### An introduction to alpaka

### performance portability with alpaka - 21<sup>st</sup> June 2023

# Andrea Bocci

CERN - EP/CMD



### who am I



- Dr. Andrea Bocci <andrea.bocci@cern.ch>, @fwyzard on Mattermost
  - applied physicist working on the CMS experiment for over 20 years
  - at CERN since 2010
  - I've held various roles related to the High Level Trigger
    - started out as the b-tagging HLT contact
    - joined as (what today is called) HLT STORM convener
    - deputy Trigger Coordinator and Trigger Coordinator
    - HLT Upgrade convener, and editor for the DAQ and HLT Phase-2 TDR
    - currently, "GPU Trigger Officer"
    - for the last 5 years, I've been working on GPUs and *performance portability* 
      - together with Matti and a few CERN colleagues
      - "Patatrack" pixel track and vertex reconstruction running on GPUs
      - R&D projects on CUDA, Alpaka, SYCL and Intel oneAPI
      - support for CUDA, HIP/ROCm, and Alpaka in CMSSW
      - Patatrack Hackathons !

#### June 21st, 2023





### summary

- this course will cover
  - what performance portability means
  - what is the Alpaka library
  - how to set up Alpaka for a simple project
  - how to compile a single source file for different back-ends
  - what are Alpaka platforms, devices, queues and events
  - how to work with host and device memory
  - how to write device functions and kernels
  - how to use an Alpaka accelerator and work division to launch a kernel
  - a complete example !



#### June 21st, 2023



## before we begin



- a C++17 compiler
  - GCC 9 installed on these machines will work fine
- the Boost libraries, version 1.74 or later
  - we have installed version 1.82 under /usr/local/boost/
- the alpaka headers
  - we have installed the current development version (as of June 20<sup>th</sup>, 2023) under /usr/local/alpaka/
- CUDA
  - CUDA 11.5 installed on these machines will work fine

to set it all up:

# disable the conda environment
conda deactivate

# alpaka requires Boost 1.74 or newer; use the prebuilt version 1.82 at
export BOOST\_BASE=/usr/local/boost

# use the development version of alpaka (June 20<sup>th</sup> 2023) at export ALPAKA\_BASE=/usr/local/alpaka

# clone the repository with the exercises on portability and alpaka
git clone https://github.com/fwyzard/intro\_to\_alpaka

June 21st, 2023



# before we begin



- a C++17 compiler
  - GCC 9 installed on these machines will work fine
- the Boost libraries, version 1.74 or later
  - we have installed version 1.82 under /usr/local/boost/
- the alpaka headers
  - we have installed the current development version (as of June 20<sup>th</sup>, 2023) under /usr/local/alpaka/

### CUDA

CUDA 11.5 installed on these machines will work fine

to set it all up:

# disable the conda environment
conda deactivate

# alpaka requires Boost 1.74 or newer; use the prebuilt version 1.82 at
export BOOST\_BASE=/usr/local/boost

# use the development version of alpaka (June 20<sup>th</sup> 2023) at export ALPAKA\_BASE=/usr/local/alpaka

# clone the repository with the exercises on portability and alpaka
git clone https://github.com/fwyzard/intro\_to\_alpaka

this part sets up the environment

make sure to do it in every session

June 21st, 2023

### performance portability



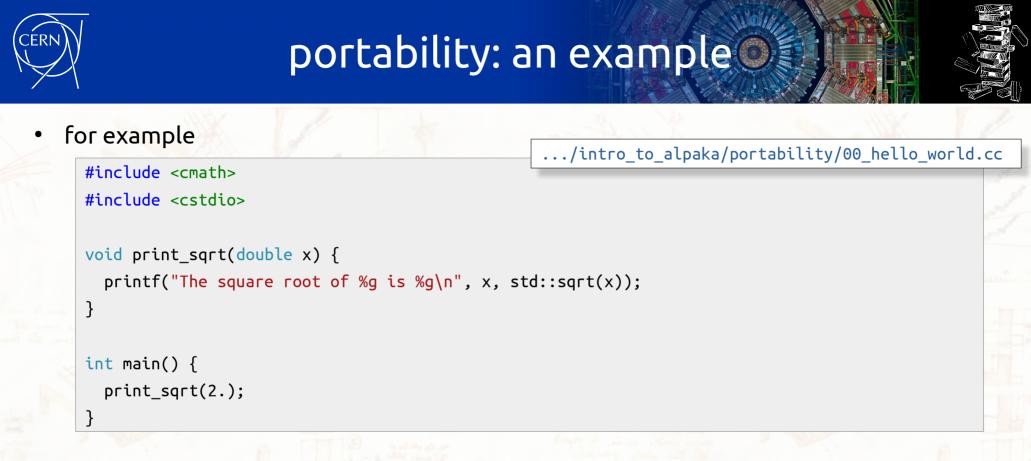
# what is *portability*?



- what do we mean by software *portability*?
  - the possibility of running a software application or library on different platforms
    - different hardware architectures, different operating systems
    - e.g. Windows running on x86, OSX running on ARM, Linux running on IBM Power, *etc*.
- how do we achieve software portability?
  - write software using a standardised language
    - C++, python, Java, *etc.*
  - use standard features
    - IEEE floating point numbers
  - use standard or portable libraries
    - C++ standard library, Boost, Eigen, *etc*.

#### June 21st, 2023





should behave in the same way on all platforms that support a standard C++ compiler:

The square root of 2 is 1.41421





### what about GPUs?



- writing a program that offloads some of the computations to a GPU is somewhat different from writing a program that runs just on the CPU
  - inside a single application ...
  - ... different hardware architectures
  - ... different memory spaces
  - ... different way to call a function or launch a task
  - ... different optimal algorithms
  - ... different compilers
  - ... different programming languages !
- sometimes it may help to think about a GPU like programming a remote machine
  - compile for completely different targets
  - launching a kernel is similar to running a complete program !





### portability: the same example

.../intro\_to\_alpaka/portability/01\_hello\_world.cu

```
#include <cmath>
#include <cstdio>
#include <cuda runtime.h>
__device__
void print sqrt(double x) {
  printf("The square root of %g is %g\n", x, std::sqrt(x));
__global__
void kernel() {
  print_sqrt(2.);
int main() {
  kernel<<<1, 1>>>();
  cudaDeviceSynchronize();
The square root of 2 is 1.41421
```

```
June 21st, 2023
```





# portability: side by side



11/78

CC

0

SA

<pre>#include <cmath></cmath></pre>	<pre>#include <cmath></cmath></pre>
<pre>#include <cstdio></cstdio></pre>	<pre>#include <cstdio></cstdio></pre>
	<pre>#include <cuda_runtime.h></cuda_runtime.h></pre>
<pre>void print_sqrt(double x) {</pre>	
<pre>printf("The square root of %g is %g\n", x, std::sqrt(x));</pre>	<pre>device</pre>
}	<pre>void print_sqrt(double x) {</pre>
	<pre>printf("The square root of %g is %g\n", x, std::sqrt(x));</pre>
<pre>int main() {</pre>	}
<pre>print_sqrt(2.);</pre>	<
}	global
TUTTERN. I STATE STATE	<pre>void kernel() {</pre>
The square root of 2 is 1.41421	<pre>print_sqrt(2.);</pre>
	<pre>int main() {</pre>
e could	kernel<<<1, 1>>>();
	<pre>cudaDeviceSynchronize();</pre>
<ul> <li>wrap the differences in a few macros or cla</li> </ul>	SSES }
<ul> <li>share the common parts</li> </ul>	
	The square root of 2 is 1.41421

June 21st, 2023



### so... are we done ?



- not really
  - trivially extending our example to an expensive computation would give horrible performance !
- why?
  - a CPU will run a single-threaded program very efficiently
  - a GPU would perform very badly
    - use a single thread out of a whole warp (32 threads): use *at most* 3% of its computing power
    - use a single block: loose any possibility of hiding memory latency
    - cannot take advantage of advanced capabilities like atomic operations, shared memory, *etc.*
  - and what about different GPU back-ends?
- what we need is *performance portability* 
  - write code in a way that can run on multiple platforms
  - leverage their potential
  - and achieve (almost) native performance on all of them

#### June 21st, 2023





### performance portability?



13/78



June 21st, 2023

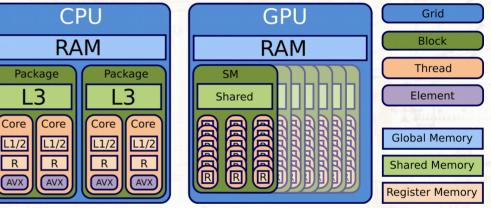
### the alpaka performance portability library



# what is alpaka?



- alpaka is a header-only C++17 abstraction library for accelerator development
  - it aims to provide *performance portability* across accelerators through the abstraction of the underlying levels of parallelism
- it currently supports
  - CPUs, with serial and parallel execution
  - GPUs by NVIDIA, with CUDA
  - GPUs by AMD, with HIP/ROCm
  - support for Intel GPUs and FPGAs is under development, based on SYCL and Intel oneAPI
- it is easy to integrate in an existing project
  - write code once, use a Makefile of CMake to build it for multiple backends
  - a *single application* can supports all the different backends *at the same time*
- the latest documentation is available at https://alpaka.readthedocs.io/en/latest/index.html



#### June 21st, 2023





### how does it work?



- Alpaka internally uses preprocessor symbols to enable the different backends:
  - ALPAKA\_ACC\_GPU\_CUDA\_ENABLED
  - ALPAKA\_ACC\_GPU\_HIP\_ENABLED
  - ALPAKA\_ACC\_CPU\_B\_SEQ\_T\_SEQ\_ENABLED

for running on NVIDIA GPUs for running on AMD GPUs for running serially on a CPU

- in this tutorial we will build separate applications from each example
  - each application is compiled with the corresponding compiler (g++, nvcc, hipcc, ...)
  - each application uses a single back-end
- it is also possible to enable more than one back-end at a time
  - however, the underlying CUDA and HIP header files will clash
  - separate the host and device parts, and use separate compilations for each backend

#### June 21st, 2023





### alsaka core concepts



### Host-side API

- initialisation and device selection: Platforms and Devices
- asynchronous operations and synchronisation: Queues and Events
- owning memory Buffers and non-owning memory Views
- submitting work to devices: work division and Accelerators

### Device-side API

- plain C++ for device functions and kernels
- shared memory, atomic operations, and memory fences
- primitives for mathematical operations
- warp-level primitives for synchronisation and data exchange (not covered)
- random number generator (not covered)

### nota bene:

• most Alpaka API objects behave like shared\_ptrs, and should be passed by value or by reference to const (*i.e.* const&)



#### June 21st, 2023



### platforms and devices



# alpaka: initialisation and device selection



### Platform and Device

- identify the type of hardware (*e.g.* host CPUs or NVIDIA GPUs) and individual devices (*e.g.* each single GPU) present on the machine
- the CPU device DevCpu serves two purposes:
  - as the "host" device, for managing the data flow (*e.g.* perform memory allocation and transfers, launch kernels, *etc.*)
  - as an "accelerator" device, for running heterogeneous code (*e.g.* to run an algorithm on the CPU)
- platforms cannot be instantiated, they are only used as a type
- devices should be created at the start of the program and used consistently
- some common cases

back end	alpaka platform	alpaka device
CPUs, serial or parallel	PltfCpu	DevCpu
NVIDIA GPU, with CUDA	PltfCudaRt	DevCudaRt
AMD GPUs, with HIP/ROCm	PltfHipRt	DevHipRt







### platforms and devices O



- Alpaka provides a simple API to enumerate the devices on a given platform:
  - alpaka::getDevCount<Platform>()
    - returns the number of devices on the given platform
  - alpaka::getDevs<Platform>()
    - initialises all the devices on the platform, and returns an std::vector<Device> with the corresponding Device objects
  - alpaka::getDevByIdx<Platform>(index)
    - initialises the index device on the platform, and returns the corresponding Device object
  - alpaka::getName(device)
    - returns the name of the given device

June 21st, 2023







#### .../intro\_to\_alpaka/alpaka/00\_enumerate.cc

```
int main() {
 // the host abstraction always has a single device
 Host host = alpaka::getDevByIdx<HostPlatform>(0u);
 std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';</pre>
 std::cout << "Found 1 device:\n";</pre>
 std::cout << " - " << alpaka::getName(host) << '\n';</pre>
  std::cout << std::endl:</pre>
  // get all the devices on the accelerator platform
 std::vector<Device> devices = alpaka::getDevs<Platform>();
  std::cout << "Accelerator platform: " << alpaka::core::demangled<Platform> << '\n';</pre>
  std::cout << "Found " << devices.size() << " device(s):\n";</pre>
 for (auto const& device : devices)
    std::cout << " - " << alpaka::getName(device) << '\n';</pre>
 std::cout << std::endl;</pre>
```

#### June 21st, 2023







22/78

.../intro\_to\_alpaka/alpaka/00\_enumerate.cc

```
int main() {
  // the host abstraction always has a single device
 Host host = alpaka::getDevByIdx<HostPlatform>(0u);
  std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';</pre>
  std::cout << "Found 1 device:\n";</pre>
  std::cout << " - " << alpaka::getName(host) << '\n';</pre>
  std::cout << std::endl:</pre>
                                                                                         these are the host and
  // get all the devices on the accelerator platform 
                                                                                         accelerator platforms
  std::vector<Device> devices = alpaka::getDevs<Platform>();
  std::cout << "Accelerator platform: " << alpaka::core::demangled<Platform> << '\n';</pre>
  std::cout << "Found " << devices.size() << " device(s):\n";</pre>
  for (auto const& device : devices)
    std::cout << " - " << alpaka::getName(device) << '\n';</pre>
  std::cout << std::endl;</pre>
```

#### June 21st, 2023





#### .../intro\_to\_alpaka/alpaka/00\_enumerate.cc

```
int main() {
 // the host abstraction always has a single device
 Host host = alpaka::getDevByIdx<HostPlatform>(0u);
 std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';</pre>
 std::cout << "Found 1 device:\n";</pre>
 std::cout << " - " << alpaka::getName(host) << '\n';</pre>
 std::cout << std::endl:</pre>
                                                                  • alpaka::core::demangled<T> is a string with
 // get all the devices on the accelerator platform
                                                                    the "human readable" name of a C++ type
 std::vector<Device> devices = alpaka::getDevs<Platform>();
 std::cout << "Accelerator platform: " << alpaka::core::demangled<Platform> << '\n';</pre>
 std::cout << "Found " << devices.size() << " device(s):\n";</pre>
 for (auto const& device : devices)
    std::cout << " - " << alpaka::getName(device) << '\n';</pre>
 std::cout << std::endl;</pre>
```







.../intro\_to\_alpaka/alpaka/00\_enumerate.cc

```
int main() {
  // the host abstraction always has a single device
  Host host = alpaka::getDevByIdx<HostPlatform>(Ou);
                                                                   • get the n<sup>th</sup> device for the given platform
  std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';</pre>
  std::cout << "Found 1 device:\n";</pre>
  std::cout << " - " << alpaka::getName(host) << '\n';</pre>
  std::cout << std::endl:</pre>
  // get all the devices on the accelerator platform
  std::vector<Device> devices = alpaka::getDevs<Platform>();

    get all devices for the given platform

  std::cout << "Accelerator platform: " << alpaka::core::demangled<Platform> << '\n';</pre>
  std::cout << "Found " << devices.size() << " device(s):\n";</pre>
  for (auto const& device : devices)
    std::cout << " - " << alpaka::getName(device) << '\n';</pre>
  std::cout << std::endl;</pre>
```

#### June 21st, 2023







25/78

.../intro\_to\_alpaka/alpaka/00\_enumerate.cc

```
int main() {
  // the host abstraction always has a single device
  Host host = alpaka::getDevByIdx<HostPlatform>(Ou);
  std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';</pre>
  std::cout << "Found 1 device:\n";</pre>
  std::cout << " - " << alpaka::getName(host) << '\n';</pre>
  std::cout << std::endl:</pre>
  // get all the devices on the accelerator platform

    get the name of each device

  std::vector<Device> devices = alpaka::getDevs<Platform>();
  std::cout << "Accelerator platform: " << alpaka::core::demangled<Platform> << '\n';</pre>
  std::cout << "Found " << devices.size() << " device(s):\n";</pre>
  for (auto const& device : devices)
    std::cout << " - " << alpaka::getName(device) << '\n';</pre>
  std::cout << std::endl;</pre>
```

#### June 21st, 2023



### some important details



\* g++ -std=c++17 -O2 -g -DALPAKA\_ACC\_CPU\_B\_SEQ\_T\_SEQ\_ENABLED -I\$BOOST\_BASE/include -I\$ALPAKA\_BASE/include 00\_enumerate.cc -o 00\_enumerate\_cpu

\* nvcc -x cu -std=c++17 -02 -g --expt-relaxed-constexpr -DALPAKA\_ACC\_GPU\_CUDA\_ENABLED -I\$B00ST\_BASE/include -I\$ALPAKA\_BASE/include 00\_enumerate.cc -o 00\_enumerate\_cuda

#include <iostream> #include <vector>

\*/

#include <alpaka/alpaka.hpp>

#include "config.h"

if you haven't done so, clone all the examples from GitHub

git clone https://github.com/fwyzard/intro\_to\_alpaka.git

June 21st, 2023





# let's build it ...



- using the CPU as the "accelerator"
  - the CPU acts as both the "host" and the "device"
  - the application runs entirely on the CPU

```
g++ -DALPAKA_ACC_CPU_B_SEQ_T_SEQ_ENABLED \
    -std=c++17 -02 -g -I$BOOST_BASE/include -I$ALPAKA_BASE/include \
    00_enumerate.cc \
    -o 00_enumerate_cpu
```

- using the CUDA GPUs as the "accelerator"
  - the CPU acts as the "host", the GPUs act as the "devices"
  - the application launches kernels that run on the GPUs

```
nvcc -x cu -expt-relaxed-constexpr -DALPAKA_ACC_GPU_CUDA_ENABLED \
    -std=c++17 -02 -g -I$BOOST_BASE/include -I$ALPAKA_BASE/include \
    00_enumerate.cc \
    -o 00_enumerate_cuda
```

#### June 21st, 2023



### ... and run it



\$ ./00\_enumerate\_cpu
Host platform: alpaka::PltfCpu
Found 1 device:

- Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz

Accelerator platform: alpaka::PltfCpu Found 1 device(s):

- Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz

#### 5 ./00\_enumerate\_cuda

Host platform: alpaka::PltfCpu

Found 1 device:

- Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz

Accelerator platform: alpaka::PltfUniformCudaHipRt<alpaka::ApiCudaRt> Found 4 device(s):

- Tesla V100-PCIE-32GB
- Tesla V100-PCIE-32GB
- Tesla V100-PCIE-32GB
- Tesla V100-PCIE-32GB

#### June 21st, 2023





### where is the magic?



#### .../intro\_to\_alpaka/alpaka/config.h

#if defined(ALPAKA\_ACC\_GPU\_CUDA\_ENABLED)

// CUDA backend

using Device = alpaka::DevCudaRt; using Platform = alpaka::Pltf<Device>;

#elif defined(ALPAKA\_ACC\_GPU\_HIP\_ENABLED)
// HIP/ROCm backend
using Device = alpaka::DevHipRt;
using Platform = alpaka::Pltf<Device>;

#elif defined(ALPAKA\_ACC\_CPU\_B\_SEQ\_T\_SEQ\_ENABLED)
// CPU serial backend
using Device = alpaka::DevCpu;

using Platform = alpaka::Pltf<Device>;

# back end alpaka platform alpaka device

CPUs, serial or parallelPltfCpuDevCpuNVIDIA GPU, with CUDAPltfCudaRtDevCudaRtAMD GPUs, with HIP/ROCmPltfHipRtDevHipRt

#### #else

// no backend specified

#error Please define one of ALPAKA\_ACC\_GPU\_CUDA\_ENABLED, ALPAKA\_ACC\_GPU\_HIP\_ENABLED, ALPAKA\_ACC\_CPU\_B\_SEQ\_T\_SEQ\_ENABLED

#### #endif

June 21st, 2023





### where is the magic?



#### .../intro\_to\_alpaka/alpaka/config.h

#if defined(ALPAKA ACC GPU CUDA ENABLED) // CUDA backend using Device = alpaka::DevCudaRt; using Platform = alpaka::Pltf<Device>; #elif defined(ALPAKA\_ACC\_GPU\_HIP\_ENABLED) // HIP/ROCm backend depending on which back-end is enabled ... using Device = alpaka::DevHipRt; using Platform = alpaka::Pltf<Device>; #elif defined ALPAKA ACC\_CPU B\_SEQ\_T\_SEQ\_ENABLED // CPU serial backend using Device = alpaka::DevCpu; using Platform = alpaka::Pltf<Device>; #else // no backend specified #error Please define one of ALPAKA ACC GPU CUDA ENABLED, ALPAKA ACC GPU HIP ENABLED, ALPAKA ACC CPU B SEQ T SEQ ENABLED

#endif

June 21st, 2023





### where is the magic?



#### .../intro\_to\_alpaka/alpaka/config.h #if defined(ALPAKA ACC GPU CUDA ENABLED) // CUDA backend using Device = alpaka::DevCudaRt; using Platform = alpaka::Pltf<Device>; #elif defined(ALPAKA ACC GPU HIP ENABLED) // HIP/ROCm backend depending on which back-end is enabled, using Device = alpaka::DevHipRt; Device and Platform are aliased to different types using Platform = alpaka::Pltf<Device>; #elif defined(ALPAKA ACC CPU B\_SEQ\_T\_SEQ\_ENABLED) // CPU serial backend using Device = alpaka::DevCpu; using Platform = alpaka::Pltf<Device>; #else

// no backend specified

#error Please define one of ALPAKA\_ACC\_GPU\_CUDA\_ENABLED, ALPAKA\_ACC\_GPU\_HIP\_ENABLED, ALPAKA\_ACC\_CPU\_B\_SEQ\_T\_SEQ\_ENABLED

#endif

#### June 21st, 2023



queues and events



### alpaka: asynchronous operations



#### Queues:

- identify a "work queue" where tasks (memory operations, kernel executions, ...) are executed in order
  - for example, a queue could represent an underlying CUDA stream or a CPU thread
  - from the point of view of the host , queues can be synchronous or asynchronous
- with a synchronous (or *blocking*) queue:
  - any operation is executed immediately, before returning to the caller
  - the host automatically waits (blocks) until each operation is complete
- with an asynchronous (or *non-blocking*) queue:
  - any operation is executed in the background, and each call returns immediately, without waiting for its completion
  - the host needs to synchronize explicitly with the queue, before accessing the results of the operations
- in general, prefer using a synchronous queue on a CPU, and an asynchronous queue on a GPU
- queues are always associated to a specific device
- most Alpaka operations (memory ops, kernel launches, etc.) are associated to a queue
- Alpaka does not provide a "default queue", create one explicitly





### common operations on queues



- creating a queue of the predefined type associated to a device is as simple as auto queue = Queue(device);
- waiting for all the asynchronous operations in a queue to complete is as simple as alpaka::wait(queue);
  - enque<mark>ue a host func</mark>tion
    - alpaka::enqueue(queue, task);
- enqueue a device function (launch a kernel)
  alpaka::exec<Accelerator>(queue, grid, kernel, arguments ...);
- allocate, set, or copy memory host and device memory alpaka::memcpy(queue, destination, source);







### alpaka: events and synchronisation



#### Events:

- events identify points in time along a work queue
- can be used to query or wait for the readiness of a task submitted to a queue
- can be used to synchronise different queues
- like queues, events are always associated to a specific device







٠

### common operations on events



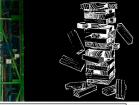
- events associated to a given device can be created with:
  - auto event = Event(device);
- events are enqueued to mark a given point along the queue:
  - alpaka::enqueue(queue, event);
    - an event is "complete" once all the work submitted to the queue before the event has been completed
- an event can be used to block the execution on the host until it is complete: alpaka::wait(event);
  - blocks the execution on the host
- or to make an other queue wait until a given event (in a different queue) is complete: alpaka::wait(other\_queue, event);
  - does not block execution on the host
  - further work submitted to other\_queue will only start after event is complete
- an event's status can also be queried without blocking the execution: alpaka::isComplete(event);

June 21st, 2023





### more magic



### .../intro\_to\_alpaka/alpaka/config.h

#### #if defined(ALPAKA\_ACC\_GPU\_CUDA\_ENABLED)

// CUDA backend

using Queue = alpaka::Queue<Device, alpaka::NonBlocking>;

```
using Event = alpaka::Event<Queue>;
```

#### #elif defined(ALPAKA\_ACC\_GPU\_HIP\_ENABLED)

// HIP/ROCm backend

```
using Queue = alpaka::Queue<Device, alpaka::NonBlocking>;
using Event = alpaka::Event<Queue>;
```

```
#elif defined(ALPAKA_ACC_CPU_B_SEQ_T_SEQ_ENABLED)
```

// CPU serial backend
using Queue = alpaka::Queue<Device, alpaka::Blocking>;
using Event = alpaka::Event<Queue>;

#### #else

// no backend specified
#error Please define one of ALPAKA\_ACC\_GPU\_CUDA\_ENABLED, ALPAKA\_ACC\_GPU\_HIP\_ENABLED, ALPAKA\_ACC\_CPU\_B\_SEQ\_T\_SEQ\_ENABLED

#endif

#### June 21st, 2023

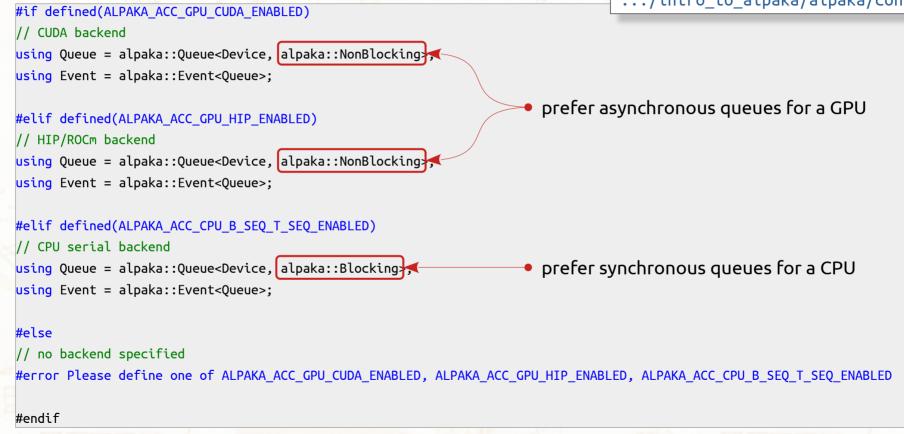




### more magic



### .../intro\_to\_alpaka/alpaka/config.h



#### June 21st, 2023





int main() {

# fun with queues



### .../intro\_to\_alpaka/alpaka/01\_blocking\_queue.cc

```
// the host abstraction always has a single device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
```

```
std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';
std::cout << "Found 1 device:\n";
std::cout << " - " << alpaka::getName(host) << '\n';
std::cout << std::endl;</pre>
```

```
// create a blocking host queue and submit some work to it
alpaka::Queue<Host, alpaka::Blocking> queue{host};
```

```
std::cout << "Enqueue some work\n";
alpaka::enqueue(queue, []() noexcept {
    std::cout << " - host task running...\n";
    std::this_thread::sleep_for(std::chrono::seconds(5u));
    std::cout << " - host task complete\n";
});
```

```
// wait for the work to complete
std::cout << "Wait for the enqueue work to complete...\n";
alpaka::wait(queue);
std::cout << "All work has completed\n";</pre>
```

### June 21st, 2023





# fun with queues





#### June 21st, 2023





int main() {

# fun with queues



### .../intro\_to\_alpaka/alpaka/01\_blocking\_queue.cc

```
// the host abstraction always has a single device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';</pre>
std::cout << "Found 1 device:\n";</pre>
std::cout << " - " << alpaka::getName(host) << '\n';</pre>
std::cout << std::endl:</pre>
// create a blocking host queue and submit some work to it
alpaka::Queue Host, alpaka::Blocking> queue{host};

    create a blocking queue on the Host

std::cout << "Enqueue some work\n";</pre>
alpaka::enqueue(queue, []() noexcept {
    std::cout << " - host task running...\n";</pre>
    std::this_thread::sleep_for(std::chrono::seconds(5u));
    std::cout << " - host task complete\n";</pre>
});
// wait for the work to complete
std::cout << "Wait for the enqueue work to complete...\n";</pre>
alpaka::wait(queue);
std::cout << "All work has completed\n";</pre>
```

#### June 21st, 2023





int main() {

# fun with queues



### .../intro\_to\_alpaka/alpaka/01\_blocking\_queue.cc

```
// the host abstraction always has a single device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
```

```
std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';
std::cout << "Found 1 device:\n";
std::cout << " - " << alpaka::getName(host) << '\n';
std::cout << std::endl;</pre>
```

```
// create a blocking host queue and submit some work to it
alpaka::Queue<Host, alpaka::Blocking> queue{host};
```

```
std::cout << "Enqueue some work\n";
alpaka::enqueue(queue, []() noexcept
std::cout << " - host task running...\n";
std::this_thread::sleep_for(std::chrono::seconds(5u));
std::cout << " - host task complete\n";
});
// wait for the work to complete</pre>
```

```
std::cout << "Wait for the enqueue work to complete...\n";
alpaka::wait(queue);
std::cout << "All work has completed\n";</pre>
```

#### June 21st, 2023





### fun with queues



43 / 78

### .../intro\_to\_alpaka/alpaka/01\_blocking\_queue.cc

### // the host abstraction always has a single device Host host = alpaka::getDevByIdx<HostPlatform>(0u);

```
std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';
std::cout << "Found 1 device:\n";
std::cout << " - " << alpaka::getName(host) << '\n';
std::cout << std::endl:</pre>
```

// create a blocking host queue and submit some work to it
alpaka::Queue<Host, alpaka::Blocking> queue{host};

```
std::cout << "Enqueue some work\n";</pre>
```

alpaka::enqueue(queue, []() noexcept {

```
std::cout << " - host task running...\n";
std::this_thread::sleep_for(std::chrono::seconds(5u));
std::cout << " - host task complete\n";</pre>
```

```
});
```

int main() {



this syntax introduces a lambda expression

that performs these operations

togethwer with alpaka::enqueue(...), this part
 - creates an object that encapsulates some operations
 - submits those opertations to run in a queue

#### June 21st, 2023



int main() {

# fun with queues



### .../intro\_to\_alpaka/alpaka/01\_blocking\_queue.cc

```
// the host abstraction always has a single device
Host host = alpaka::getDevByIdx<HostPlatform>(Ou);
```

```
std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';
std::cout << "Found 1 device:\n";
std::cout << " - " << alpaka::getName(host) << '\n';
std::cout << std::endl;</pre>
```

```
// create a blocking host queue and submit some work to it
alpaka::Queue<Host, alpaka::Blocking> queue{host};
```

```
std::cout << "Enqueue some work\n";
alpaka::enqueue(queue, []() noexcept {
    std::cout << " - host task running...\n";
    std::this_thread::sleep_for(std::chrono::seconds(5u));
    std::cout << " - host task complete\n";
});
```

```
// wait for the work to complete
std::cout << "Wait for the enqueue work to complete...\n";
alpaka::wait(queue);
std::cout << "All work has completed\n";</pre>
```

wait for the enqueued operations to complete

#### June 21st, 2023





# let's build it and run it O



- in this example we are not making use of any accelerator
  - let's build it only for the CPU back-end

```
g++ -DALPAKA_ACC_CPU_B_SEQ_T_SEQ_ENABLED \
    -std=c++17 -02 -g -I$BOOST_BASE/include -I$ALPAKA_BASE/include \
    01_blocking_queue.cc \
    -o 01_blocking_queue_cpu
```

### and run it

\$ ./01\_blocking\_queue\_cpu

Host platform: alpaka::PltfCpu

Found 1 device:

- Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz

#### Enqueue some work

- host task running...
- host task complete
- Wait for the enqueue work to complete...

```
All work has completed
```

### June 21st, 2023





int main() {

### an async example



### .../intro\_to\_alpaka/alpaka/02\_nonblocking\_queue.cc

### // the host abstraction always has a single device Host host = alpaka::getDevByIdx<HostPlatform>(0u);

```
std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';
std::cout << "Found 1 device:\n";
std::cout << " - " << alpaka::getName(host) << '\n';
std::cout << std::endl;</pre>
```

### // create a non-blocking host queue and submit some work to it alpaka::Queue<Host, alpaka::NonBlocking> queue{host};

```
std::cout << "Enqueue some work\n";
alpaka::enqueue(queue, []() noexcept {
    std::cout << " - host task running...\n";
    std::this_thread::sleep_for(std::chrono::seconds(5u));
    std::cout << " - host task complete\n";
});
```

```
// wait for the work to complete
std::cout << "Wait for the enqueue work to complete...\n";
alpaka::wait(queue);
std::cout << "All work has completed\n";</pre>
```

### June 21st, 2023





int main() {

### an async example



### .../intro\_to\_alpaka/alpaka/02\_nonblocking\_queue.cc

```
// the host abstraction always has a single device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
std::cout << "Host platform: " << alpaka::core::demangled<HostPlatform> << '\n';</pre>
std::cout << "Found 1 device:\n";</pre>
std::cout << " - " << alpaka::getName(host) << '\n';</pre>
std::cout << std::endl:</pre>
// create a non-blocking host queue and submit some work to it
alpaka::Queue Host, alpaka::NonBlocking> queue{host};

    create a non-blocking queue on the Host

std::cout << "Enqueue some work\n";</pre>
alpaka::enqueue(queue, []() noexcept {
    std::cout << " - host task running...\n";</pre>
    std::this_thread::sleep_for(std::chrono::seconds(5u));
    std::cout << " - host task complete\n";</pre>
});
// wait for the work to complete
std::cout << "Wait for the enqueue work to complete...\n";</pre>
alpaka::wait(queue);
std::cout << "All work has completed\n";</pre>
```

#### June 21st, 2023





# let's build it and run it O



- in this example, too, we are not making use of any accelerator
  - let's build it only for the CPU back-end with POSIX threads

```
g++ -DALPAKA_ACC_CPU_B_SEQ_T_SEQ_ENABLED \
    -std=c++17 -02 -g -I$BOOST_BASE/include -I$ALPAKA_BASE/include -pthread \
    02_nonblocking_queue.cc \
    -o 02_nonblocking_queue_cpu
```

### and run it

\$ ./02\_nonblocking\_queue\_cpu

Host platform: alpaka::PltfCpu

Found 1 device:

- Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz

#### Enqueue some work

Wait for the enqueue work to complete...

- host task running...
- host task complete

All work has completed

### June 21st, 2023





### blocking vs non-blocking



\$ ./01\_blocking\_queue\_cpu
Host platform: alpaka::PltfCpu
Found 1 device:

- Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz

Enqueue some work

- host task running...
- host task complete

Wait for the enqueue work to complete...

All work has completed

\$ ./02\_nonblocking\_queue\_cpu
Host platform: alpaka::PltfCpu
Found 1 device:

- Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz

Enqueue some work

Wait for the enqueue work to complete...

- host task running...
- host task complete
- All work has completed

June 21st, 2023





# blocking vs non-blocking



\$ ./01\_blocking\_queue\_cpu
Host platform: alpaka::PltfCpu
Found 1 device:

- Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz

Enqueue some work

- host task running...
- host task complete

Wait for the enqueue work to complete... All work has completed

- with a synchronous (or *blocking*) queue:
  - any operation is executed immediately, before returning to the caller
  - the host automatically waits (blocks) until each operation is complete
- with an asynchronous (or *non-blocking*) queue:
  - any operation is executed in the background, and each call returns immediately, without waiting for its completion
  - the host needs to synchronize explicitly with the queue, before accessing the results of the operations

Host platform: alpaka::PltfCpu
Found 1 device:
 - Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz
Enqueue some work
Wait for the enqueue work to complete...
 - host task running...

- host task complete

\$ ./02 nonblocking queue cpu

All work has completed

June 21st, 2023

### A. Bocci - An introduction to Alpaka



50

memory operations



### memory in alpaka



### Buffers and Views

- can refer to memory on the host or on any device
  - general purpose host memory (e.g. as returned by malloc or new)
  - pinned host memory, visible by devices on a given platform (e.g. as returned by cudaMallocHost)
  - global device memory (e.g. as returned by cudaMalloc)
- can have arbitrary dimensions
- 0-dimensional buffers and views wrap and provide access to a single element:

float x = \*buffer;
float y = buffer->pt();

• 1-dimensional buffers and views wrap and provide access to an array of elements:

```
float x = buffer[i];
```

N-dimensional buffers and views wrap arbitrary memory areas:

```
float* p = std::data(buffer);
```

expect a nicer accessor syntax with c++23 std::mdspan and improved operator[]

### June 21st, 2023





# memory buffers



- buffers own the memory they point to
  - a host memory buffer can use either standard host memory, or pinned host memory mapped to be visible by the GPUs in a given platform
  - a buffer knows what device the memory is on, and how to free it
- buffers have shared ownership of the memory
  - like shared\_ptr<T>
  - making a copy of a buffer creates a second handle to the same underlying memory
  - the memory is automatically freed when the last buffer object is destroyed (*e.g.* goes out of scope)
    - with queue-ordered semantic, memory is freed when the work submitted to the queue associate to the buffer is complete
- note that buffers always allow modifying their content
  - a Buffer<const T> would not be useful, because its contents could never be set
  - a const Buffer<T> does not prevent changes to the contents, as they can be modified through a copy





•

### allocating memory



- buffer allocations and deallocations can be immediate or queue-ordered
  - immediate operations
    - allocate and free the memory immediately
    - may result in a device-wide synchronisation
    - e.g. malloc / free or cudaMalloc / cudaFree

```
// allocate an array of "size" floats in standard host memory
auto buffer = alpaka::allocBuf<float, uint32_t>(host, size);
```

```
// allocate an array of "size" floats in pinned host memory
// mapped to be efficiently copieable to/from all the devices on the Platform
auto buffer = alpaka::allocMappedBuf<Platform, float, uint32_t>(host, size);
```

```
// alloca an array of "size" floats in global device memory
auto buffer = alpaka::allocBuf<float, uint32_t>(device, size);
```

- queue-ordered operations are usually asynchronous, and may cache allocations
  - guarantee that the memory is allocated before any further operations submitted to the queue are executed
  - guarantee that the memory will be freed once all pending operation in the queue are complete
  - e.g. cudaMallocAsync / cudaFreeAsync

// allocate an array of "size" floats in global gpu memory, ordered along queue
auto buffer = alpaka::allocAsyncBuf<float, uint32\_t>(queue, size);

available only on device that support it (CPUs, NVIDIA CUDA  $\geq$  11.2, AMD ROCm  $\geq$  5.4)

### June 21st, 2023





### using buffers



### .../intro\_to\_alpaka/alpaka/03\_memory.cc

```
// require at least one device
std::size t n = alpaka::getDevCount<Platform>():
if (n == 0) {
 exit(EXIT FAILURE);
```

```
// use the single host device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
std::cout << "Host: " << alpaka::getName(host) << '\n';</pre>
```

```
// allocate a buffer of floats in host memory, mapped to ... the device
uint32 t size = 42;
auto host buffer =
```

alpaka::allocMappedBuf<Platform, float, uint32 t>(host, Vec1D{size}); std::cout << "pinned host memory buffer at " << std::data(host buffer) << "\n\n";</pre>

```
// fill the host buffers with values
for (uint32 t i = 0; i < size; ++i) {</pre>
 host buffer[i] = i;
```

```
// use the first device
Device device = alpaka::getDevByIdx<Platform>(0u);
std::cout << "Device: " << alpaka::getName(device) << '\n';</pre>
```

// create a work queue Queue queue{device};

#### // ...

// allocate a buffer of floats in global device memory, asynchronously auto device buffer = alpaka::allocAsyncBuf<float, uint32 t>(queue, Vec1D{size}); std::cout << "memory buffer on " << alpaka::getName(alpaka::getDev(device buffer))</pre> << " at " << std::data(device buffer) << "\n\n";

// set the device memory to all zeros (byte-wise, not element-wise) alpaka::memset(queue, device buffer, 0x00);

// copy the contents of the device buffer to the host buffer alpaka::memcpy(queue, host buffer, device buffer);

// the device buffer goes out of scope, but the memory is freed only // once all enqueued operations have completed

// wait for all operations to complete alpaka::wait(queue);

// read the content of the host buffer for (uint32 t i = 0; i < size; ++i) {</pre> std::cout << host\_buffer[i] << ' ';</pre> std::cout << '\n';</pre>

#### June 21st, 2023

#### A. Bocci - An introduction to Alpaka

}

}







#### .../intro\_to\_alpaka/alpaka/03\_memory.cc // require at least one device std::size t n = alpaka::getDevCount<Platform>(): // ... **if** (n == 0) { exit(EXIT FAILURE); // allocate a buffer of floats in global device memory, asynchronously auto device buffer = alpaka::allocAsyncBuf<float, uint32 t>(queue, Vec1D{size}); // use the single host device std::cout << "memory buffer on " << alpaka::getNam (alpaka::getDev(device buffer))</pre> Host host = alpaka::getDevByIdx<HostPlatform>(0u); << " at " << std::data(device buffer) << "\n\n"; std::cout << "Host: " << alpaka::getName(host) << '\n';</pre> // set the device memory to all zeros (byte-wise, not element-wise) // allocate a buffer of floats in host memory, mapped to ... the device alpaka::memset(queue, device buffer, 0x00); uint32 t size = 42; // copy the contents of the device buffer to the host buffer auto host buffer = alpaka::allocMappedBuf<Platform, float, uint32\_t>(host, Vec1D{size}); alpaka::memcpy(queue, host buffer, device buffer); std::cout << "pinned host memory buffe at " << std::data(host buffer) << "\n\n";</pre> // the device buffer goes out of scope, but the memory is freed only // fill the host buffers with values // once all enqueued operations have completed for (uint32 t i = 0; i < size; ++i) {</pre> } host buffer[i] = i; // wait for all operations to complete allocate buffers alpaka::wait(queue); // use the first device Device device = alpaka::getDevByIdx<Platform>(0u); // read the content of the host buffer std::cout << "Device: " << alpaka::getName(device) << '\n';</pre> for (uint32 t i = 0; i < size; ++i) {</pre> std::cout << host\_buffer[i] << ' ';</pre> // create a work queue std::cout << '\n';</pre> Queue queue{device};

#### June 21st, 2023







#### .../intro\_to\_alpaka/alpaka/03\_memory.cc // require at least one device std::size t n = alpaka::getDevCount<Platform>(): // ... **if** (n == 0) { exit(EXIT FAILURE); // allocate a buffer of floats in global device memory, asynchronously auto device buffer = alpaka::allocAsyncBuf<float, uint32 t>(queue, Vec1D{size}); // use the single host device std::cout << "memory buffer on " << alpaka::getName(alpaka::getDev(device buffer))</pre> Host host = alpaka::getDevByIdx<HostPlatform>(0u); << " at " << std::data(device buffer) << "\n\n"; std::cout << "Host: " << alpaka::getName(host) << '\n';</pre> // set the device memory to all zeros (byte-wise, not element-wise) // allocate a buffer of floats in host memory, mapped to ... the device alpaka::memset(queue, device\_buffer, 0x00); uint32 t size = 42; auto host buffer = // copy the contents of the device buffer to the host buffer alpaka::allocMappedBuf<Platform, float, uint32 t>(host, Vec1D{size}); alpaka::memcpy(queue, host\_buffer, device\_buffer); std::cout << "pinned host memory buffer at " << std::data(host buffer) << "\n\n";</pre> // the device buffer goes out of scope, but the memory is freed only // fill the host buffers with values *H* once all enqueued operations have completed for (uint32 t i = 0; i < size; ++i) {</pre> } host\_buffer[i] = i; // wait for all operations to complete get the buffers' memory addresses alpaka::wait(queue); // use the first device Device device = alpaka::getDevByIdx<Platform>(0u); // read the content of the host buffer std::cout << "Device: " << alpaka::getName(device) << '\n';</pre> for (uint32 t i = 0; i < size; ++i) {</pre> std::cout << host\_buffer[i] << ' ';</pre> // create a work queue std::cout << '\n';</pre> Queue queue{device};

#### June 21st, 2023





# using buffers



### .../intro\_to\_alpaka/alpaka/03\_memory.cc

```
// require at least one device
std::size_t n = alpaka::getDevCount<Platform>();
if (n == 0) {
    exit(EXIT_FAILURE);
}
```

```
// use the single host device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
std::cout << "Host: " << alpaka::getName(host) << '\n';</pre>
```

```
// allocate a buffer of floats in host memory, mapped to ... the device
uint32_t size = 42;
auto host buffer =
```

alpaka::allocMappedBuf<Platform, float, uint32\_t>(host, Vec1D{size}); std::cout << "pinned host memory buffer at " << std::data(host\_buffer) << "\n\n";</pre>

```
// fill the host buffers with values
for (uint32 t i = 0; i < size; ++i) {
    host_buffer[i] = i;
    write to and read from
    the host buffer like a vector
// use the first device
Device device = alpaka::getDevByIdx<Platform>(0u);
std::cout << "Device: " << alpaka::getName(device) << '\n';</pre>
```

// create a work queue
Queue queue{device};

#### // ...

// set the device memory to all zeros (byte-wise, not element-wise)
alpaka::memset(queue, device\_buffer, 0x00);

// copy the contents of the device buffer to the host buffer alpaka::memcpy(queue, host\_buffer, device\_buffer);

// the device buffer goes out of scope, but the memory is freed only
// once all enqueued operations have completed

// wait for all operations to complete
alpaka::wait(queue);

```
// read the content of the host buffer
for (uint32_t i = 0; i < size; ++i) {
   std::cout host_buffer[i] << ' ';
}
std::cout << '\n';</pre>
```

June 21st, 2023







### .../intro\_to\_alpaka/alpaka/03\_memory.cc

// require at least one device
std::size\_t n = alpaka::getDevCount<Platform>();
if (n == 0) {
 exit(EXIT\_FAILURE);
}

// use the single host device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
std::cout << "Host: " << alpaka::getName(host) << '\n';</pre>

// allocate a buffer of floats in host memory, mapped to ... the device
uint32\_t size = 42;
auto host buffer =

alpaka::allocMappedBuf<Platform, float, uint32\_t>(host, Vee1D{size}); std::cout << "pinned host memory buffer at " << std::data(host\_buffer) << "\n\n";</pre>

// fill the host buffers with values
for (uint32\_t i = 0; i < size; ++i) {
 host\_buffer[i] = i;
}
memse</pre>

memset and memcpy operations are always asynchronous

// use the first device

Device device = alpaka::getDevByIdx<Platform>(0u); std::cout << "Device: " << alpaka::getName(device) << '\n';</pre>

// create a work queue
Queue queue{device};

#### // ...

// set the device memory to all zeros (byte-wise, not element-wise)
alpaka::memset(queue, device\_buffer, 0x00);

// copy the contents of the device buffer to the host buffer
alpaka::memcpy(queue, host\_buffer, device\_buffer);

// the device buffer goes out of scope, but the memory is freed only // once all enqueued operations have completed

// wait for all operations to complete
alpaka::wait(queue);

```
// read the content of the host buffer
for (uint32_t i = 0; i < size; ++i) {
   std::cout << host_buffer[i] << ' ';
}
std::cout << '\n';</pre>
```

#### June 21st, 2023





### memory views



- views wrap memory allocated by some other mechanism to provide a common interface
  - e.g. a local variable on the stack, or memory owned by an std::vector
  - views *do not own* the underlying memory
  - the lifetime of a view should not exceed that of the memory it points to

```
float* data = new float[size];
auto view = alpaka::ViewPlainPtr<float, uint32_t>(data, host, Vec1D{size}); // define a view for a C++ array
alpaka::memcpy(queue, view, device_buffer); // copy the data to the array
```

- views to standard containers
  - Alpaka provides adaptors and can automatically use std::array<T, N> and std::vector<T> as views

std::vector<float> data(size);
alpaka::memcpy(queue, data, device\_buffer);

// copy the data to the vector

- using views to emulate buffers to constant objects
  - buffers always allow modifying their content
  - but we can wrap them in a constant view: alpaka::ViewConst<Buffer<T>>

auto const\_view = alpaka::ViewConst(device\_buffer); alpaka::memcpy(queue, host\_buffer, const\_view);

// copy the data to the host

### June 21st, 2023







```
// require at least one device
std::size_t n = alpaka::getDevCount<Platform>();
if (n == 0) {
    exit(EXIT_FAILURE);
}
```

```
// use the single host device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
std::cout << "Host: " << alpaka::getName(host) << '\n';</pre>
```

```
// allocate a buffer of floats in host memory, mapped to ... the device
uint32_t size = 42;
std::vector<float> host_data(size);
std::cout << "host vector at " << std::data(host_data) << "\n\n";</pre>
```

```
// fill the host buffers with values
for (uint32_t i = 0; i < size; ++i) {
    host_data[i] = i;
}</pre>
```

```
// use the first device
Device device = alpaka::getDevByIdx<Platform>(0u);
std::cout << "Device: " << alpaka::getName(device) << '\n';</pre>
```

// create a work queue
Queue queue{device};

### .../intro\_to\_alpaka/alpaka/04\_views.cc

// set the device memory to all zeros (byte-wise, not element-wise)
alpaka::memset(queue, device\_buffer, 0x00);

// create a read-only view to the device data
auto const\_view = alpaka::ViewConst(device\_buffer);

// copy the contents of the device buffer to the host buffer alpaka::memcpy(queue, host\_data, const\_view);

// the device buffer goes out of scope, but the memory is freed only // once all enqueued operations have completed

// wait for all operations to complete
alpaka::wait(queue);

// read the content of the host buffer
for (uint32\_t i = 0; i < size; ++i) {
 std::cout << host\_data[i] << ' ';
}
std::cout << '\n';</pre>

#### June 21st, 2023





### using views



# // require at least one device std::size\_t n = alpaka::getDevCount<Platform>(); if (n == 0) { exit(EXIT\_FAILURE); }

```
// use the single host device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
std::cout << "Host: " << alpaka::getName(host) << '\n';</pre>
```

```
// allocate a buffer of floats in host memory, mapped to ... the device
uint32 t size = 42;
std::vector<float> host_data(size);
std::cout << "host vector at " << std::data(host data) << "host";</pre>
```

```
// fill the host buffers with values
for (uint32_t i = 0; i < size; ++i) {
    host_data[i] = i;
}</pre>
```

use a vector directly

```
// use the first device
Device device = alpaka::getDevByIdx<Platform>(0u);
std::cout << "Device: " << alpaka::getName(device) << '\n';</pre>
```

```
// create a work queue
Queue queue{device};
```

### .../intro\_to\_alpaka/alpaka/04\_views.cc

// set the device memory to all zeros (byte-wise, not element-wise)
alpaka::memset(queue, device\_buffer, 0x00);

// create a read-only view to the device data
auto const\_view = alpaka::ViewConst(device\_buffer);

```
// copy the contents of the device buffer to the host buffer
_alpaka::memcpy(que); host_data, const_view);
```

```
// the device buffer goes out of scope, but the memory is freed only // once all enqueued operations have completed
```

```
// wait for all operations to complete
alpaka::wait(queue);
```

```
// read the content of the host buffer
for (uint32_t i = 0; i < size; ++i) {
   std::cout << host_data[i] << ' ';
}
std::cout << '\n';</pre>
```

#### June 21st, 2023





### using views



```
// require at least one device
std::size_t n = alpaka::getDevCount<Platform>();
if (n == 0) {
    exit(EXIT_FAILURE);
}
```

```
// use the single host device
Host host = alpaka::getDevByIdx<HostPlatform>(0u);
std::cout << "Host: " << alpaka::getName(host) << '\n';</pre>
```

```
// allocate a buffer of floats in host memory, mapped to ... the device
uint32_t size = 42;
std::vector<float> host_data(size);
std::cout << "host vector at " << std::data(host_data) << "\n\n";</pre>
```

```
// fill the host buffers with values
for (uint32_t i = 0; i < size; ++i) {
    host_data[i] = i;
}</pre>
```

copy from a const view to garantee not changing the device buffer

// use the first device
Device device = alpaka::getDevByIdx<Platform>(0u);
std::cout << "Device: " << alpaka::getName(device) << '\n';</pre>

```
// create a work queue
Queue queue{device};
```

### .../intro\_to\_alpaka/alpaka/04\_views.cc

63 / 78

// set the device memory to all zeros (byte-wise, not element-wise)
alpaka::memset(queue, device\_buffer, 0x00);

// create a read-only view to the device data
auto const\_view = alpaka::ViewConst(device\_buffer);

```
// copy the contents of the device buffer to the host buffer
_alpaka::memcpy(queue, host_data, _____nst_view);
```

// the device buffer goes out of scope, but the memory is freed only // once all enqueued operations have completed

```
// wait for all operations to complete
alpaka::wait(queue);
```

```
// read the content of the host buffer
for (uint32_t i = 0; i < size; ++i) {
   std::cout << host_data[i] << ' ';
}
std::cout << '\n';</pre>
```

#### June 21st, 2023





### alpaka device functions



### device functions

device functions are marked with the ALPAKA\_FN\_ACC macro

ALPAKA\_FN\_ACC
float my\_func(float arg) { ... }

- backend-specific functions
  - if the implementation of a device function may depend on the backend or on the work division into groups and threads, it should be templated on the Accelerator type, and take an Accelerator object

```
template <typename TAcc>
ALPAKA_FN_ACC
float my_func(TAcc const& acc, float arg) { ... }
```

- the availability of C++ features depends on the backend and on the device compiler
  - dynamic memory allocation is (partially) supported, but strongly discouraged
  - c++ std containers should be avoid
  - exceptions are usually not supported
  - recursive functions are supported only by some backends (CUDA: yes, but often inefficient; SYCL: no)
  - c++20 is available in CUDA code only starting from CUDA 12.0
  - etc.

### June 21st, 2023





### alsaka device functions



examples:

- mathematical operations are similar to what is available in the c++ standard:
  - e.g.
    - alpaka::math::sin(acc, arg)
- atomic operations are similar to what is available in CUDA and HIP
  - e.g.

alpaka::atomicAdd(acc, T\* address, T value, alpaka::hierarchy::Blocks)

- warp-level functions are similar to what is available in CUDA and HIP
  - e.g.

alpaka::warp::ballot(acc, arg)

June 21st, 2023









### kernels

- are implemented as an ALPAKA\_FN\_ACC void operator()(...) const function of a dedicated struct or class
  - kernels never return anything: -> void
  - kernels cannot change any data member on the host: must be declared const
- are always templated on the accelerator type, and take an accelerator object as the first argument

```
struct Kernel {
  template <typename TAcc>
  ALPAKA_FN_ACC void operator()(
    TAcc const& acc,
    float const* in1, float const* in2, float* out, size_t size) const
  {
    ...
  }
};
```

the TAcc acc argument identifies the backend and provides the details of the work division





# alpaka: grids, blocks, threads...



- alpaka maintains the work division into blocks and threads used in CUDA and OpenCL:
  - a kernel launch is divided into a grid of **blocks** 
    - the various block are scheduled independently, so they may be running concurrently or at different times
    - operations in different blocks cannot be synchronised
    - operations in different blocks can communicate only through the device global memory
  - each block is composed of threads running in parallel
    - threads in a block tend to run concurrently, but may diverge or be scheduled independently from each other
    - operations in a block can be synchronised, *e.g.* with alpaka::syncBlockThreads(acc);
    - operations in a block can communicate through shared memory
  - blocks can be decomposed into sub-groups, *i.e.* warps
    - threads in the same warp can synchronise and exchange data using more efficient primitives





### ... and elements ?



- to support efficient algorithms running on a CPU, alpaka introduces an additional level in the execution hierarchy: elements
  - each thread in a block may run on multiple consecutive elements
  - CPU backends usually run with multiple elements per thread
    - a good choice might be 16 elements, so 16 consecutive integers or floats can be loaded into a cache line
    - in principle, this could allow a host compiler to auto-vectorise the code, but more testing and development is needed !
  - GPU backends usually run with a single element per thread
    - memory accesses are already coalesced at the warp level
    - in principle, 2 elements per thread could be used with short or float16 data
- kernel should be written to allow for different number of elements per thread
  - a common approach is to use
    - N blocks, M threads per block, 1 element per thread on a GPU
    - N blocks, 1 thread per block, M elements per thread on a CPU





### a simple strided loop



- we provide a helper to implement a simple N-dimensional strided loop
  - the launch grid is tiled and repeated as many times as needed to cover the problem size
  - this tends to be the most efficient approach when all threads can work independently

```
struct Kernel {
  template <typename TAcc>
  ALPAKA_FN_ACC void operator()(
   TAcc const& acc,
   float const* in1, float const* in2, float* out, size_t size) const
  {
   for (auto index : elements_with_stride(acc, size)) {
     out[index] = in1[index] + in2[index];
   }
  }
};
```

for more complicated cases, use the alpaka::getWorkDiv and alpaka::getIdx functions

#include "workdivision.h"



### launching kernels



### alpaka: work submission



### Accelerator

- describes "how" a kernel runs on a device
  - N-dimensional work division (1D, 2D, 3D, ...)
  - on the CPU, serial vs parallel execution at the thread and block level (single thread, multi-threads, TBB tasks, ...)
  - implementation of shared memory, atomic operations, *etc*.
- accelerators are created only when a kernel is executed, and can only be accessed in device code
  - each device function can (should) be templated on the accelerator type, and take an accelerator as its first argument
  - the accelerator object can be used to extract the execution configuration (blocks, threads, elements)
  - the accelerator type can be used to implement per-accelerator behaviour
- for example, an algorithm can be implemented in device code using a parallel approach for a GPU-based accelerator, and a serial approach for a CPU-based accelerator





### launching a kernel



- a kernel launch requires
  - the type of the accelerator where the kernel will run
  - the queue to submit the work to
  - the work division into blocks, threads, and elements
  - an instance of the type that implements the kernel
  - the arguments to the kernel function
- we provide some helper types and functions
  - config.h includes the aliases Acc1D, Acc2D, Acc3D for 1D, 2D and 3D kernels
  - workdivision.h provides the helper function make\_workdiv<TAcc>(blocks, threads\_or\_elements)

// launch a 1-dimensional kernel with 32 groups of 32 threads (GPU) or elements (CPU)
auto grid = make\_workdiv<Acc1D>(32, 32);
alpaka::exec<Acc1D>(queue, grid, Kernel{}, a.data(), b.data(), sum.data(), size);





### a complete al saka example



### a complete al saka example



### .../intro\_to\_alpaka/alpaka/05\_kernel.cc

• running on the CPU

\$ ./05\_kernel\_cpu Host: Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz Device: Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz Testing VectorAddKernel with scalar indices with a grid of (32) blocks x (1) threads x (32) elements... success Testing VectorAddKernel1D with vector indices with a grid of (32) blocks x (1) threads x (32) elements... success Testing VectorAddKernel3D with vector indices with a grid of (5, 5, 1) blocks x (1, 1, 1) threads x (4, 4, 4) elements... success

running on the GPU

\$ ./05\_kernel\_cuda Host: Intel(R) Xeon(R) Gold 6252 CPU @ 2.10GHz Device: Tesla V100-PCIE-32GB Testing VectorAddKernel with scalar indices with a grid of (32) blocks x (32) threads x (1) elements... success Testing VectorAddKernel1D with vector indices with a grid of (32) blocks x (32) threads x (1) elements... success Testing VectorAddKernel3D with vector indices with a grid of (5, 5, 1) blocks x (4, 4, 4) threads x (1, 1, 1) elements... success





### summary





June 21st, 2023

- we have covered
  - what performance portability means what is the Alpaka library
  - how to set up Alpaka for a simple project
  - how to compile a single source file for different back-ends
  - what are Alpaka platforms, devices, queues and events
  - how to work with host and device memory
  - how to write device functions and kernels
  - how to use an Alpaka accelerator and work division to launch a kernel
  - and a complete example !
- congratulations!
  - now you can write *portable* and *performant* applications





(more) questions?



Copyright CERN 2023

Creative Commons 4.0 Attribution-ShareAlike International - CC BY-SA 4.0