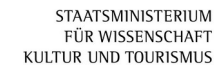




Plasma-PEPSC Webinar 28 May 2024

alpaka Parallel Programming Library



alpaka – Abstraction Library for Parallel Kernel Acceleration

alpaka is...

- A parallel programming library: Accelerate your code by exploiting your hardware's parallelism!
- An abstraction library independent of hardware ecosystem: Create portable code that runs on CPUs and GPUs!
- Free & open-source software

The word "alpaka" is written in a blue, lowercase, sans-serif font. The letter 'p' is stylized with an orange outline that forms a shape resembling a alpaca's head and neck.

Problem of HPC Systems?

Heterogenous Hardware Ecosystem!

TOP500

- 1 Frontier(USA) 1.194 Exaflop/s,
AMD EPYC CPU + AMD Instinct GPU
- 2 Aurora(USA) 585 Petaflop/s,
Intel Xeon CPU + Intel GPU Max
- 3 Eagle(USA) 561 Petaflop/s,
Intel Xeon CPU + Nvidia GPU H100
- 4 Fugaku(Japan) 442 Petaflop/s,
Fujitsu A64FX CPU
- 5 Lumi(Finland) 380 Petaflop/s,
AMD EPYC CPU + AMD Instinct GPU

www.top500.org

THE LIST

11/2023 Highlights

The 62nd edition of the TOP500 shows five new or upgraded entries in the top 10 but the Frontier system still remains the only true exascale machine with an HPL score of 1.194 Exaflop/s.

The Frontier system at the Oak Ridge National Laboratory, Tennessee, USA remains the No. 1 system on the TOP500 and is still the only system reported with an HPL performance exceeding one Exaflop/s. Frontier brought the pole position back to the USA one year ago on the June 2022 listing and has since been remeasured with an HPL score of 1.194 Exaflop/s.

Frontier is based on the latest HPE Cray EX235a architecture and is equipped with AMD EPYC 64C 2GHz processors. The system has 8,699,904 total cores, a power efficiency rating of 52.59 gigaflops/watt, and relies on HPE's Slingshot 11 network for data transfer.

The Aurora system at the Argonne Leadership Computing Facility, Illinois, USA is currently being commissioned and will at full scale exceed Frontier with a peak performance of 2 Exaflop/s. It was submitted with a measurement on half of the final system achieving 585 Petaflop/s on the HPL benchmark which secured the No. 2 spot on the TOP500.

Aurora is built by Intel based on the HPE Cray EX - Intel Exascale Compute Blade which uses Intel Xeon CPU Max Series processors and Intel Data Center GPU Max Series accelerators which communicate through HPE's Slingshot-11 network interconnect.

The Eagle system installed in the Microsoft Azure cloud in the USA is newly listed as No. 3. This Microsoft NDv5 system is based on Intel Xeon Platinum 8480C processors and NVIDIA H100 accelerators and achieved an HPL score of 561 Pflop/s.

[read more »](#)

List Statistics

Vendors System Share

1 **Frontier** - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE

2 **Aurora** - HPE Cray EX - Intel Exascale Compute Blade, Xeon CPU Max 9470 52C 2.4GHz, Intel Data Center GPU Max, Slingshot-11, Intel

3 **Eagle** - Microsoft NDv5, Xeon Platinum 8480C 48C 2GHz, NVIDIA H100, NVIDIA Infiniband NDR, Microsoft

4 **Supercomputer Fugaku** - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu

5 **LUMI** - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE

6 **Leonardo** - BullSequana XH2000, Xeon Platinum 8358 32C 2.6GHz, NVIDIA A100 SXM4 64 GB, Quad-rail NVIDIA HDR100 Infiniband, EVIDEN

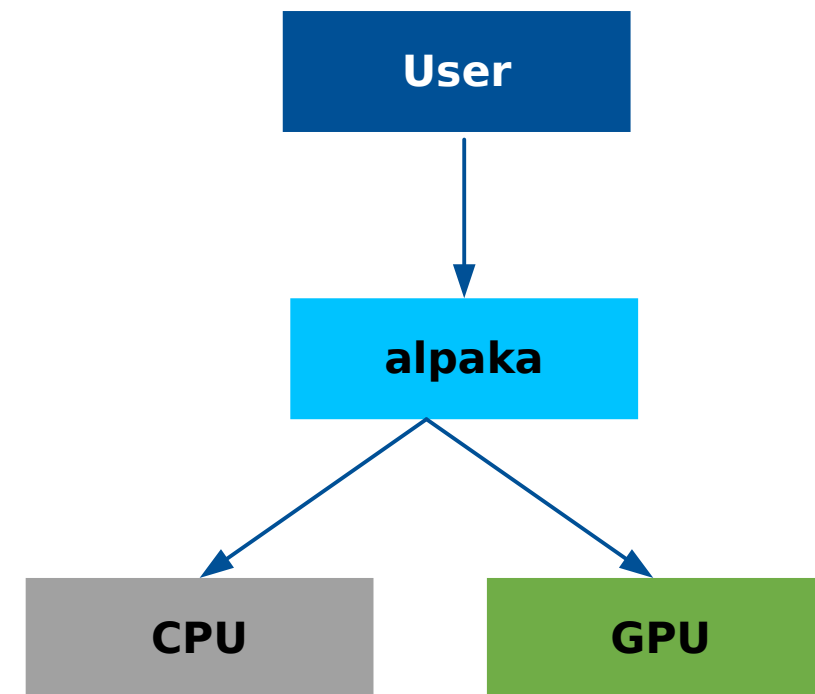
7 **Summit** - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-

Without alpaka

- Hardware ecosystem is heterogenous, platforms are not inter-operable → parallel programs not easily portable

alpaka: one API to rule them all

- **Abstraction** (not hiding!) of the underlying hardware, compiler and OS
 - No default device, built-in functions, language extensions
- **Easy change of the backend in code**
- **Direct usage of vendor APIs, not depend on “unified APIs”**
 - GPU Backends: Hip (AMD), Cuda (NVidia), SYCL (Intel GPUs)
 - One can use vendor profilers and debuggers (Cuda,HIP...) for alpaka code!
 - CPU Backends: OpenMp, Threads, TbbBlocks
- **Zero abstraction overhead for Kernel execution!**
- **Heterogenous Programming:** Using different backends in a synchronized manner.



Find us on GitHub!

alpaka library: <https://www.github.com/alpaka-group/alpaka>

- Full source code and many examples, Issue tracker

The documents: <https://alpaka.readthedocs.io/en/latest/>

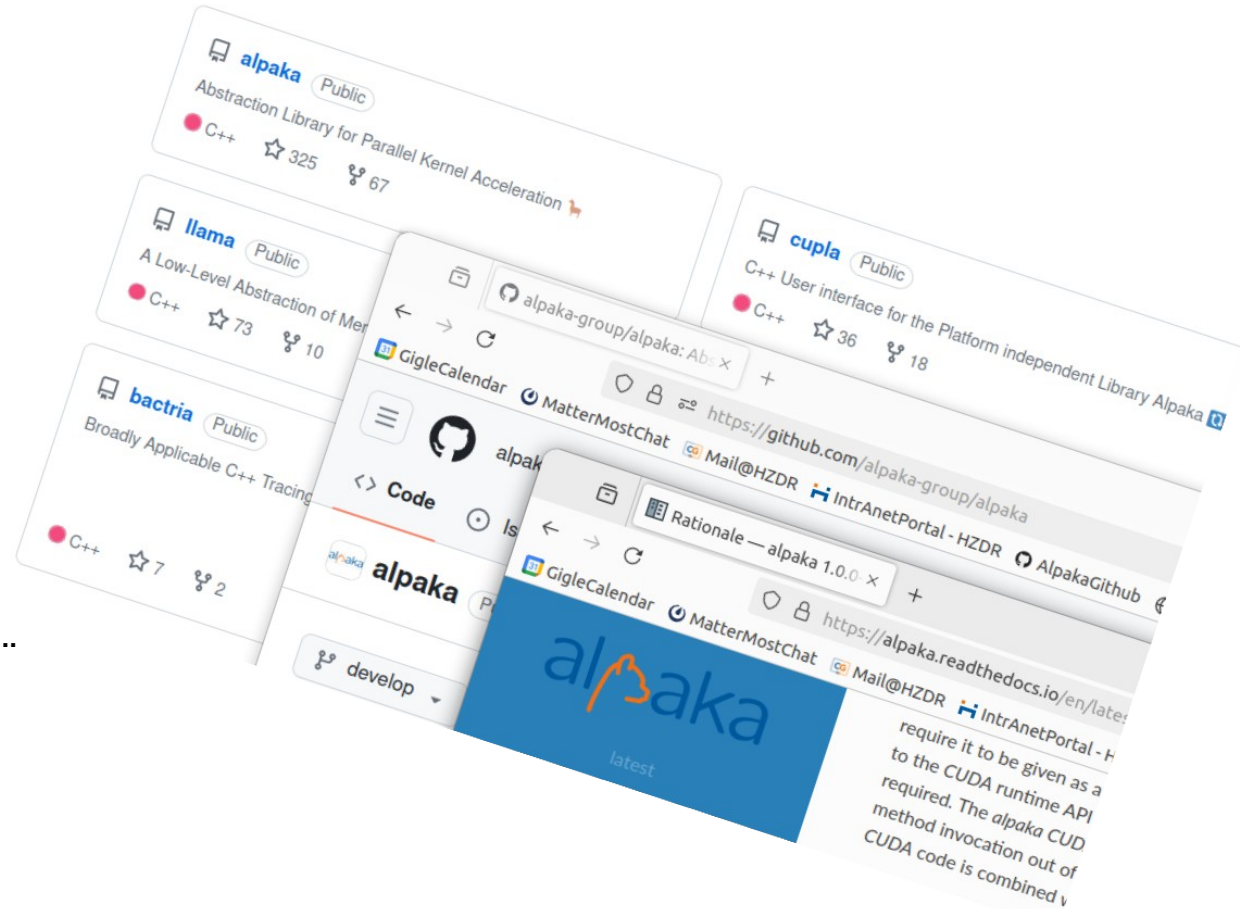
- Installation guide
- Cheatsheet
- Abstraction model and the rationale behind alpaka

Project group: <https://www.github.com/alpaka-group>

- Contains all alpaka-related projects, documentation, samples, ...

alpaka is a free software (MPL 2.0)

moz://a Public License



Programming with alpaka

- C++ only!
- alpaka is written entirely in C++17. Coming soon: C++20.
- Header-only library. No additional runtime dependency.
`#include <alpaka/alpaka.hpp>` is enough!
- Supports a wide range of modern C++ compilers (g++, clang++, Apple LLVM, MSVC)
- Portable across operating systems: Linux, macOS, Windows



Installation and Examples

I-Install Dependencies

- alpaka requires **Boost**, **Cmake** and a modern C++ compiler (g++, clang++, Visual C++, ...)
 - Linux: `sudo apt install libboost-all-dev` (DEB)
 - MacOS: `brew install boost` (using homebrew, <https://brew.sh>)
 - Windows: `vcpkg install boost` (using vcpkg, <https://github.com/microsoft/vcpkg>)
- Depending on your target platform you may need additional packages
 - NVIDIA GPUs: CUDA Toolkit (<https://developer.nvidia.com/cuda-toolkit>)
 - AMD GPUs: ROCm / HIP (<https://rocm.docs.amd.com/en/latest/index.html>)
 - Intel GPUs: OneAPI Toolkit (<https://www.intel.com/content/www/us/en/developer/tools/oneapi/base-toolkit.html#gs.9x3lnh>)
- CMake is the preferred system for building and installing
 - Linux: `sudo apt install cmake` (DEB). macOS and Windows: Download the installer from <https://cmake.org/download/>

II - Compiling and running examples

- You can build all examples at once from your build directory:
 - configure the build with setting some cmake variables according to your system

```
cmake -Dalpaka_BUILD_EXAMPLES=ON -DCMAKE_BUILD_TYPE=Release -Dalpaka_ACC_CPU_B_SEQ_T_SEQ_ENABLE=ON -Dalpaka_ACC_GPU_CUDA_ENABLE=ON ..
```
 - build the examples

```
cmake --build . --config Release
```
 - alpaka/build/example/ directory will include compiled examples.
e.g. alpaka/build/example/vectorAdd directory will include the executable `vectorAdd`
- Run all examples from the build directory of alpaka

```
ctest example/
```
- Run all tests from the build directory of alpaka if `BUILD_TESTING` cmake variable is set to ON.

```
ctest test/
```
- Examples can be re-compiled and run in their corresponding directories under build directory if there is a code change in the source tree.

```
cd alpaka/build/example/vectorAdd  
cmake --build . (or run the make command if make file is there )
```

```
1 git clone https://github.com/alpaka-group/alpaka.git
2 cd alpaka
3 mkdir build
4 cd build/
5 cmake -Dalpaka_BUILD_EXAMPLES=ON -DCMAKE_BUILD_TYPE=Release -Dalpaka_
  ACC_CPU_B_SEQ_T_SEQ_ENABLE=ON -Dalpaka_ACC_GPU_CUDA_ENABLE=ON ..
6 cmake --build . --config Release
7 cd example/vectorAdd/
8 ./vectorAdd
```


III - Install alpaka Library

Download alpaka: `git clone -b develop https://github.com/alpaka-group/alpaka.git`

- In the terminal/powershell, switch to the downloaded alpaka directory:
`cd /path/to/alpaka`
- Create a build directory and switch to it:
`mkdir build`
`cd build`
- Configure build directory (If default directories is ok for you or you are planning to use alpaka from build directory; you can omit the install prefix cmake variable)
`cmake -DCMAKE_INSTALL_PREFIX=/some/other/path/ ..`
- **Install alpaka directly, compilation is not needed since it is a header only library!**
alpaka installation will reside in `/some/other/path/`.
`cmake --install .`
- You should now have a complete alpaka installation in the directory you chose earlier.

```
git clone https://github.com/alpaka-group/alpaka.git
cd alpaka
mkdir build
cd build
cmake -DCMAKE_INSTALL_PREFIX=/some/other/path/ ..
cmake --install .
```

For Detailed information:

<https://github.com/alpaka-group/alpaka-workshop-slides/tree/develop>

IV - Create your first alpaka project

• Option 1

1. After installation, create your cmake file as the one below in a new project directory. Use

```
find_package(alpaka REQUIRED)
```

2. Create your cpp code (possibly under `src` directory of the project directory)

3. Inside the directory of your project run commands:

```
mkdir build
```

```
cd build
```

4. Configure the example according to your system:

```
cmake -DCMAKE_BUILD_TYPE=Release  
-Dalpaka_ACC_CPU_B_SEQ_T_SEQ_ENABLE=ON -Dalpaka_ACC_GPU_CUDA_ENABLE=ON ..
```

5. Build your code

```
cmake --build . --config Release
```

```
##  
# CMakeLists.txt for the myHelloWorld example  
cmake_minimum_required(VERSION 3.22)  
  
set(_TARGET_NAME myHelloWorld)  
project(${_TARGET_NAME} LANGUAGES CXX)  
# Find alpaka.  
find_package(alpaka REQUIRED)  
# Add executable.  
alpaka_add_executable(  
    ${_TARGET_NAME}  
    src/myHelloWorld.cpp)  
  
target_link_libraries(  
    ${_TARGET_NAME}  
    PUBLIC alpaka::alpaka)
```

```
myHelloWorld/  
├── build  
│   ├── CMakeCache.txt  
│   ├── CMakeFiles  
│   ├── cmake_install.cmake  
│   ├── Makefile  
│   └── myHelloWorld  
├── CMakeLists.txt  
├── src  
│   └── myHelloWorld.cpp
```

Option 2

1. Copy one of the examples to a folder outside the alpaka source tree. The folder will be your project directory.

2. Remove `ExampleDefaultAcc` and `getAccName` from the code. Select a backend (accelerator) inside the code instead of `ExampleDefaultAcc`.

3. Clone alpaka under your project-directory

4. Add alpaka directory to `CmakeLists.txt` by `add_subdirectory`

5. Repeat steps 3-5 of **Option1** under your project directory.

```
# CMakeLists.txt for using alpaka by adding alpaka  
# as a subdirectory to your project.  
cmake_minimum_required(VERSION 3.22)  
  
# Project.  
set(_TARGET_NAME myVectorAddProject)  
project(${_TARGET_NAME} LANGUAGES CXX)  
  
# add alpaka directory name(cloned into your project)  
add_subdirectory(./alpaka)  
  
# Add executable.  
alpaka_add_executable(  
    ${_TARGET_NAME}  
    src/vectorAdd.cpp)  
target_link_libraries(  
    ${_TARGET_NAME}  
    PUBLIC alpaka::alpaka)
```

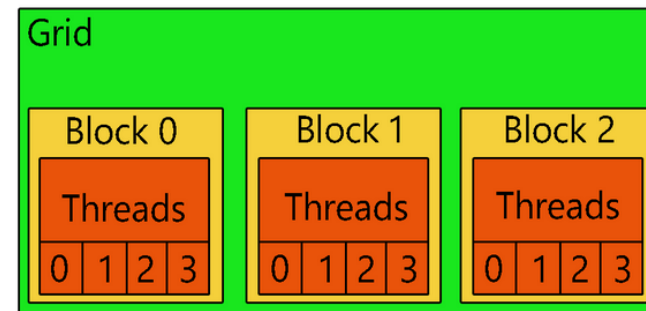
```
myVectorAddProject/  
├── alpaka  
│   ├── benchmarks  
│   ├── CHANGELOG.md  
│   ├── cmake  
│   ├── CMakeLists.txt  
│   ├── CMakePresets.json  
│   ├── CONTRIBUTING.md  
│   ├── docs  
│   ├── example  
│   ├── include  
│   ├── LICENSE  
│   ├── README.md  
│   ├── README_SYCL.md  
│   ├── script  
│   ├── test  
│   └── thirdParty  
├── build  
│   ├── alpaka  
│   ├── CMakeCache.txt  
│   ├── CMakeFiles  
│   ├── cmake_install.cmake  
│   ├── Makefile  
│   └── myVectorAddProject  
├── CMakeLists.txt  
├── src  
│   └── vectorAdd.cpp
```

Tenets of Thread-Parallel Programming

- Grid, block and thread based paralelisation model.**

The model is instantiated differently on different processors, because of cache size and speed, the synchronization mechanism of the hardware, or simply the CPU-GPU difference. By using this abstraction the execution can be optimally adapted to the available hardware.

- Large number of threads** should run the same code (**kernel**) on different data in parallel.
- Indexing of threads.** Each thread should work on a different data portion or do a specific task, therefore each thread has an index accessible in kernel.



Grid-Block Extent: $\text{Vec}\{3\}$
Block-Thread Extent: $\text{Vec}\{4\}$



Grid-Block Extent: $\text{Vec}\{3\}$ or $\text{Vec}\{1,3\}$
Block-Thread Extent: $\text{Vec}\{4,5\}$

- **Extent:** A vector representing the sizes along dimensions.

In 3d an extent is {Width,Length,Height}

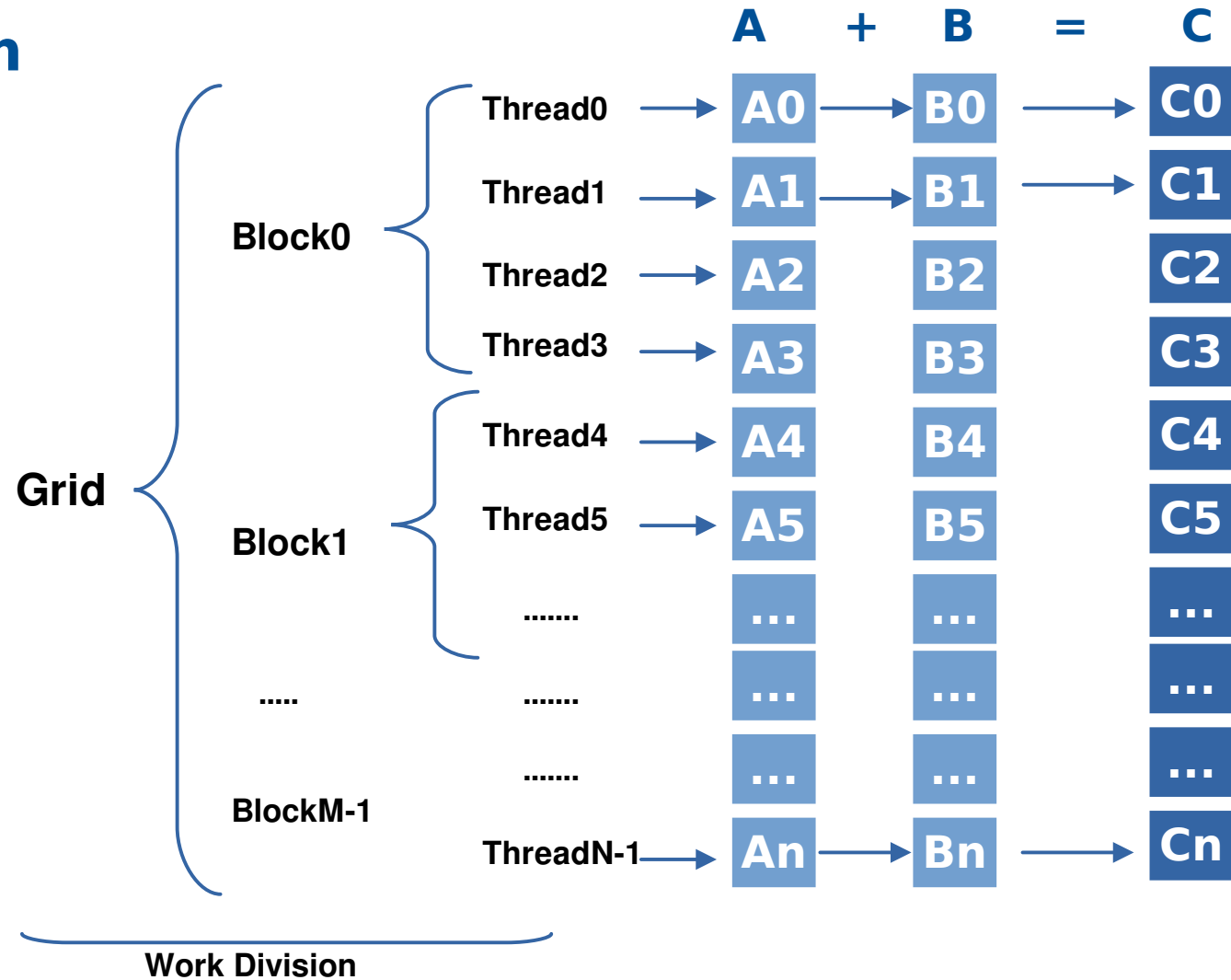


- **Dimensions:** Set of dimension names. {X-dimension, Y-dimension, Z-dimension}

- **Number of Dimensions**

Vector Addition Problem

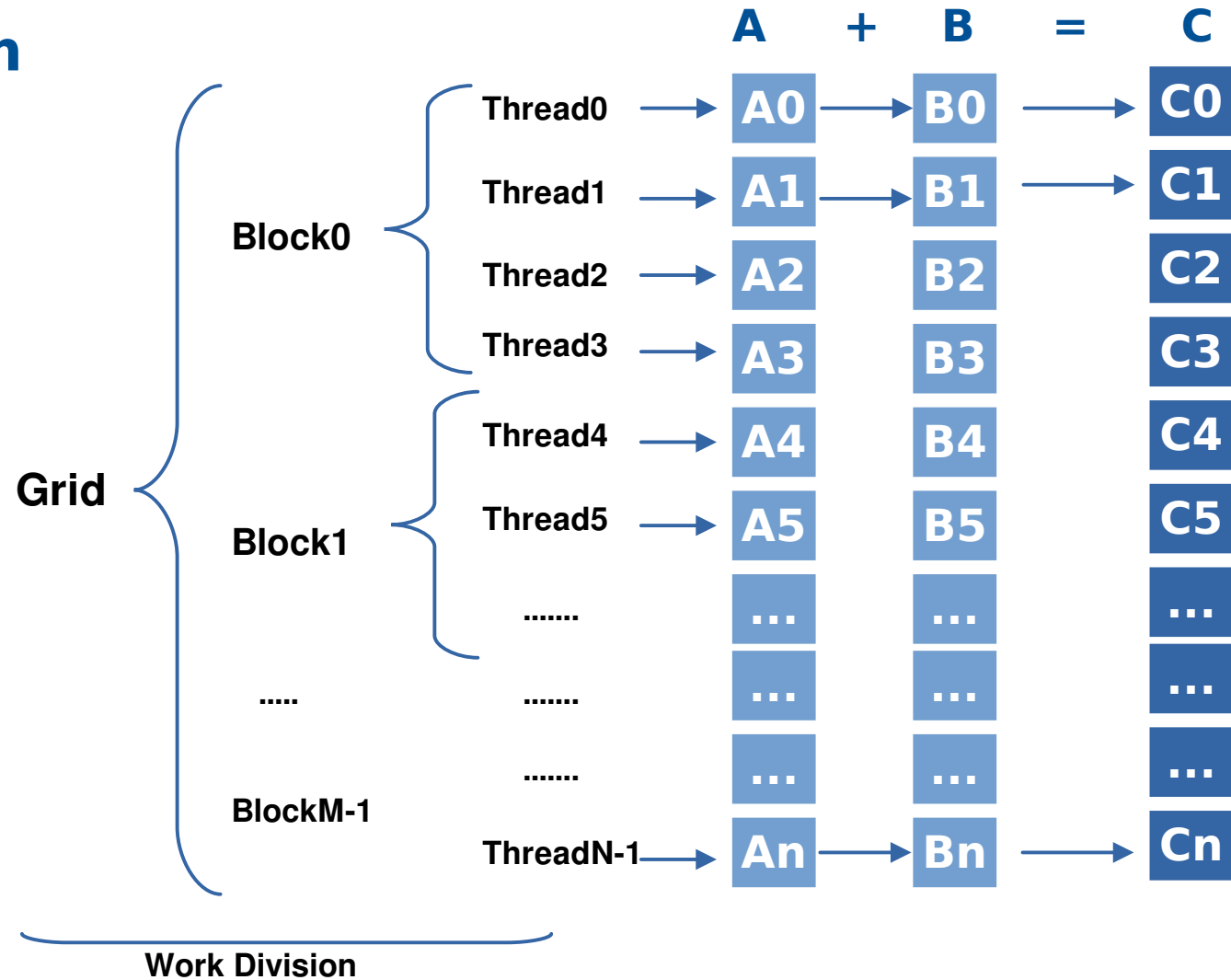
$$C = A + B$$



$$C = A + B$$

One to one, Identity

II- Grouping threads



Vector Addition Code Steps

1. Create Kernel.
2. Decide where will the paralel and non-parallel parts of the code run.
3. Decide how to parallelise (number of blocks and threads).
4. Allocate host and device memory for A,B and C.
5. Copy the memory to the device.
6. Run the kernel
7. Copy the result data back to the host.

1. Define the alpaka Kernel

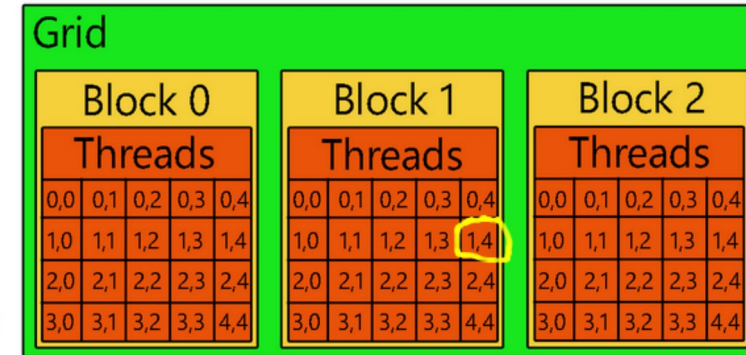
- Contains the algorithm that is run by each thread
- alpaka Kernels are functors (function-like C++ structs / classes) or lambdas
- Arguments can be pointers and trivially copyable types
- Agnostic to device details

```
// Single header library
#include <alpaka/alpaka.hpp>

#include <iostream>
//! An example kernel: vector addition
class VectorAddKernel
{
public:
    ALPAKA_NO_HOST_ACC_WARNING
    template<typename TAcc, typename TElem, typename TIdx>
    ALPAKA_FN_ACC auto operator()(
        TAcc const& acc, // the accelerator
        TElem const* const A, TElem const* const B, TElem* const C,
        TIdx const& numElements ) const -> void
    {
        static_assert(alpaka::Dim<TAcc>::value == 1, "The kernel expects 1-dimensional indices!");
        // Get thread index
        TIdx const gridThreadId(alpaka::getIdx<alpaka::Grid, alpaka::Threads>(acc)[0u]);

        if(gridThreadId < numElements)
        {
            // Use thread index as the data index
            C[gridThreadId] = A[gridThreadId] + B[gridThreadId];
        }
    }
};
```

Obtaining the indices of threads/blocks inside the Kernel



- Index of Thread on the Grid:

```
auto gridThreadIndex = alpaka::getIdx<alpaka::Grid, alpaka::Threads>(acc);  
// gridThreadIndex is {1,9}
```
- Index of Thread on a Block:

```
auto threadBlockIndex = alpaka::getIdx<alpaka::Block, alpaka::Threads>(acc);  
// threadBlockIndex is {1,4}
```
- Index of Block on the Grid:

```
auto blockGridIndex = alpaka::getIdx<alpaka::Grid, alpaka::Blocks>(acc);  
// the blockGridIndex is {1}
```


2. Select the Accelerator, Platform and Device

- alpaka provides a number of pre-defined **Accelerators**.

- Acc**Gpu**CudaRt for **Nvidia** GPUs
- Acc**Gpu**HipRt for **AMD**, **Intel** and **Nvidia** GPUs
- Acc**Gpu**SyclIntel for **AMD**, **Intel** and **Nvidia** GPUs

- Acc**Cpu**Fibers based on Boost.fiber
- Acc**Cpu**Omp2Blocks based on OpenMP 2.x
- Acc**Cpu**Omp4 based on OpenMP 4.x
- Acc**Cpu**TbbBlocks based on TBB
- Acc**Cpu**Threads based on `std::thread`
- Acc**Cpu**Sycl

- Acc**Fpga**SyclIntel

- Device** instance represents a single physical device

```
auto main() -> int
{
    using Dim = alpaka::DimInt<1u>; // Define the index domain
    using Idx = std::size_t; // Index type of the threads and buffers
    using DataType = std::uint32_t; // Define the buffer element type

    // Define the accelerator: AccGpuCudaRt, AccGpuHipRt,
    // AccCpuThreads, AccCpuOmp2Threads, AccCpuOmp2Blocks, AccCpuTbbBlocks AccCpuSerial
    using Acc = alpaka::AccGpuCudaRt<Dim, Idx>;
    using DevAcc = alpaka::Dev<Acc>;

    // Select a device from platform of Acc
    auto const platform = alpaka::Platform<Acc>{};
    auto const devAcc = alpaka::getDevByIdx(platform, 0);
}
```

3. How to parallelise?

I- Get a valid work division from alpaka

Use `getValidWorkDiv` function

- The function divides the full grid-thread extent into blocks.
- Inputs:
 - Full grid-thread extent. (User provides total number of threads needed.)
 - Elements per thread extent
- Most probable `workDiv` in the code will be `{Vec{numElements/1024}, Vec{1024}, Vec{1}}`

II - Determine the workdivision manually

- WorkDivision data structure consists 3 vectors:
 - Grid block extent.
`Vec{numElements/1024}` or `Vec{1,1,numElements/2014}`
 depending on the number of dimensions.
 - Block thread extent.
`Vec{1024}` or `Vec{1,1,1024}`
 - Elements per thread is `Vec{1}` or `Vec{1,1,1}`

```
auto const platform = alpaka::Platform<Acc>{};
auto const devAcc = alpaka::getDevByIdx(platform, 0);

// Define the work division depending on the data
Idx const numElements(100000);
Idx const elementsPerThread(1u);
alpaka::Vec<Dim, Idx> const extent(numElements);

// Let alpaka calculate good block and grid sizes given our full problem extent
alpaka::WorkDivMembers<Dim, Idx> const workDiv = alpaka::getValidWorkDiv<Acc>(
    devAcc, // device
    extent, // {length, height, depth} of grid. For 1D only length of the vector: {length!
    elementsPerThread,
    false, alpaka::GridBlockExtentSubDivRestrictions::Unrestricted);

.....

.....

.....
// Instantiate the kernel function object
VectorAddKernel kernel;
alpaka::exec<Acc>( // Run the kernel execution task
    queue,
    workDiv,
    kernel, alpaka::getPtrNative(bufAccA), alpaka::getPtrNative(bufAccB),
    alpaka::getPtrNative(bufAccC),
    numElements);
// Copy back the result
alpaka::memcpy(queue, bufHostC, bufAccC); // bufHostC includes the result!
alpaka::wait(queue);
}
```

4. Allocate data vectors A and B on host and device.

- **alpaka::Buf** is multi-dimensional dynamic array.

It contains

- *memory*,
- *size*,
- the *device* the memory is allocated in!
- **alpaka::allocBuf()** allocates memory to the given device.
- Alternatively **alpaka::View** is used to adapt already allocated memory. if we already have an STL vector at host for example we don't need to create Buf.

```
using Acc = alpaka::AccGpuCudaRt<Dim, Idx>;
using DevAcc = alpaka::Dev<Acc>;

// Select a device from platform of Acc
auto const platform = alpaka::Platform<Acc>{};
auto const devAcc = alpaka::getDevByIdx(platform, 0);

// Define the work division depending on the data
Idx const numElements(100000);
Idx const elementsPerThread(1u);
alpaka::Vec<Dim, Idx> const extent(numElements);

....
// Get the host device for allocating memory on the host.
auto const platformHost = alpaka::PlatformCpu{};
// Get the device directly from CPU platform not from the platform of Acc
auto const devHost = alpaka::getDevByIdx(platformHost, 0);

// Host device type is needed, because it is not known (for the backend it is known in Acc)
using DevHost = alpaka::DevCpu;
// Allocate 3 host memory buffers
using BufHost = alpaka::Buf<DevHost, Data, Dim, Idx>;
BufHost bufHostA(alpaka::allocBuf<Data, Idx>(devHost, extent));
BufHost bufHostB(alpaka::allocBuf<Data, Idx>(devHost, extent));
BufHost bufHostC(alpaka::allocBuf<Data, Idx>(devHost, extent));

// Fill the host buffers
for(Idx i(0); i < numElements; ++i)
{
    bufHostA[i] = randomA; bufHostB[i] = randomB; bufHostC[i] = 0;
}

// Allocate 3 buffers on the accelerator
using BufAcc = alpaka::Buf<DevAcc, Data, Dim, Idx>;
BufAcc bufAccA(alpaka::allocBuf<Data, Idx>(devAcc, extent));
BufAcc bufAccB(alpaka::allocBuf<Data, Idx>(devAcc, extent));
BufAcc bufAccC(alpaka::allocBuf<Data, Idx>(devAcc, extent));
```

5.1 Create the Queue for memcpy and kernel task

- alpaka::Queue is “a queue of tasks”
- Queue is always FIFO, everything is sequential inside the queue.
- *and more*
 - *Different queues run in parallel for many devices*
 - *Used for synchronization*
 - *Accelerator back-ends can be mixed (used in interleaves) within a device queue.*
 - ...

```
// Create a queue on the device, define the synchronization behaviour
alpaka::Queue<Acc, alpaka::Blocking> queue(devAcc);

// Copy from Host to Acc
alpaka::memcpy(queue, bufAccA, bufHostA);
alpaka::memcpy(queue, bufAccB, bufHostB);
alpaka::memcpy(queue, bufAccC, bufHostC);

// Instantiate the kernel function object
VectorAddKernel kernel;
alpaka::exec<Acc>( // Run the kernel execution task
    queue,
    workDiv,
    kernel, alpaka::getPtrNative(bufAccA), alpaka::getPtrNative(bufAccB),
    alpaka::getPtrNative(bufAccC),
    numElements);
// Copy back the result
alpaka::memcpy(queue, bufHostC, bufAccC); // bufHostC includes the result!
alpaka::wait(queue);
```


5.2 Copy data vectors to the Device

- **alpaka::memcpy** copies the data from one buffer/view to another buffer or view.
- **alpaka::Buf** knows the device it belongs to.

```
// Allocate 3 buffers on the accelerator
using BufAcc = alpaka::Buf<DevAcc, Data, Dim, Idx>;
BufAcc bufAccA(alpaka::allocBuf<Data, Idx>(devAcc, extent));
BufAcc bufAccB(alpaka::allocBuf<Data, Idx>(devAcc, extent));
BufAcc bufAccC(alpaka::allocBuf<Data, Idx>(devAcc, extent));

// Define the synchronization behavior of a queue
// choose between Blocking and NonBlocking
// Create a queue on the device
alpaka::Queue<Acc, alpaka::Blocking> queue(devAcc);

// Copy from Host to Acc
alpaka::memcpy(queue, bufAccA, bufHostA);
alpaka::memcpy(queue, bufAccB, bufHostB);
alpaka::memcpy(queue, bufAccC, bufHostC);
```

6. Execute the kernel

- Call **alpaka::exec** function
- The result is stored in an **alpaka::Buf**

7. Copy result back

- Copy the result in device to the host

```
// Instantiate the kernel function object
VectorAddKernel kernel;

alpaka::exec<Acc>(// Run the kernel execution task
    queue,
    workDiv,
    kernel,
    alpaka::getPtrNative(bufAccA),
    alpaka::getPtrNative(bufAccB),
    alpaka::getPtrNative(bufAccC),
    numElements);
// Copy back the result
alpaka::memcpy(queue, bufHostC, bufAccC); // bufHostC includes the result!
```

Parallel vector addition code

```
// Single header library
#include <alpaka/alpaka.hpp>

#include <iostream>

//! An example kernel: vector addition
class VectorAddKernel
{
public:
    ALPAKA_NO_HOST_ACC_WARNING
    template<typename TAcc, typename TElem, typename TIdx>
    ALPAKA_FN_ACC auto operator()(
        TAcc const& acc, // the accelerator
        TElem const* const A,
        TElem const* const B,
        TElem* const C,
        TIdx const& numElements) const -> void
    {
        static_assert(alpaka::Dim<TAcc>::value == 1, "Kernel expects 1-dimensional indices!");
        // Get thread index
        TIdx const gridThreadIdIdx(alpaka::getIdx<alpaka::Grid, alpaka::Threads>(acc)[0u]);

        if(gridThreadIdIdx < numElements)
        {
            // Use thread index as the data index
            C[gridThreadIdIdx] = A[gridThreadIdIdx] + B[gridThreadIdIdx];
        }
    }
};

auto main() -> int
{
    using Dim = alpaka::DimInt<1u>; // Define the index domain
    using Idx = std::size_t; // Index type of the threads and buffers
    using DataType = std::uint32_t; // Define the buffer element type

    // Define the accelerator: AccGpuCudaRt, AccGpuHipRt,
    // AccCpuThreads, AccCpuOmp2Threads, AccCpuOmp2Blocks, AccCpuTbbBlocks, AccCpuSerial
    using Acc = alpaka::AccGpuCudaRt<Dim, Idx>;
    using DevAcc = alpaka::Dev<Acc>;

    // Select a device from platform of Acc
    auto const platform = alpaka::Platform<Acc>{};
    auto const devAcc = alpaka::getDevByIdx(platform, 0);

    // Define the work division depending on the data
    Idx const numElements(100000);
    Idx const elementsPerThread(1u);
    alpaka::Vec<Dim, Idx> const extent(numElements);
```

```
// Let alpaka calculate good block and grid sizes given our full problem extent
alpaka::WorkDivMembers<Dim, Idx> const workDiv = alpaka::getValidWorkDiv<Acc>({
    devAcc, // device
    extent, // {length, height, depth} of grid. For 1D only length of the vector!
    elementsPerThread,
    false, alpaka::GridBlockExtentSubDivRestrictions::Unrestricted});

// Get the host device for allocating memory on the host.
auto const platformHost = alpaka::PlatformCpu{};
// Get the device directly from CPU platform not from the platform of Acc
auto const devHost = alpaka::getDevByIdx(platformHost, 0);

// Host device type is needed, because it is not known (for the backend it is known in Acc)
using DevHost = alpaka::DevCpu;
// Allocate 3 host memory buffers
using BufHost = alpaka::Buf<DevHost, DataType, Dim, Idx>;
BufHost bufHostA(alpaka::allocBuf<DataType, Idx>(devHost, extent));
BufHost bufHostB(alpaka::allocBuf<DataType, Idx>(devHost, extent));
BufHost bufHostC(alpaka::allocBuf<DataType, Idx>(devHost, extent));

// Fill the buffers
for(Idx i(0); i < numElements; ++i)
{ bufHostA[i] = randomA; bufHostB[i] = randomB; bufHostC[i] = 0; }

// Allocate 3 buffers on the accelerator
using BufAcc = alpaka::Buf<DevAcc, DataType, Dim, Idx>;
BufAcc bufAccA(alpaka::allocBuf<DataType, Idx>(devAcc, extent));
BufAcc bufAccB(alpaka::allocBuf<DataType, Idx>(devAcc, extent));
BufAcc bufAccC(alpaka::allocBuf<DataType, Idx>(devAcc, extent));

// Create a queue on the device, define the synchronization behaviour
alpaka::Queue<Acc, alpaka::Blocking> queue(devAcc);

// Copy from Host to Acc
alpaka::memcpy(queue, bufAccA, bufHostA);
alpaka::memcpy(queue, bufAccB, bufHostB);
alpaka::memcpy(queue, bufAccC, bufHostC);

// Instantiate the kernel function object
VectorAddKernel kernel;
alpaka::exec<Acc>({ // Run the kernel execution task
    queue,
    workDiv,
    kernel, alpaka::getPtrNative(bufAccA), alpaka::getPtrNative(bufAccB),
    alpaka::getPtrNative(bufAccC),
    numElements});
// Copy back the result
alpaka::memcpy(queue, bufHostC, bufAccC); // bufHostC includes the result!
}
```


Parallel vector addition code

```
// Single header library
#include <alpaka/alpaka.hpp>

#include <iostream>

//! An example kernel: vector addition
class VectorAddKernel
{
public:
    ALPAKA_NO_HOST_ACC_WARNING
    template<typename TAcc, typename TElem, typename TIdx>
    ALPAKA_FN_ACC auto operator()(
        TAcc const& acc, // the accelerator
        TElem const* const A,
        TElem const* const B,
        TElem* const C,
        TIdx const& numElements) const -> void
    {
        static_assert(alpaka::Dim<TAcc>::value == 1, "Kernel expects 1-dimensional indices!");
        // Get thread index
        TIdx const gridThreadIdx(alpaka::getIdx<alpaka::Grid, alpaka::Threads>(acc)[0u])

        if(gridThreadIdx < numElements)
        {
            // Use thread index as the data index
            C[gridThreadIdx] = A[gridThreadIdx] + B[gridThreadIdx];
        }
    }
};

auto main() -> int
{
    using Dim = alpaka::DimInt<1u>; // Define the index domain
    using Idx = std::size_t; // Index type of the threads and buffers
    using DataType = std::uint32_t; // Define the buffer element type

    // Define the accelerator: AccGpuCudaRt, AccGpuHipRt,
    // AccCpuThreads, AccCpuOmp2Threads, AccCpuOmp2Blocks, AccCpuTbbBlock, ...
    using Acc = alpaka::AccGpuCudaRt<Dim, Idx>;
    using DevAcc = alpaka::Dev<Acc>;

    // Select a device from platform of Acc
    auto const platform = alpaka::Platform<Acc>{};
    auto const devAcc = alpaka::getDevByIdx(platform, 0);

    // Define the work division depending on the data
    Idx const numElements(100000);
    Idx const elementsPerThread(1u);
    alpaka::Vec<Dim, Idx> const extent(numElements);
```

Single header

kernel

Select accelerator and the device (GPU)

```
// Let alpaka calculate good block and grid sizes given our full problem extent
alpaka::WorkDivMembers<Dim, Idx> const workDiv = alpaka::getValidWorkDiv<Acc>(*
    devAcc, // device
    extent, // {length, height, depth} of grid. For 1D only length of the vector!
    elementsPerThread,
    false, alpaka::GridBlockExtentSubDivRestrictions::Unrestricted);

// Get the host device for allocating memory on the host.
auto const platformHost = alpaka::PlatformCpu{};
// Get the device directly from CPU platform not from the platform of Acc
auto const devHost = alpaka::getDevByIdx(platformHost, 0);

// Host device type is needed, because it is not known (for the backend it is known in Acc)
using DevHost = alpaka::DevCpu;
// Allocate 3 host memory buffers
using BufHost = alpaka::Buf<DevHost, DataType, Dim, Idx>;
BufHost bufHostA(alpaka::allocBuf<DataType, Idx>(devHost, extent));
BufHost bufHostB(alpaka::allocBuf<DataType, Idx>(devHost, extent));
BufHost bufHostC(alpaka::allocBuf<DataType, Idx>(devHost, extent));

// Fill the buffers
for(Idx i(0); i < numElements; ++i)
{ bufHostA[i] = randomA; bufHostB[i] = randomB; bufHostC[i] = 0; }

// Allocate 3 buffers on the accelerator
using BufAcc = alpaka::Buf<DevAcc, DataType, Dim, Idx>;
BufAcc bufAccA(alpaka::allocBuf<DataType, Idx>(devAcc, extent));
BufAcc bufAccB(alpaka::allocBuf<DataType, Idx>(devAcc, extent));
BufAcc bufAccC(alpaka::allocBuf<DataType, Idx>(devAcc, extent));

// Create a queue on the device, define the synchronization behaviour
alpaka::Queue<Acc, alpaka::Blocking> queue(devAcc);

// Copy from Host to Acc
alpaka::memcpy(queue, bufAccA, bufHostA);
alpaka::memcpy(queue, bufAccB, bufHostB);
alpaka::memcpy(queue, bufAccC, bufHostC);

// Instantiate the kernel function object
VectorAddKernel kernel;
alpaka::exec<Acc>(&queue, // Run the kernel execution task
    workDiv,
    kernel, alpaka::getPtrNative(bufAccA), alpaka::getPtrNative(bufAccB),
    alpaka::getPtrNative(bufAccC),
    numElements);

// Copy back the result
alpaka::memcpy(queue, bufHostC, bufAccC); // bufHostC includes the result!
```

Get Work division

Get host device (CPU)

allocate memory at host (CPU)

allocate memory at device (GPU)

Copy vectors to device (GPU)

Execute kernel

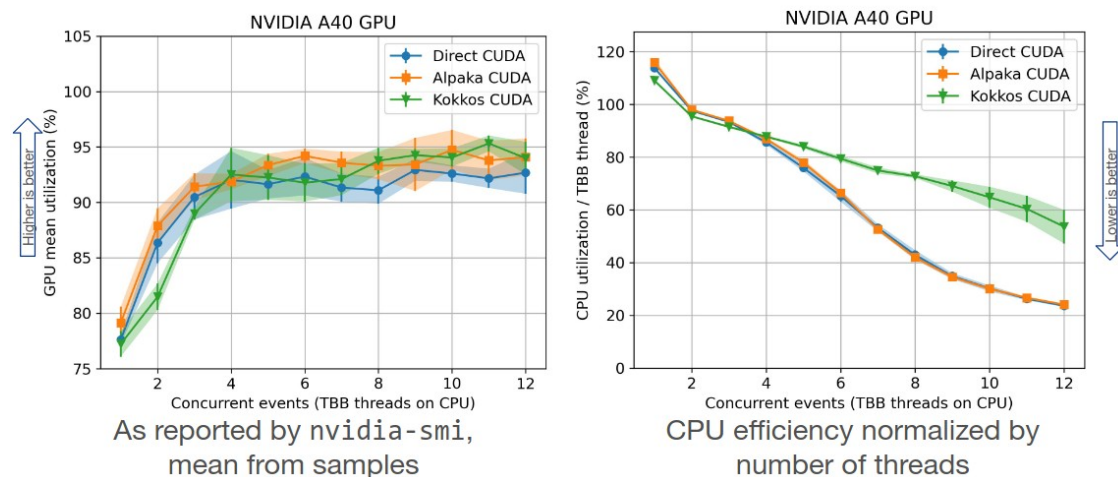
Copy result to host (CPU)

Programing Tips

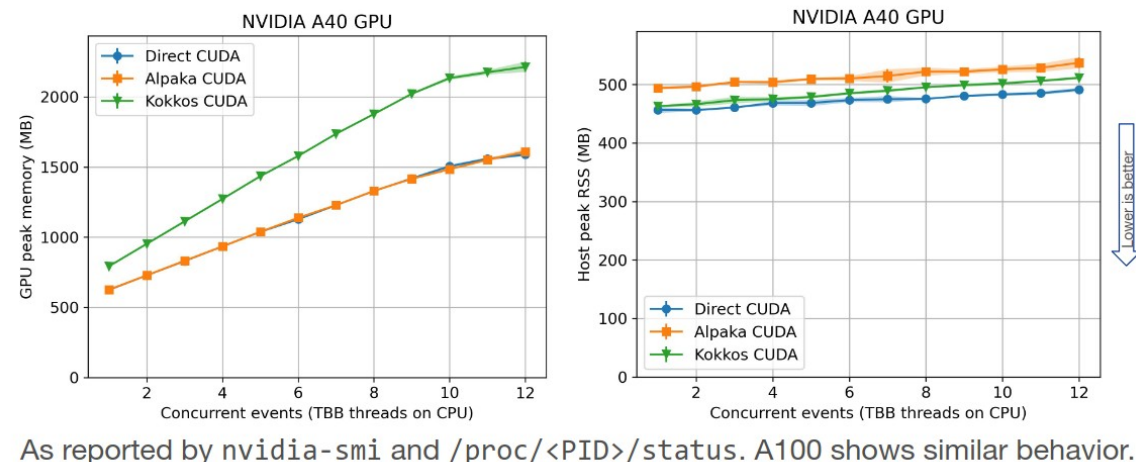
- You can do printf debugging; but can **not use** `std::cout` in alpaka Kernel
- If you want to pass multi-dimensional data to kernel, use `mdspan` (enable it via cmake option).
(If you don't use `mdspan`; you will need to take care of alignment/pitch values. Pass the pointer, extents and the pitch.)
- A kernel can be run directly by **exec** function or can be **enqueued** as a task.
- If there are unused number of dimensions in workdiv; use 1, for that dimension.
`auto blockThreadExtent = alpaka::Vec<TDim3D,Idx>{1u,1u,128u};`
- Vendor specific profiling and debugging tools can be used directly on compiled alpaka (e.g. `nsys`, `rocprof` ...)
- If you debug GPU code try to compile your code for CPU; and use CPU debugger tools
(Change acc type to CPU accelerators then debug using `gdb` and similar tools.)

alpaka Performance

Mean GPU and CPU utilization on NVIDIA A40 GPU



Peak memory usage on NVIDIA A40 GPU



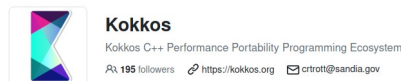
Source: Evaluating Performance Portability with the CMS Heterogeneous Pixel Reconstruction code

N. Andriotis¹, A. Bocci², E. Cano², L. Cappelli³, M. Dewing⁴, T. Di Pilato^{5,6}, J. Esseiva⁷, L. Ferragina⁸, G. Hugo², M. Kortelainen⁹, M. Kwok⁹, J. J. Olivera Loyola¹⁰, F. Pantaleo², A. Perego¹¹, W. Redjeb^{2,12}

¹BSC ²CERN ³INFN Bologna ⁴ANL ⁵CASUS ⁶University of Geneva ⁷LBNL ⁸University of Bologna ⁹FNAL ¹⁰ITESM ¹¹University of Milano Bicocca ¹²RWTH

CHEP 2023

https://indico.jlab.org/event/459/contributions/11824/attachments/9281/14171/20230511-CHEaP23_CMSPortability.pdf



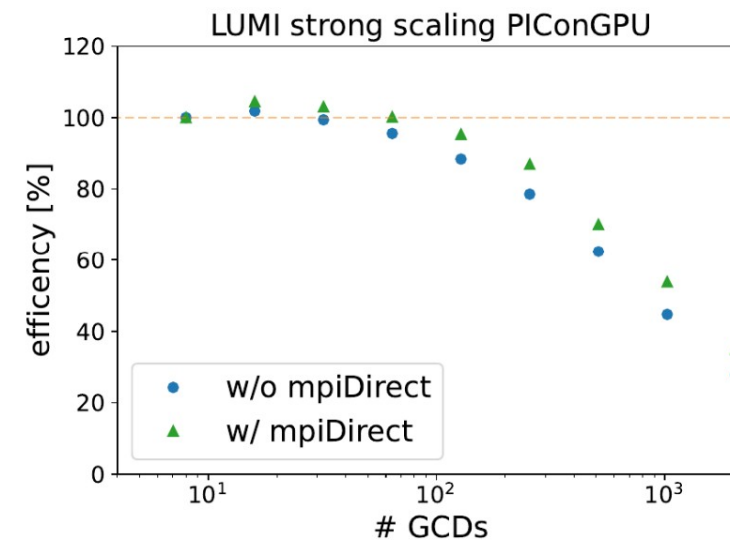
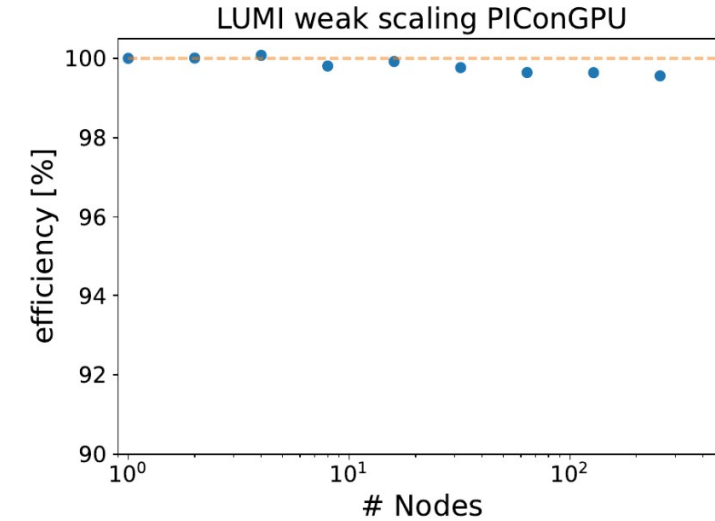
alpaka in the wild...

PIConGPU:

- Fully relativistic, manycore, 3D3V particle-in-cell (PIC) code
- Implements central algorithms in plasma physics
- Scalable to more than 18,000 GPUs
- Developed at Helmholtz-Zentrum Dresden-Rossendorf



<https://github.com/ComputationalRadiationPhysics/picongpu>



How to start using alpaka

- Don't write code initially on **cuda** because **alpaka** is already low level!
- **Use alpaka directly** by using examples and the cheat-sheet.
- **BUT if you already have a codebase in cuda**, converting to **cupla** can be a fast solution to benefit from alpaka features! Cupla is a member of alpaka group of softwares.

cupla - C++ User interface for the Platform Independent Library



<https://github.com/alpaka-group/cupla>

Cuda to portable C++ code

- Change the suffix `*.cu` of the CUDA source files to `*.cpp`
- Remove `#include <cuda_runtime.h>` and other cuda specific include files.
- Add `#include <cuda_to_cupla.hpp>`

Cuda

Kernel Function

```
template<int blockSize>

__global__ void fooKernel(int * ptr, float value)

{ // ... }
```

Kernel call at host

```
dim3 gridSize(42,1,1);
dim3 blockSize(256,1,1);
fooKernel<16><<< gridSize, blockSize, 0, 0>>>(ptr, 23);
```

Device function

```
template<typename TElem>
__device__ int deviceFunction(TElem x)
{ // ... }

// call

auto result = deviceFunction(x);
```

Shared memory

```
__shared__ int foo;

__shared__ int fooCArray2D[4][32];
```

Cupla

Kernel Functor

```
template<int blockSize>
struct fooKernel {
    template<typename TAcc>
    ALPAKA_FN_ACC void operator()(TAcc const & acc, int *
    const ptr, float const value) const

    { // ... }

};
```

Kernel call at host

```
dim3 gridSize(42,1,1); dim3 blockSize(256,1,1);

CUPLA_KERNEL(fooKernel<16>)(gridSize, blockSize, 0, 0)(ptr,
23);
```

Device function

```
template< typename TAcc, typename TElem >
ALPAKA_FN_ACC int deviceFunction( TAcc const & acc, TElem
x)
{ // ... }

// call

auto result = deviceFunction(acc, x);
```

Shared memory

```
sharedMem(foo, int);

sharedMem(fooCArray2D, cupla::Array< cupla::Array<int,4>, 32>);
```

Community and Long Term Support

- Partners using and contributing to alpaka



- alpaka is a part of Helmholtz Roadmap 2027-2034

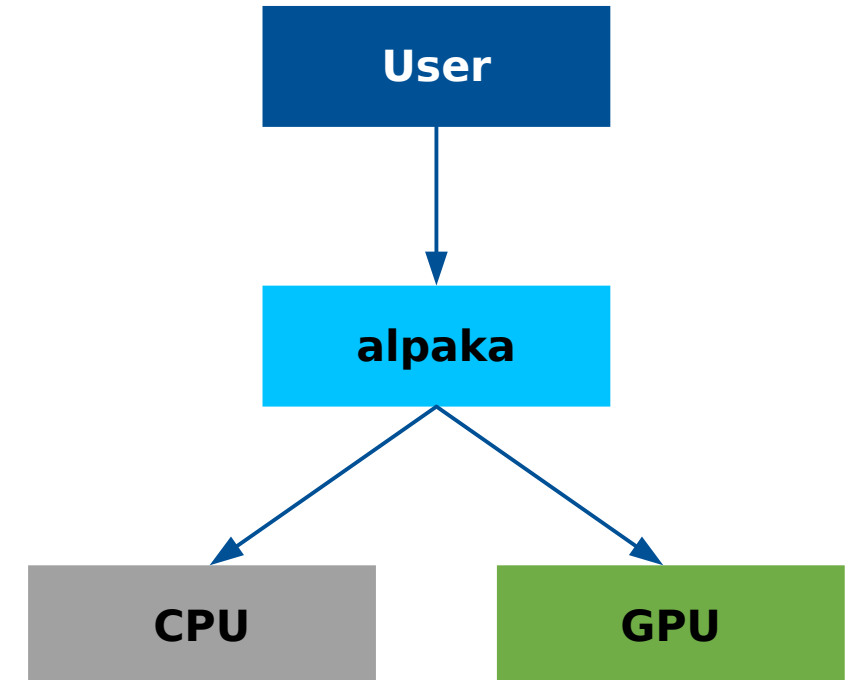
As a summary

Without alpaka

- In HPC world, Multiple hardware types are available from different vendors (CPUs, GPUs, ...)
- Increasingly heterogeneous hardware configurations available
- Platforms are not inter-operable → parallel programs are not easily portable

alpaka: one API to rule them all

- **Abstraction** (not hiding!) of the underlying hardware & software platforms
 - AMD, Nvidia, Intel GPUs, Different CPU parallelisations like TbbBlocks, OpenMP, Threads
- **Easy change of the backend in Code**
- **Built down to the same machine code with the vendor solutions**
- **Zero abstraction overhead for Kernel execution!**
- **Heterogenous Programming:** Using different backends in a synchronized manner



If you use alpaka for your research, please cite one of the following publications:

Matthes A., Widera R., Zenker E., Worpitz B., Huebl A., Bussmann M. (2017): Tuning and Optimization for a Variety of Many-Core Architectures Without Changing a Single Line of Implementation Code Using the alpaka Library. In: Kunkel J., Yokota R., Taufer M., Shalf J. (eds) High Performance Computing. ISC High Performance 2017. Lecture Notes in Computer Science, vol 10524. Springer, Cham, DOI: [10.1007/978-3-319-67630-2_36](https://doi.org/10.1007/978-3-319-67630-2_36).

E. Zenker et al., “alpaka – An Abstraction Library for Parallel Kernel Acceleration”, 2016 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW), Chicago, IL, 2016, pp. 631 – 640, DOI: [10.1109/IPDPSW.2016.50](https://doi.org/10.1109/IPDPSW.2016.50).

Worpitz, B. (2015, September 28). Investigating performance portability of a highly scalable particle-in-cell simulation code on various multi-core architectures. Zenodo. DOI: [10.5281/zenodo.49768](https://doi.org/10.5281/zenodo.49768).

Thank you!

You can contact us for any of your requests or questions about alpaka!

APPENDICES

Appendix-1

Summary of alpaka Structures

- **Accelerator** provides abstract view of all capable physical devices
- **Device** represents a single physical device
- **Queue** enables communication between the host and a single Device
- **Platform** is a union of Accelerator, Device and Kernel
- **Task** is a device-side operation (e.g kernel, memory operation)
- Others: **Event**, **Buffer** (dynamic array), **Vector** (static array)
- **Question**: How is portability between back-ends achieved?

Appendix-2: Programming Heterogeneous Systems-I

How to use multiple backends in parallel?

- Acquire at least one Device per Accelerator
- Create one Queue per Device

```
// Define Accelerators
using AccCpu = alpaka::AccCpuOmp2Blocks<Dim, Idx>;
using AccGpu = alpaka::AccGpuCudaRt<Dim, Idx>;

// Acquire at least one Device per Accelerator
auto devCpu = alpaka::getDevByIdx<AccCpu>(0u);
auto devGpu = alpaka::getDevByIdx<AccGpu>(0u);

// Create one queue per device

using QueueProperty = alpaka::NonBlocking;
using QueueCpu = alpaka::Queue<AccCpu, QueueProperty>;
using QueueGpu = alpaka::Queue<AccGpu, QueueProperty>;
auto queueCpu = QueueCpu{devCpu};
auto queueGpu1 = QueueGpu{devGpu};
auto queueGpu2 = QueueGpu{devGpu};

// Run tasks in parallel
alpaka::enqueue(queueCpu, taskCpu);
alpaka::enqueue(queueGpu1, taskGpu1);
alpaka::enqueue(queueGpu2, taskGpu2);
// Make sure all non-blocking queue tasks finished
// before the main thread ends
alpaka::wait(queueCpu);
alpaka::wait(queueGpu1);
alpaka::wait(queueGpu2);
```

Appendix-3

Programming Heterogeneous Systems-II

Communication by Buffers

- Buffers are defined and created per Device
- Buffers can be copied between different Devices
- **Notice:** CPU to GPU copies (and vice versa) require GPU queue

```
// Allocate buffers
auto bufCpu = alpaka::allocBuf<float, Idx>(devCpu, extent);
auto bufGpu = alpaka::allocBuf<float, Idx>(devGpu, extent);

/* Initialization ... */

// Copy buffer from CPU to GPU - destination comes first
alpaka::memcpy(gpuQueue, bufGpu, bufCpu, extent);

// Execute GPU kernel
alpaka::enqueue(gpuQueue, someKernelTask);

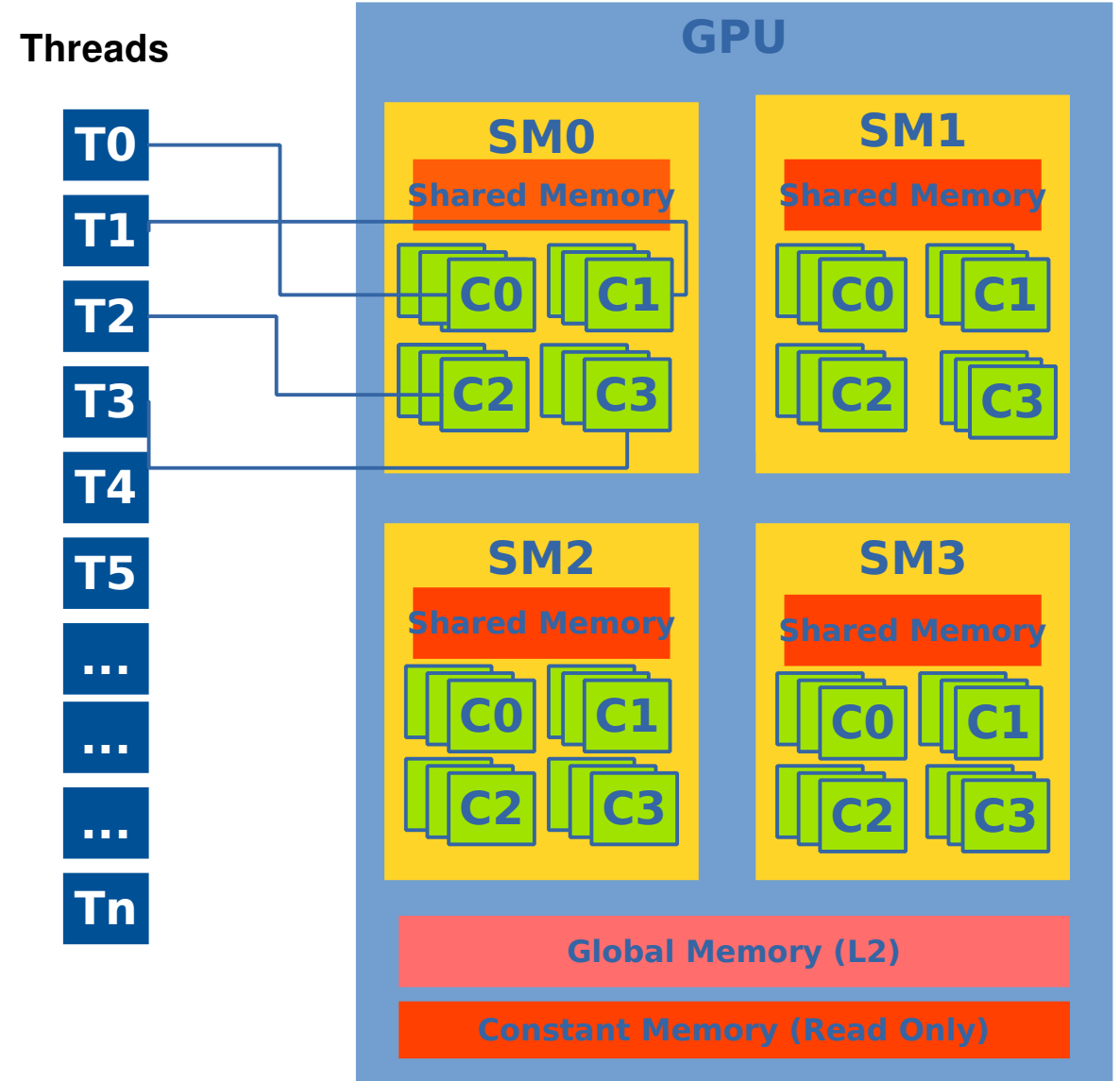
// Copy results back to CPU and wait for completion
alpaka::memcpy(gpuQueue, bufCpu, bufGpu, extent);

// Wait for GPU, then execute CPU kernel
alpaka::wait(cpuQueue, gpuQueue);
alpaka::enqueue(cpuQueue, anotherKernelTask);
```

Appendix-4: Determining the WorkDivision-I

How to distribute threads among SMs?

- Main determinants of mapping threads to the SMs:
 - Number of cores SM,
 - Warpsize,
 - Register and local memory (lm) usage of each thread,
 - Shared memory usage of each thread,
 - Threads per SM, Threads per Block, Blocks per Device
- Memory latencies: Global Memory and Constant memory has different latencies.
- Memory sizes: Size of shared memory used by threads in a block or blocks assigned to an SM



Appendix-5

1D WorkDivision

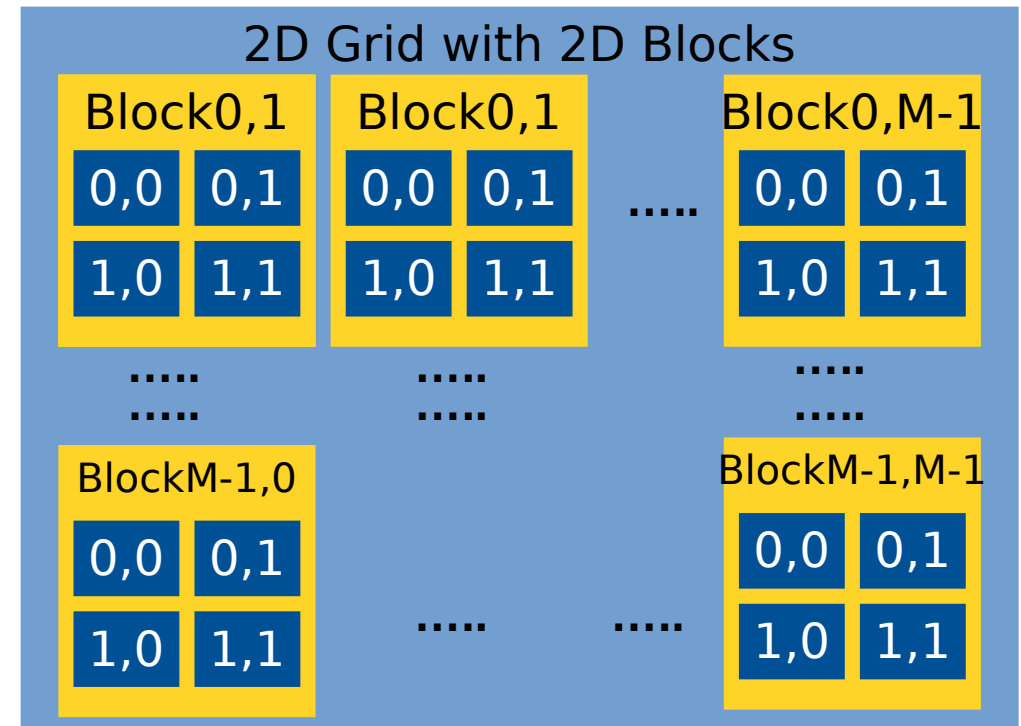
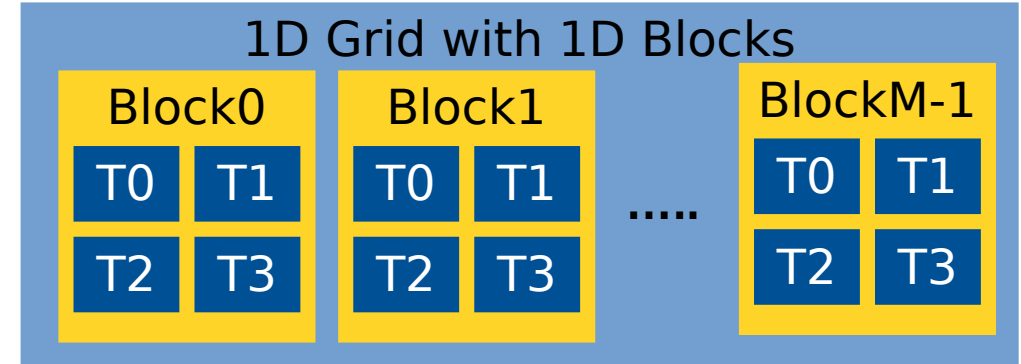
Grid Block Extent 1D vector = $\text{Vec}\{M\}$
 Block Thread Extent 1D vector = $\text{Vec}\{4u\}$ // CHANGE THE IMAGE!!!
 Thread Elem Extent 1D vector = $\text{Vec}\{1u\}$
 Grid Thread Extent 1D vector = $\text{Vec}\{4 * M\}$ or $\text{Vec}\{N\}$

1D WorkDivision $\{\{M\}, \{4u\}, \{1u\}\}$
 // if Dim is 3 then fill with 1u
 1D WorkDivision $\{\{1u, 1u, M\}, \{1u, 1u, 4\}, \{1u, 1u, 1u\}\}$

2D WorkDivision!

Grid Block Extent 2D vector = $\text{Vec}\{M, M\}$
 Block Thread Extent 2D vector = $\text{Vec}\{2u, 2u\}$
 Thread Elem Extent 2D vector = $\text{Vec}\{1u, 1u\}$
 Grid Thread Extent 2D vector = $\text{Vec}\{4u * M, 4u * M\}$ or $\text{Vec}\{N, N\}$

2D WorkDivision $\{\{M, M\}, \{2u, 2u\}, \{1u, 1u\}\}$
 2D WorkDivision $\{\{1u, M, M\}, \{1u, 2u, 2u\}, \{1u, 1u, 1u\}\}$

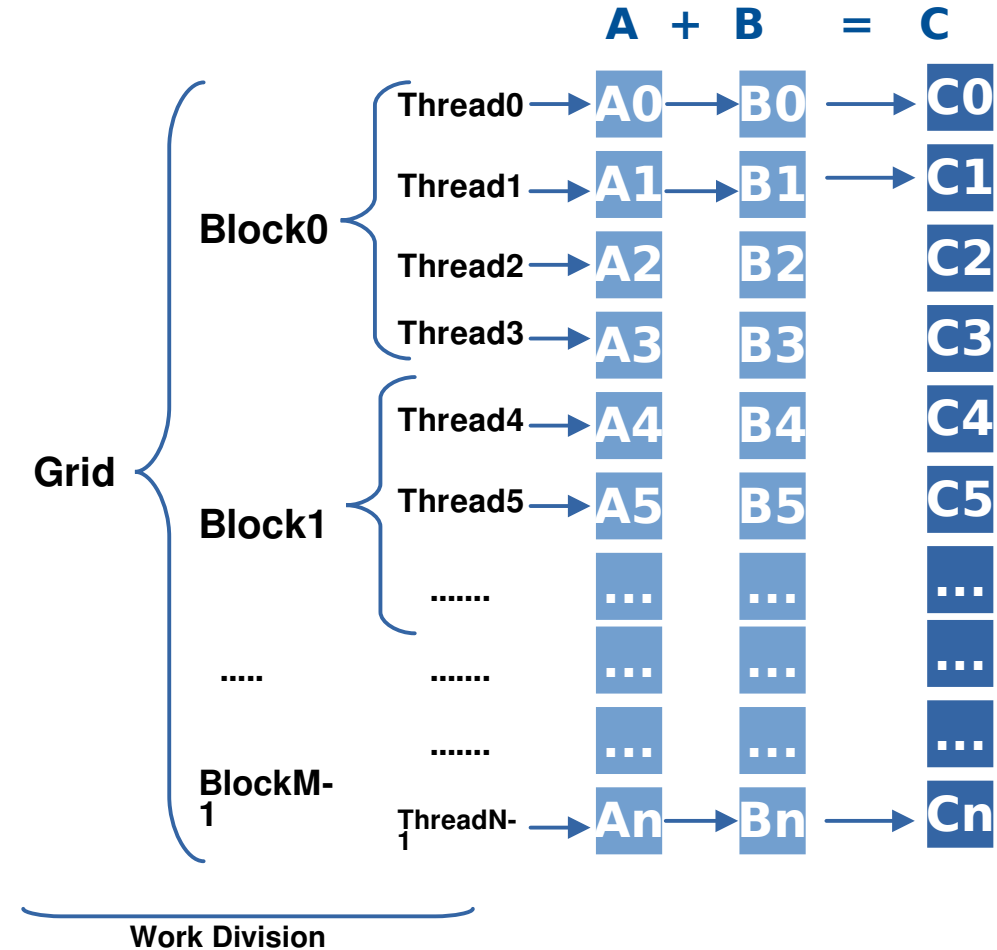


Appendix-6

Set the workdivision manually

- WorkDivision data structure consists 3 vectors:
 - Grid block extent.
`Vec{M}` or `Vec{1, 1, M}` depending on the number of dimensions.
 - Block thread extent.
`Vec{4}` or `Vec{1, 1, 4}`
 - Elements per thread
- Setting work div manually

```
using Dim1D = alpaka::DimInt<1>; //Set number of dims to 1
using Vec1D = alpaka::Vec<Dim1D, Idx>; //Define alias
auto workDiv1D = alpaka::WorkDivMembers(Vec1D{M}, Vec1D{4u}, Vec1D{1u});
// alternatively
using Dim3D = alpaka::DimInt<3>; //Set number of dims to 3
using Vec3D = alpaka::Vec<Dim3D, Idx>; //Define alias
auto workDiv3D = alpaka::WorkDivMembers(Vec3D{1,1,M}, Vec3D{1,1,4u}, Vec3D{1,1,1u});
```



Appendix-7 Tasks and Events

- Device-side operations (kernels, memory operations) can be wrapped in tasks.
- Tasks are executed by **enqueue()** function.
- Tasks on the same queue are executed in order (FIFO principle)


```
alpaka::enqueue(queueA, task1);
alpaka::enqueue(queueA, task2); // task2 starts
after task1 has finished, even queueA is non-blocking
```
- Order of tasks in different queues is unspecified
 - ```
alpaka::enqueue(queueA, task1);
```

```
alpaka::enqueue(queueB, task2); // task2 starts
before, after or in parallel to task1
```
- For easier synchronization, alpaka Events can be inserted before, after or between Tasks:

```
auto mvEvent =
alpaka::Event<alpaka::Queue>(mvDev);
alpaka::enqueue(queueA, mvEvent);

alpaka::wait(queueB, mvEvent); // queueB will only
resume after queueA reached mvEvent
```

```
// Create a queue on the device
QueueAcc queue(devAcc);

....
// Instantiate the kernel function object
VectorAddKernel kernel;

// Create the kernel execution task.
auto const taskKernel = alpaka::createTaskKernel<Acc>(
 workDiv,
 kernel,
 alpaka::getPtrNative(bufAccA),
 alpaka::getPtrNative(bufAccB),
 alpaka::getPtrNative(bufAccC),
 numElements);

alpaka::enqueue(queue, taskKernel);
alpaka::wait(queue); // wait in case we are using an asynchronous queue
```



## APPENDIX-8

# Accelerator Details

- Accelerator chosen by the programmer and **hides hardware specifics** behind alpaka's abstract API

```
using Acc = acc::AccGpuCudaRt<Dim, Idx>;
```

- **Inside Kernel:** contains thread state, provides access to alpaka's device-side API

- **The Accelerator provides the means to access to the indices**

```
// get thread index on the grid
auto gridThreadIdx = alpaka::getIdx<Grid, Threads>(acc);
// get block index on the grid
auto gridBlockIdx = alpaka::getIdx<Grid, Blocks>(acc);
```

- **The Accelerator gives access to alpaka's shared memory** (for threads inside the same block)

```
// allocate a variable in block shared static memory
auto & mySharedVar = block::shared::st::allocVar<int, __COUNTER__>(acc);
```

```
// get pointer to the block shared dynamic memory
float * mySharedBuffer = block::shared::dyn::getMem<float>(acc);
```

- **It also enables synchronization on the block level**

```
// synchronize all threads within the block
block::sync::syncBlockThreads(acc);
```

- **Internally, the accelerator maps all device-side functions to their native counterparts**

- Device-side functions require the accelerator as first argument:

```
math::sqrt(acc, /* ... */); time::clock(acc);
atomic::atomicOp<atomic::op::Or>(acc, /* ... */, hierarchy::Grids); (Atomics)
```

- **On Host:** Meta-parameter for choosing correct physical device and dependent types

## APPENDIX-9 Device information and device management

- Each alpaka **Device** represents a single physical device;
- Contains device information:
  - `auto const name = alpaka::getName(myDev);` // Back-end-defined device name
  - `auto const bytes = alpaka::getMemBytes(myDev);` // Size of device memory
  - `auto const free = alpaka::getFreeMemBytes(myDev);` // Size of available device memory
- Provides the means for device management:
  - `alpaka::reset(myDev);` // Reset GPU device state
- Encapsulates back-end device:
  - `auto nativeDevice = alpaka::getDev(myDev);` // nativeDevice is not portable!

## APPENDIX-10 Queue operations

- Queues execute Tasks (see next slide):

- `alpaka::enqueue(myQueue, taskRunKernel);`

- Check for completion:

- `bool done = alpaka::empty(myQueue);`

- Wait for completion, Events (see next slide), or other Queues:

- `alpaka::wait(myQueue);` // blocks caller until all operations have completed

- `alpaka::wait(myQueue, myEvent);` // blocks myQueue until myEvent has been reached

- `alpaka::wait(myQueue, otherQueue);` // blocks myQueue until otherQueue's ops have completed

## Appendix11 - Changing the target platform by changing accelerator

```
using namespace alpaka;

using Dim = dim::DimInt<1u>;
using Idx = std::size_t;

/** BEFORE **/
using Acc = alpaka::AccCpuOmp2Blocks<Dim, Idx>;

/** AFTER **/
using Acc = alpaka::AccGpuHipRt<Dim, Idx>;

/* No change required - dependent types and variables are automatically changed */
auto myDev = alpaka::getDevByIdx<Acc>(0u);

using Queue = alpaka::Queue<Acc, queue::NonBlocking>;
auto myQueue = Queue{myDev};
```