



# Alpaka Parallel Programming Library In a Nutshell

## 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 H100 GPU
- 4 Fugaku (Japan) 442 Petaflop/s, Fujitsu A64FX CPU
- 5 Lumi (Finland) 380 Petaflop/s, AMD EPYC CPU + AMD Instinct GPU
- 6 Leonardo (Italy) 238.7 Petaflop/s, Intel Xeon CPU + NVIDIA A100 GPU
- 7 MareNostrum (Spain) 183.2 Petaflops/s, Intel Xeon CPU + NVIDIA H100 GPU
- 8 Summit(USA) 148.3 Petaflops/s, IBM Power9 CPU + NVIDIA Volta GPU
- 9 Eos(USA) 121.4 Petaflops/s, Intel Xeon CPU + NVIDIA H100 GPU
- 10 Sierra(USA) 94 Petaflops/s, IBM POWER9 CPU + NVIDIA Tesla V100 GPU

### THE LIST

[www.top500.org](http://www.top500.org)

#### 11/2023 Highlights

g

The 62<sup>nd</sup> 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

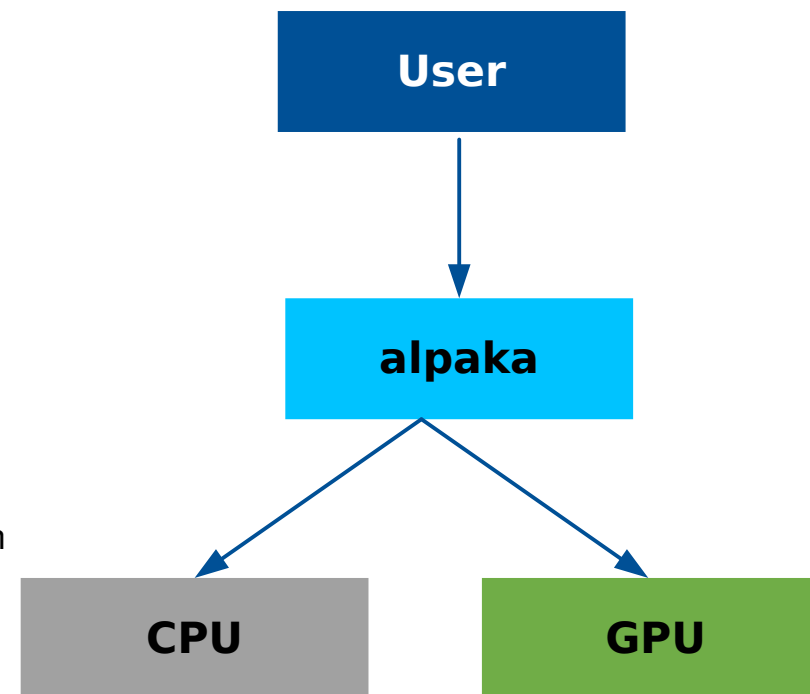
7 **Summit** - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-

## 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 platforms
  - No default device, default queue, built-in functions, language extensions
- **Easy change of the backend**
  - Code needs only minor adjustments to support different accelerators
- **Direct usage of vendor APIs, not depend on “unified APIs”**
  - GPU Backends: Hip backend uses HIP API, Cuda backend uses Cuda API, SYCL backend can be used for Intel GPUs
  - Cuda profilers and HIP profilers can be used to profile Alpaka generated kernel code.
  - CPU Backends: OpenMp, Threads, TbbBlocks
- **Easy indexing of threads in kernels**
- **Alpaka proposes and validates the type of parallelism**
- **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/>

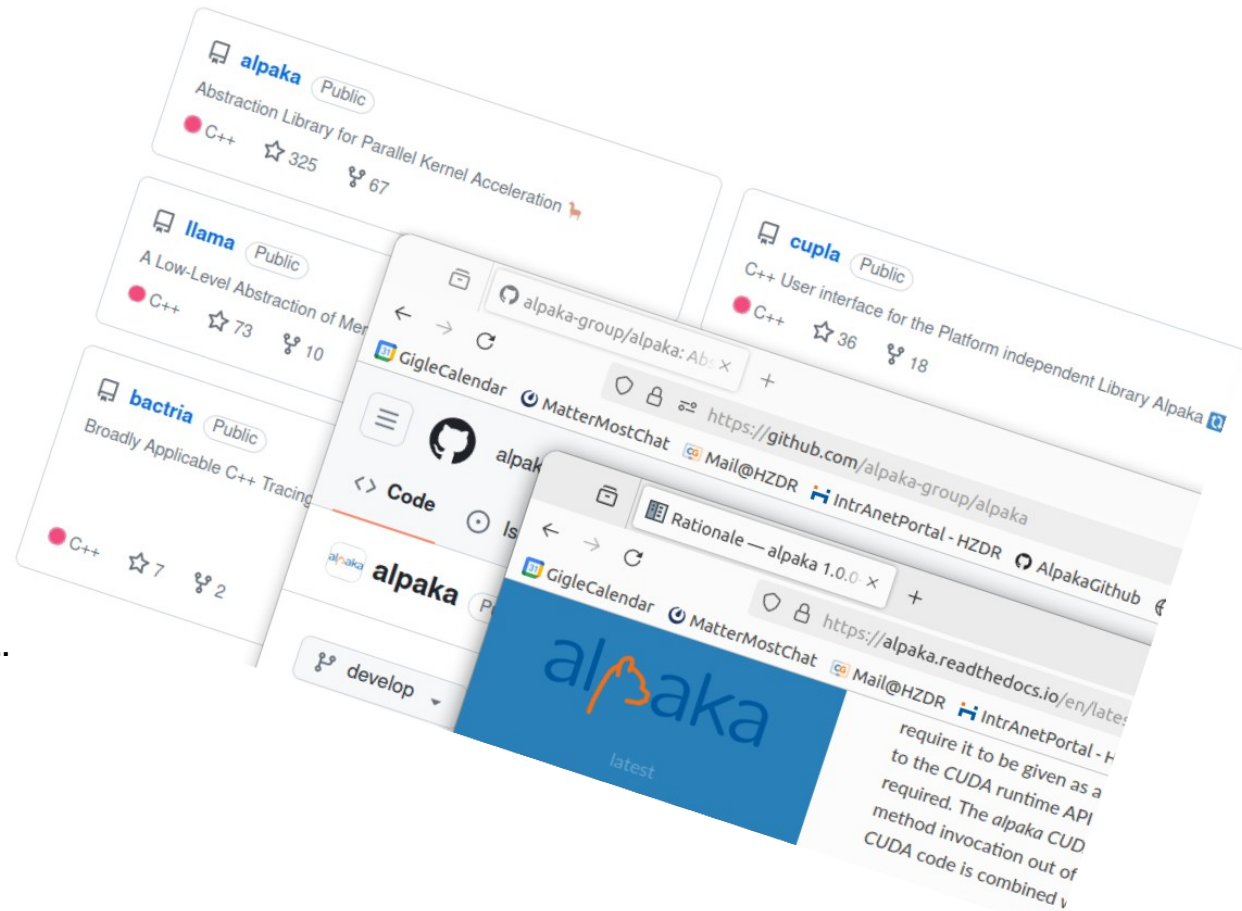
- Installation guide
- Cheatsheet
- Abstraction model and the rationale behind alpaka

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, MS VS)
- Portable across operating systems: Linux, macOS, Windows are supported





## Installation and building examples\*

### Install Alpaka 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>)
- 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/>
- Download alpaka:** `git clone -b develop https://github.com/alpaka-group/alpaka.git`
- In the terminal, 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/ ..`
- Optionally use ccmake tool to show/change the cmake variables.  
`ccmake ../`
- Select your backends from the ccmake parameters list and configure again. e.g turn ON  
`alpaka_ACC_CPU_B_SEQ_T_SEQ_ENABLE` and `alpaka_ACC_GPU_CUDA_ENABLE` for CPU and Cuda backends.

### Two options: Install alpaka or use alpaka from your build tree

- Install Alpaka:** Your alpaka installation will later reside in `/some/other/path/`
- Execute the installation(note the dot):  
`cmake --install .`
- You should now have a complete alpaka installation in the directory you chose earlier.

(\*)Detailed information: <https://github.com/alpaka-group/alpaka-workshop-slides/tree/develop>

```
BUILD_TESTING ON
Boost_ATOMIC_LIBRARY_RELEASE /usr/lib/x86_64-linux-gnu/libboost_atomic.so.1.74.0
Boost_INCLUDE_DIR /usr/include
CMAKE_BUILD_TYPE Release
CMAKE_CUDA_ARCHITECTURES 52
CMAKE_INSTALL_PREFIX /install
FETCHCONTENT_BASE_DIR /home/yusuf081/projects/alpaka-dir/alpaka/build/_deps
FETCHCONTENT_FULLY_DISCONNECTED OFF
FETCHCONTENT_QUIET ON
FETCHCONTENT_SOURCE_DIR_MDSPAN
FETCHCONTENT_UPDATES_DISCONNECTED OFF
FETCHCONTENT_UPDATES_DISCONNECTED OFF
MDSPAN_CXX_STANDARD DETECT
MDSPAN_ENABLE_BENCHMARKS OFF
MDSPAN_ENABLE_COMP_BENCH OFF
MDSPAN_ENABLE_CONCEPTS ON
MDSPAN_ENABLE_CUDA OFF
MDSPAN_ENABLE_EXAMPLES OFF
MDSPAN_ENABLE_OPENMP ON
MDSPAN_ENABLE_TESTS OFF
MDSPAN_USE_SYSTEM_GTEST OFF
RT_LIBRARY /usr/lib/x86_64-linux-gnu/librt.a
TBB_DIR /usr/lib/x86_64-linux-gnu/cmake/TBB
alpaka_ACC_CPU_B_OMP2_T_SEQ_EN OFF
alpaka_ACC_CPU_B_SEQ_T_OMP2_EN OFF
alpaka_ACC_CPU_B_SEQ_T_SEQ_ENA ON
alpaka_ACC_CPU_B_SEQ_T_THREADS OFF
alpaka_ACC_CPU_B_TBB_T_SEQ_ENA ON
alpaka_ACC_CPU_DISABLE_ATOMIC_ OFF
alpaka_ACC_GPU_CUDA_ENABLE ON
alpaka_ACC_GPU_CUDA_ONLY_MODE OFF
alpaka_ACC_GPU_HIP_ENABLE OFF
alpaka_ACC_GPU_HIP_ONLY_MODE OFF
alpaka_ACC_SYCL_ENABLE OFF
alpaka_ASSERT_ACC_ENABLE ON
alpaka_BLOCK_SHARED_DYN_MEMBER 47
alpaka_BUILD_BENCHMARKS ON
alpaka_BUILD_EXAMPLES ON
alpaka_CHECK_HEADERS OFF
alpaka_CUDA_EXPT_EXTENDED_LAMB ON
alpaka_CUDA_KEEP_FILES DEFAULT
alpaka_CUDA_SHOW_CODELINES OFF
alpaka_CUDA_SHOW_REGISTER DEFAULT
alpaka_CXX_STANDARD 17
```

## Compiling and running examples

- You can build all examples at once from your build directory:  
configure the build with setting some cmake variables  
`cmake -Dalpaka_BUILD_EXAMPLES=ON -DCMAKE_BUILD_TYPE=Release ..`  
build the examples  
`cmake --build . --config Release`  
alpaka/build/example/ directory will include compiled examples.  
alpaka/build/example/vectorAdd directory will include the executable `vectorAdd`
- Run all examples from the build directory of Alpaka  
`ctest example/`
- Run all tests  
`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 )

```
git clone https://github.com/alpaka-group/alpaka.git
cd alpaka/
mkdir build
cd build/
cmake -Dalpaka_BUILD_EXAMPLES=ON -DCMAKE_BUILD_TYPE=Release ..
// use ccmake to set other variables
cmake --build . --config Release
ctest example
./example/vectorAdd/vectorAdd
```



## Tenets of Thread-Parallel Programming

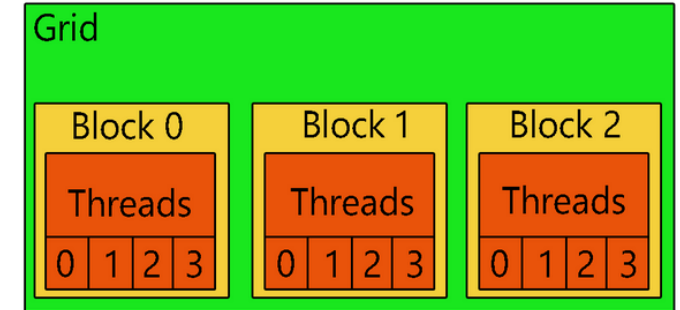
- **Massive number of threads.** Large number of threads should run the same code separately! This code is called “Kernel”.
- **Indexing of threads.** Each thread should work on a different data portion or do a specific task. Hence, each thread has an index, inside the kernel this index is accessible.
- **Grouping threads to map them to the cores.** Threads are grouped in blocks *to map them to the cores to increase occupancy*.

The warpsize of the GPUs and memory(shared memory,registers) usage of kernel affects the performance. Therefore grid of threads are divided into thread blocks.

**Extent:** A vector representing the sizes along each dimension. In 3d, a vector of 3 items {height,width,depth}

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

**Number of Dimensions**



**Grid Block Extent: Vec{3}**  
**Block Thread Extent: Vec{4}**



**Grid Block Extent: Vec{1,3}**  
**Block Thread Extent: Vec{5,3}**

# Vector Addition Problem

-Fix the number of data dimensions and the index type

-Create A and B vectors on Host  
Use dynamic array **alpaka::Buf**

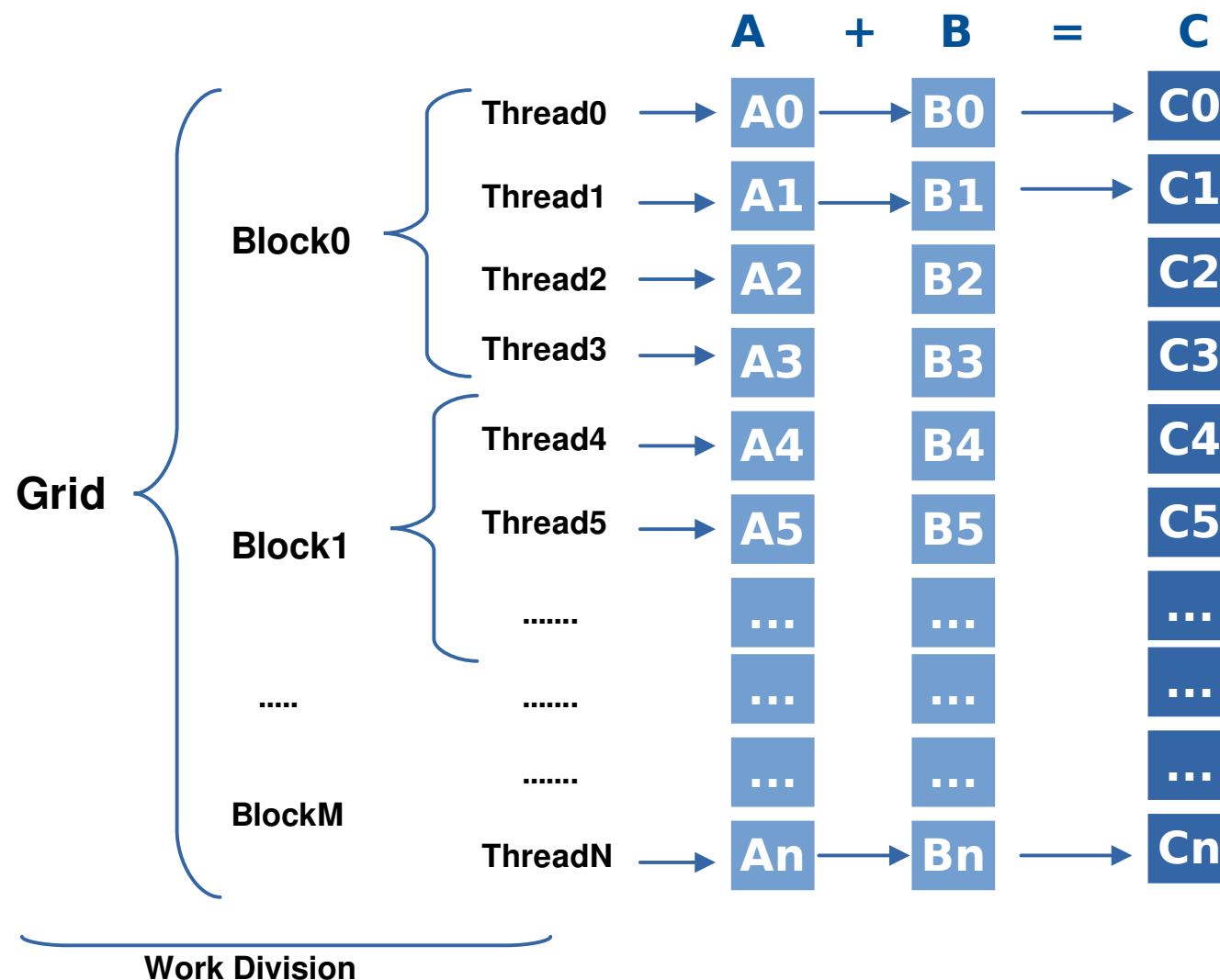
-Create empty C vector on Host  
Use **alpaka::Buf**

-Copy A and B Vectors to Device  
Use **alpaka::memcpy**

-Determine the Parallelism: Work Division.  
Extents of threads per block, blocks per grids, elements per thread.

-Create the Queue  
Use **alpaka::Queue**

-Execute the kernel on device with the workdiv on the queue. Pass the A,B,C vectors to the executer. Use **alpaka::exec**



# Parallel vector addition code

- Define kernel
- Select the backend depending on your system, the accelerator!
- Select the device (e.g your specific GPU)
- Decide on proper WorkDivision
- Create input vectors A and B and result vector C on host (e.g CPU)
- Create input vectors A and B and result vector C on device (e.g CPU)
- Copy the vectors to the device memory
- Run kernel
- Copy the result vector back to the host vector

```
// 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, 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!
    alpaka::wait(queue);
}
```

## Define the Alpaka Kernel

- Contains the algorithm
- Written on per-data-element basis
- It allows access to the thread index directly
- alpaka Kernels are functors (function-like C++ structs / classes)
- `operator()` is annotated with `ALPAKA_FN_ACC` specifier
- `operator()` must return `void`
- `operator()` must be `const`
- Arguments can be pointers and trivially copyable types
- Agnostic to device details
- Code piece that is run by each thread

```
// 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 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];
        }
    }
};
```

## Obtaining the indices and extents in Kernel

- **alpaka provides API functions for obtaining indices**

- Index of Thread on the Grid: `alpaka::getIdx<alpaka::Grid, alpaka::Threads>(acc)[dim];`
- Index of Thread on a Block: `alpaka::getIdx<alpaka::Block, alpaka::Threads>(acc)[dim];`
- Index of Block on the Grid: `alpaka::getIdx<alpaka::Grid, alpaka::Blocks>(acc)[dim];`

- **alpaka provides API functions for obtaining extents of the grid and the blocks**

- Number of Threads on the Grid: `alpaka::getWorkDiv<alpaka::Grid, alpaka::Threads>(acc)[dim];`
- Number of Threads on a Block: `alpaka::getWorkDiv<alpaka::Block, alpaka::Threads>(acc)[dim];`
- Number of Blocks on the Grid: `alpaka::getWorkDiv<alpaka::Grid, alpaka::Blocks>(acc)[dim];`



# alpaka in a Nutshell: Accelerator

## Select the Accelerator

- alpaka provides a number of pre-defined Accelerators.
- Accelerator is a type that is only instantiated on Device not on the host.
- Accelerator is an abstraction that makes the kernel 'device agnostic'.  
( Gives kernel means the access thread index. Responsible for the thread state, memory management, and synchronization ??)
- For GPUs:
  - AccGpuCudaRt for NVIDIA GPUs
  - AccGpuHipRt for AMD, Intel and NVIDIA GPUs
- For CPUs
  - AccCpuFibers based on Boost.fiber
  - AccCpuOmp2Blocks based on OpenMP 2.x
  - AccCpuOmp4 based on OpenMP 4.x
  - AccCpuTbbBlocks based on TBB
  - AccCpuThreads based on std::thread

```
// Easy change of the backend type!  
// Example: CUDA GPU accelerator  
using Acc = alpaka::AccGpuCudaRt<Dim, Idx>;
```

```
// Switch to AMD: HIP GPU accelerator  
using Acc = alpaka::AccGpuHipRt<Dim, Idx>;
```

```
// Switch to CPU Omp: CPU accelerator  
using Acc = alpaka::AccCpuOmp2Blocks<Dim, Idx>;
```

```
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  
    using DevHost = alpaka::DevCpu;  
    // Allocate 3 host memory buffers  
    using BufHost = alpaka::Buf<DevHost, DataType, Dim, Idx>;  
    ...  
    ...  
    // Instantiate the kernel function object  
    VectorAddKernel kernel;  
    alpaka::exec<Acc>( // Run the kernel execution task  
        queue,
```



## Get device instance

### alpaka Devices

- Each alpaka Device instance represents a single physical device
- Determined by programmer's Accelerator choice
- Easy management of physical devices
- Contains device information

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

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

auto const name = alpaka::getName(devAcc); // Back-end-defined device name
auto const bytes = alpaka::getMemBytes(devAcc); // Size of device memory
auto const free = alpaka::getFreeMemBytes(devAcc); // Size of available device memory
// Provides the means for device management:
alpaka::reset(devAcc);
```

```
// Select a device
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.
using DevHost = alpaka::DevCpu;
auto const platformHost = alpaka::PlatformCpu{};
auto const devHost = alpaka::getDevByIdx(platformHost, 0);

// 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));

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

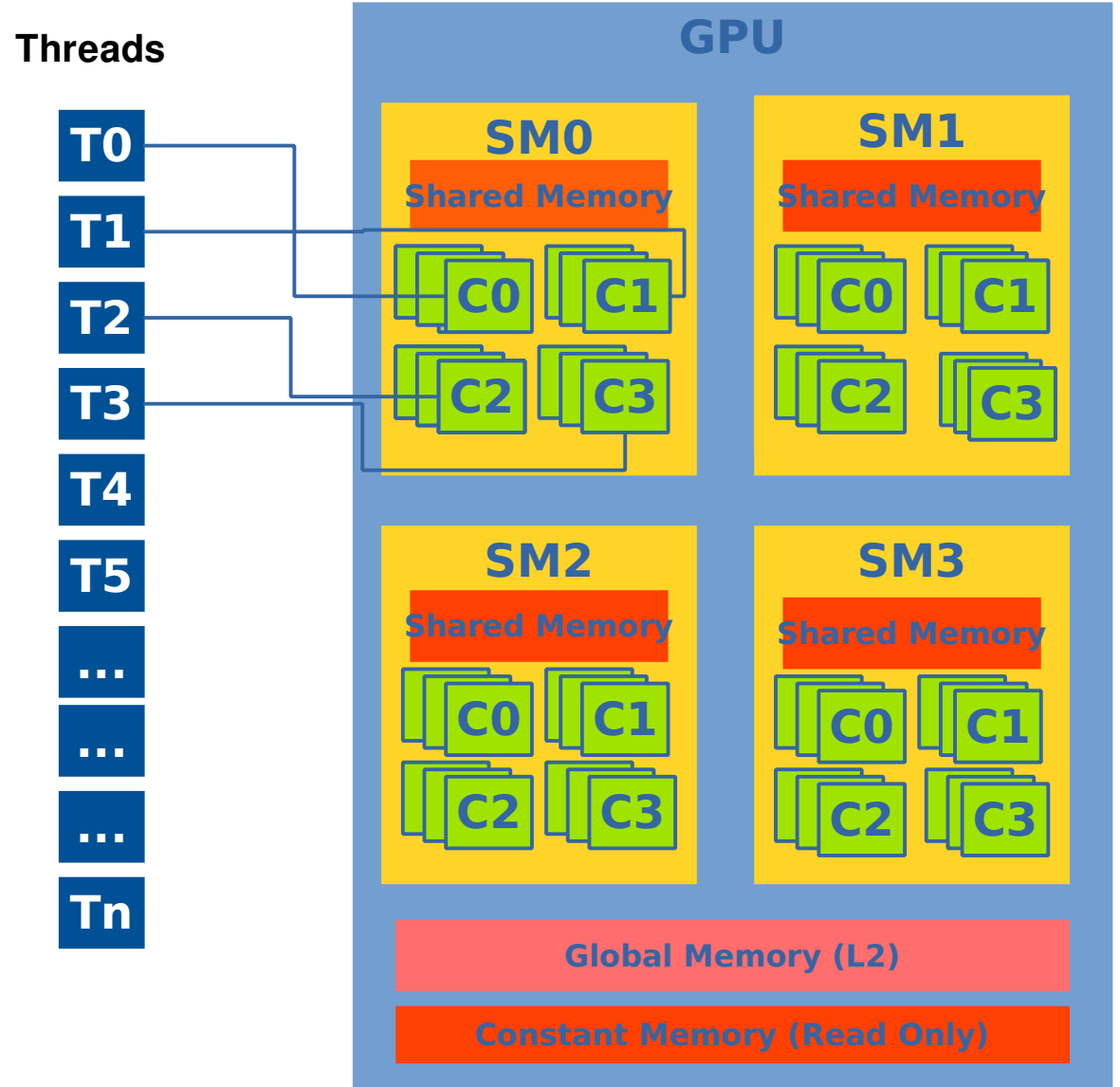
// Allocate 3 buffers on the accelerator
using BufAcc = alpaka::Buf<DevAcc, Data, Dim, Idx>;
BufAcc bufAccA(alpaka::allocBuf<Data, Idx>(devAcc, extent));
BufAcc bufAccB(alpaka::allocBuf<Data, Idx>(devAcc, extent));
BufAcc bufAccC(alpaka::allocBuf<Data, Idx>(devAcc, extent));

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

## Determine the WorkDivision-I

### How to distribute threads among SMs?

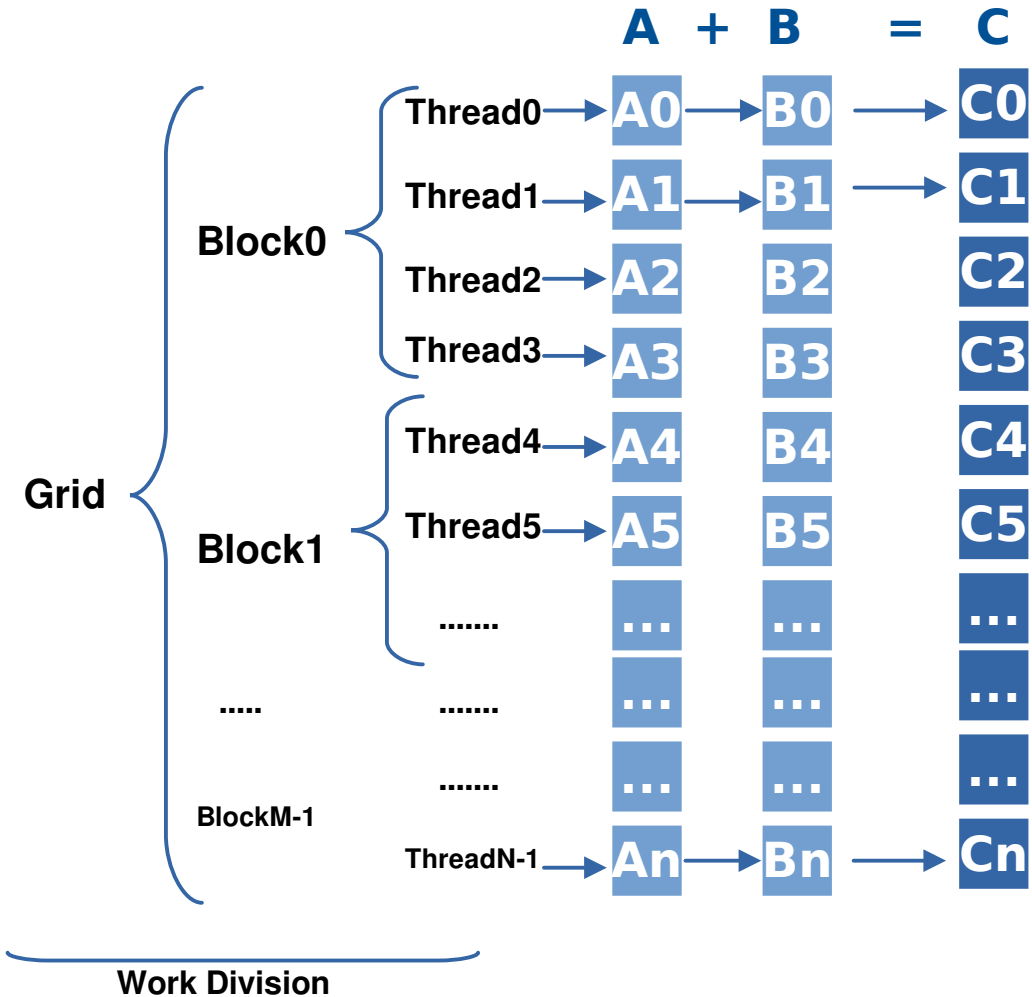
- Main determinants of mapping threads to the SMs:
  - Number of cores SM,
  - Warpsize,
  - Register and local memory (lm) usage of each thread,
  - Shared memory usage of each thread,
  - Threads per SM, Threads per Block, Blocks per Device
- Memory latencies: Global Memory and Constant memory has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



## Determine the workdivision-II

- WorkDivision data structure consists 3 vectors:
  - Grid block extent.  
 $\text{Vec}\{M\}$  or  $\text{Vec}\{1, 1, M\}$  depending on the number of dimensions.
  - Block thread extent.  
 $\text{Vec}\{4\}$  or  $\text{Vec}\{1, 1, 4\}$ .
  - 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});
```



## 1D WorkDivision!

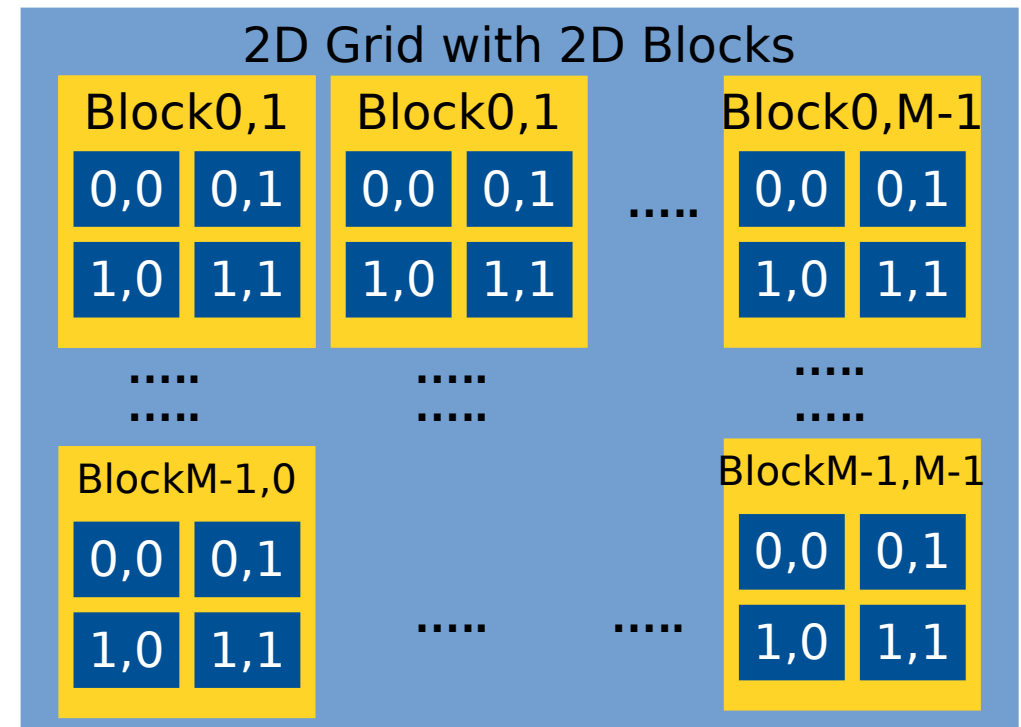
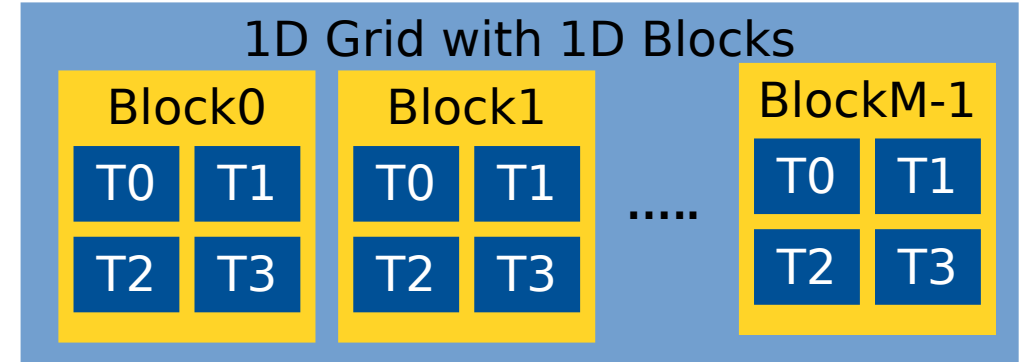
Grid Block Extent 1D vector =  $\text{Vec}\{M\}$   
 Block Thread Extent 1D vector =  $\text{Vec}\{4u\}$   
 Thread Elem Extent 1D vector =  $\text{Vec}\{1u\}$   
 Grid Thread Extent 1D vector =  $\text{Vec}\{4*M\}$  or  $\text{Vec}\{N\}$

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



## Determine the WorkDivision-III

- Use `getValidWorkDiv` function
- Inputs are
  - Full grid-thread extent. Hence we guarantee total number of threads needed.
  - Elements per thread extent
- `getValidWorkDiv` generates divides the full grid-thread extent to blocks of given dimension.

```
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);
}
```



## Allocate data vectors on host and device. Copy data vectors to the Device

- **alpaka::Buf** is multi dimensional dynamic array. Knows it's data type, index type, extent and device! Copyable among host and accelerator devices.

- **alpaka::allocBuf** allocates memory to the given device.

- **alpaka::memcpy** copies the data from one buffer/view to another. User should NOT have to specify which buffer resides on which device but should provide a queue.

Memory copy needs to be synchronized, therefore needs a queue

- How to **copy an STL container** from a device to another?

1. Convert the container to **alpaka::view** by attaching a device!

```
std::array<Data, nElements> array;
```

```
auto hostViewPtr = alpaka::createView(devHost, array.data(), extents);
```

2. Copy using a **alpaka::memcpy**.

```
alpaka::memcpy(devQueue, deviceBuffer1, hostViewPtr);
```

```
// 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);
```



## Create the Queue for memcpy and kernel executions

- Alpaka::Queue provides communication between the host and the device
- Queue of tasks on the device
- Queue is always FIFO
- Create **alpaka::Queue** using the device instance and the accelerator type (e.g GPU)
- Two queue types: blocking and non-blocking
- Blocking queues block the Host until Device-side command returns.
- Non-blocking queues return control to Host immediately, Device-side command runs asynchronously

```
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!
alpaka::wait(queue);
// bufHostC is filled with the result!
```

## Execute the kernel using the workdiv, copy result back to the Host

- **alpaka::exec** function is used
- Result is stored in an **alpaka::Buf**
- Copy the result in device to the host
- If the queue not blocking then we can use **alpaka::wait** on the queue for synchronization.

```
// 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!
alpaka::wait(queue);
```

## On Tasks and Events

- Device-side operations (kernels, memory operations) are called Tasks
- 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
```

- 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 myEvent = alpaka::Event<alpaka::Queue>(myDev);  
  
alpaka::enqueue(queueA, myEvent);  
alpaka::wait(queueB, myEvent); // queueB will only resume after queueA reached myEvent
```

## Programming Heterogeneous Systems-I

### Using multiple Platforms Synchronously

- Alpaka enables easy heterogeneous programming!
- Create one Accelerator per back-end
- 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 Devices
auto devCpu = alpaka::getDevByIdx<AccCpu>(0u);
auto devGpu = alpaka::getDevByIdx<AccGpu>(0u);

// Create Queues
using QueueProperty = alpaka::NonBlocking;
using QueueCpu = alpaka::Queue<AccCpu, QueueProperty>;
using QueueGpu = alpaka::Queue<AccGpu, QueueProperty>;

auto queueCpu = QueueCpu{devCpu};
auto queueGpu = QueueGpu{devGpu};
```

## Programming Heterogeneous Systems-II

### Communication by Buffers

- Buffers are defined and created per Device
- Buffers can be copied between different Devices / Queues
- Not restricted to a single platform!
- **Restriction:** 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);
```



## Matrix-Vector Multiplication

- using global memory
- using constant memory
- using using shared memory

```
static constexpr auto numberOfColumns = 5000u;
using Elem = std::uint32_t;
// const memory for vector allocation.
ALPAKA_STATIC_ACC_MEM_CONSTANT alpaka::DevGlobal<TAcc, Elem[1u][numberOfColumns]>
  g_constantMemory1DUninitialized;

class MatVecMulKernelUsingConstantMemory
{
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 C,
        TIdx const& numRows,
        TIdx const& numCol) 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 <= numRows)
        {
            TElem dotProd = 0;
            for(size_t j = 0; j < numCol; ++j)
            {
                dotProd += A[gridThreadId * numCol + j] *
  g_constantMemory1DUninitialized<TAcc>.get()[0][j];
            }
            C[gridThreadId] = dotProd;
        }
    }
};

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

    ....
    using BufHost = alpaka::Buf<DevHost, Data, Dim, Idx>;
    BufHost bufHostB(alpaka::allocBuf<Data, Idx>(devHost, extentVector));
    ....
    // Create a queue on the device
    alpaka::Queue<Acc, QueueProperty> queue(devAcc);
    ....
    // Used if the kernel using const memory is enqueued
    alpaka::memcpy(queue, g_constantMemory1DUninitialized<Acc>, bufHostB);
    alpaka::wait(queue);
}
```



## Tips

- 
- **If you want to use data more than 1D structures you need to take care of alignment/pitch values.**
- If there are unused number of dimensions in workdiv use 1, for that dimension.  
`auto threadElemExtent = alpaka::Vec<Tdim3D,Idx>{1,1,128};`
- CPU to GPU copies require a GPU Queue
- 2 queues synchronization

## 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?

## [ToBeChanged] CUPLA: Converting CUDA code to Alpaka

- Change the suffix \*.cu of the CUDA source files to \*.cpp
- Remove cuda specific includes on top of your header and source files
- Add **#include <cuda\_to\_cupla.hpp>** remove **#include <cuda\_runtime.h>**
- Transform kernels ( \_\_global\_\_ functions) to functors
- The functor's **operator()** must be qualified as **const**
- Add the function prefix **ALPAKA\_FN\_ACC** to the operator() const
- Add as first (templated) kernel parameter the accelerator with the name acc (it is important that the accelerator is named acc because all cupla-to-alpaka replacements use the variable acc)
- If the kernel calls other functions you must pass the accelerator **acc** to each call
- Add the qualifier **const** to each parameter which is not changed inside the kernel

### • CUDA kernel

```
template< int blockSize >
__global__ void fooKernel( int * ptr, float value )
{
    // ...
}
```

### Cupla kernel

```
template< int blockSize >
struct fooKernel
{
    template< typename T_Acc >
    ALPAKA_FN_ACC
    void operator()( T_Acc const & acc, int * const ptr, float const value) const
    {
        // ...
    }
};
```

### CUDA kernel call at host

```
dim3 gridSize(42,1,1);
dim3 blockSize(256,1,1);
// extern shared memory and stream is optional
fooKernel< 16 ><<< gridSize, blockSize, 0, 0 >>>( ptr, 23 );
```

### CUPLA kernel call at host:

```
dim3 gridSize(42,1,1);
dim3 blockSize(256,1,1);
// extern shared memory and stream is optional
CUPLA_KERNEL(fooKernel< 16 >)( gridSize, blockSize, 0, 0 )( ptr, 23 );
```

### CUDA shared memory (in kernel)

```
__shared__ int foo;
__shared__ int fooCArray[32];
__shared__ int fooCArray2D[4][32];
// extern shared memory (size was defined during the host side kernel call)
extern __shared__ float fooPtr[];
```

```
int bar = fooCArray2D[ threadIdx.x ][ 0 ];
```

### CUPLA shared memory (in kernel)

```
sharedMem( foo, int );
sharedMem( fooCArray, cupla::Array< int, 32 > );
sharedMem( fooCArray, cupla::Array< cupla::Array< int, 4 >, 32 > );
sharedMemExtern( fooPtr, float );
```

```
int bar = fooCArray2D[ threadIdx.x ][ 0 ]
```

### Use ALPAKA\_FN\_ACC in device function definitions and add an acc parameter.

(Note that to be exact the acc parameter is only necessary when alpaka functions like blockIdx or atomicMax etc are used.)

### • CUDA Device Function

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

### CUPLA

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

### • CUDA Device Function Call

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

### CUPLA

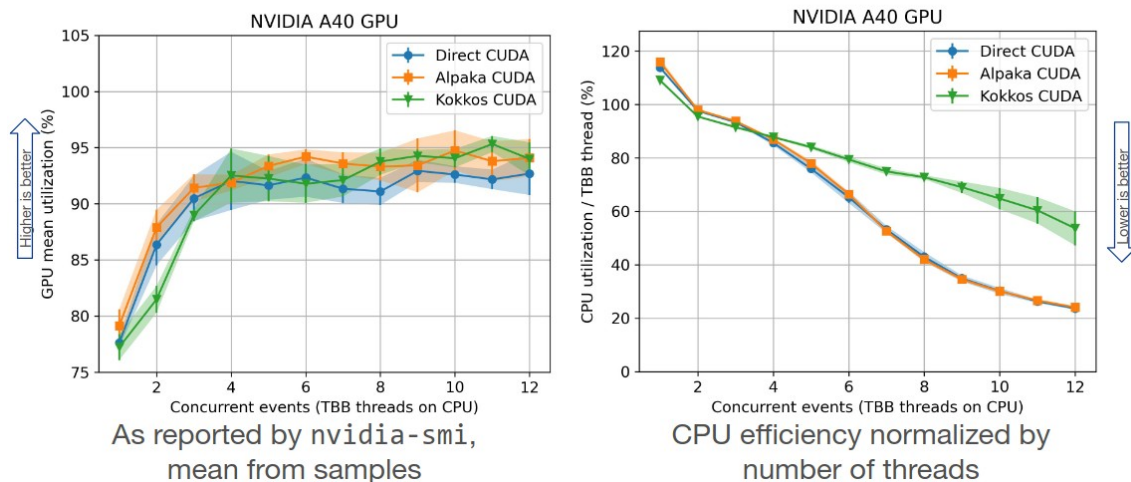
```
auto result = deviceFunction( acc, x );
```

## How to start using Alpaka

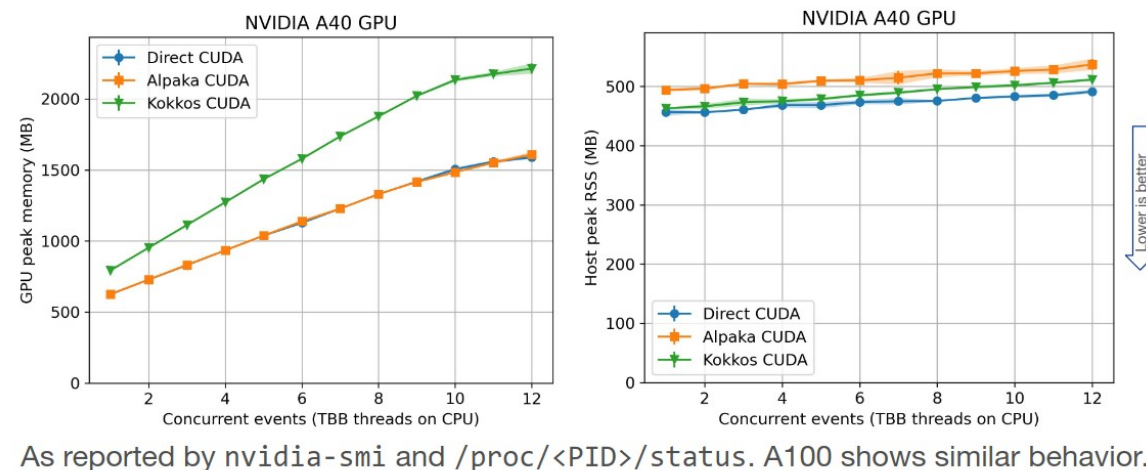
- Use Alpaka directly using examples and the cheat-sheet
- Don't write code initially on Cuda because alpaka is already low level, like Cuda.
- If you have a codebase in cuda, converting to Cupla can be a fast solution to benefit from alpaka features!

## 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. Andriotis<sup>1</sup>, A. Bocci<sup>2</sup>, E. Cano<sup>2</sup>, L. Cappelli<sup>3</sup>, M. Dewing<sup>4</sup>, T. Di Pilato<sup>5,6</sup>, J. Esseiva<sup>7</sup>, L. Ferragina<sup>8</sup>, G. Hugo<sup>2</sup>,

M. Kortelainen<sup>9</sup>, M. Kwok<sup>9</sup>, J. J. Olivera Loyola<sup>10</sup>, F. Pantaleo<sup>2</sup>, A. Perego<sup>11</sup>, W. Redjeb<sup>2,12</sup>

<sup>1</sup>BSC <sup>2</sup>CERN <sup>3</sup>INFN Bologna <sup>4</sup>ANL <sup>5</sup>CASUS <sup>6</sup>University of Geneva <sup>7</sup>LBNL <sup>8</sup>University of Bologna

<sup>9</sup>FNAL <sup>10</sup>ITESM <sup>11</sup>University of Milano Bicocca <sup>12</sup>RWTH

CHEP 2023

[https://indico.jlab.org/event/459/contributions/11824/attachments/9281/14171/20230511-CHEaP23\\_CMSPortability.pdf](https://indico.jlab.org/event/459/contributions/11824/attachments/9281/14171/20230511-CHEaP23_CMSPortability.pdf)

## Community and Long Term Support

- Partners using and contributing to Alpaka



**HZB** Helmholtz  
Zentrum Berlin

---

**HZDR**  
HELMHOLTZ ZENTRUM  
DRESDEN ROSSENDORF

- Alpaka is a part of Helmholtz Roadmap 2027-2034



## Alpaka in the wild...

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

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



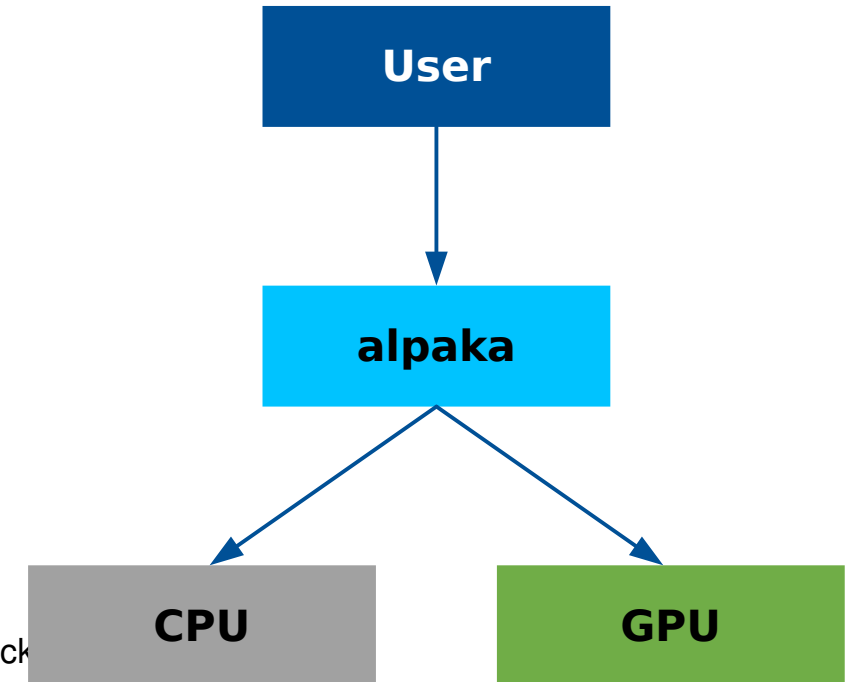
## As a summary

### Without alpaka

- Multiple hardware types are available from different vendors (CPUs, GPUs, ...)
- Increasingly heterogeneous hardware configurations available
- Platforms not inter-operable → parallel programs 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**
  - Code needs only minor adjustments to support different accelerators
- **Easy indexing of threads in kernels**
- **Easy setup of the type of parallelism by WorkDivision** (Block sizes in grid, Thread sizes in block)
- **Heterogenous Programming**: Using different backends in a synchronized manner



## Thank you! 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).

## Appendix

## [Appendix1] WorkDivision Example

- Determines the number of kernel instantiations
- Determines the type of parallelism
  - Dimensions of a grid in terms of blocks,
  - Dimensions of a block in terms of threads
  - Elements per thread

```
// Define the work division
// The workdiv is divided in three levels of parallelization:
// - grid-blocks:      The number of blocks in the grid
// - block-threads:    The number of threads per block (parallel, synchronizable).
// - thread-elements:  The number of elements per thread (sequential, not synchronizable)
//                      Each kernel has to execute its elements sequentially.

using Vec = alpaka::Vec<Dim, Idx>;
auto const elementsPerThread = Vec::all(static_cast<Idx>(1));
auto const threadsPerGrid = Vec{4, 2, 4};
using WorkDiv = alpaka::WorkDivMembers<Dim, Idx>;
WorkDiv const workDiv = alpaka::getValidWorkDiv<Acc>( devAcc, threadsPerGrid,
| elementsPerThread, false, alpaka::GridBlockExtentSubDivRestrictions::Unrestricted);

// Instantiate the kernel function object
HelloWorldKernel helloWorldKernel;

// Run the kernel
// To execute the kernel, you have to provide the
// work division as well as the additional kernel function parameters.
alpaka::exec<Acc>(queue, workDiv, helloWorldKernel/* put kernel arguments here */);
```

## [Appendix2] Matrix Vector Multiplication Example

-



## APPENDIX-5

- 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 gridThreadId = 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-6 Physical device information and management by “alpaka Device”

- 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-7 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

## Appendix8 - Changing the target platform

```
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};
```

## APPENDIX 9

### Acc versus Device

Acc is a host side compile time construct but instantiated in Kernel  
Abstraction that make the kernel device agnostic