](sycl::handler& h) {
20        sycl::accessor x(dx, h, sycl::read_only);
21        sycl::accessor y(dy, h, sycl::read_only);
22        sycl::accessor z(dz, h, sycl::read_write);
23        h.parallel_for<class saxpy3>( sycl::range<1>{length},
          ↪   [=] (sycl::id<1> it) {
24          const size_t i = it[0];
25          z[I] += 3.1415 * x[i] + y[I];
26        });
27      });
28      q.wait();
29    }
30    catch (sycl::exception & e) {
31      std::cout << e.what() << std::endl;
32      return 1;
33    }
34    return 0;
35  }
```