



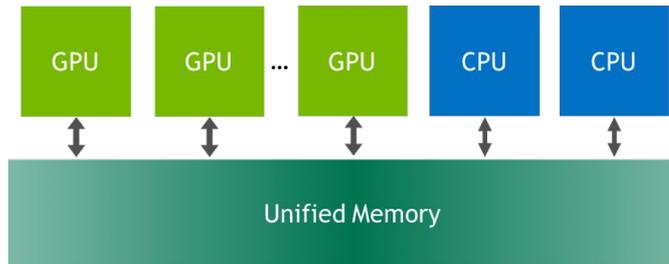
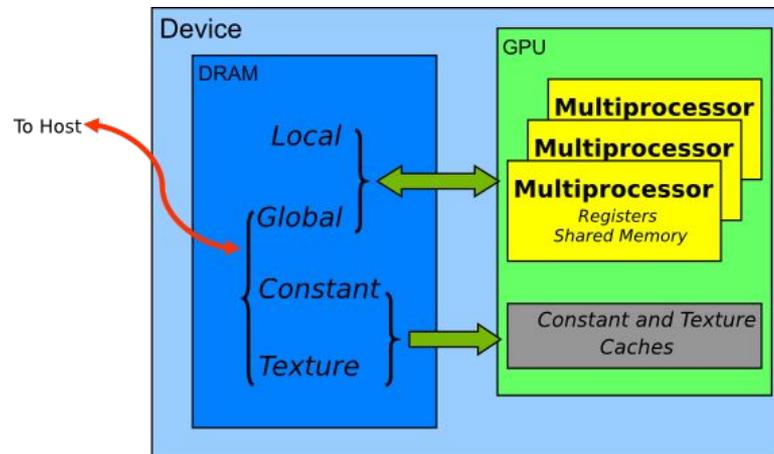
Heterogeneous Memory Management with VecMem

Attila Krasznahorkay, Stephen Swatman, Paul Gessinger

- Memory management on accelerators
- C++ memory resources
- Memory management with memory resources
 - Organising the memory handling in host code
 - Organising the memory handling in device code
- Uses of VecMem in algorithmic code

Accelerator Memory Management

- Modern CPUs have a very complicated memory management system
 - Which we can in most cases avoid knowing about
- GPUs have a complicated system of their own
 - However this we can not avoid knowing more about to use GPUs efficiently 😞
 - Most importantly, caching is much less automated than on modern CPUs
- In some cases however you can get away with not knowing everything
 - For a performance penalty...



Memory Management APIs



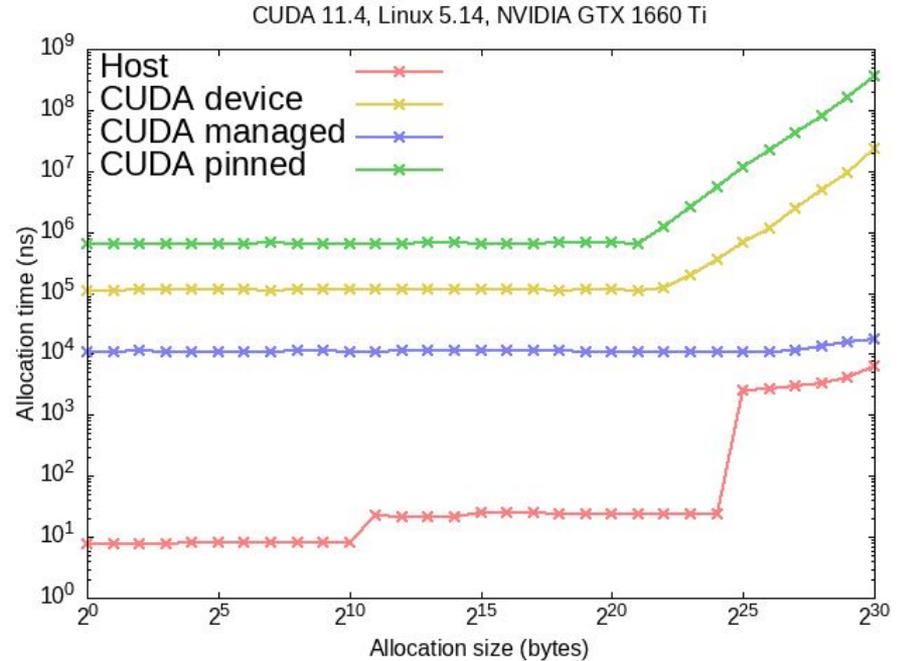
```
const std::size_t arraySize = ...;
float *hostArray = nullptr,
      *deviceArray = nullptr;
hostArray =
    (float*) malloc( arraySize * sizeof(float) );
CUDA_CHECK(
    cudaMalloc( &deviceArray,
               arraySize * sizeof(float) ) );
...
CUDA_CHECK(
    cudaMemcpy( deviceArray, hostArray,
               arraySize * sizeof(float),
               cudaMemcpyHostToDevice ) );
doSomething<<<...>>>( ..., deviceArray, ... );
...
free( hostArray );
CUDA_CHECK( cudaFree( deviceArray ) );
```

- All “heterogeneous platforms” provide memory management functions through a C API
 - Even in [oneAPI/SYCL](#) this is becoming the dominant way for managing memory
- This is very understandable, as it allows for interfacing with practically any other language
- However it is very different from how you are supposed to manage memory in modern C++
 - You shouldn't even be using [new/delete](#), let alone [std::malloc\(...\)/std::free\(...\)](#)

Device Memory Allocation Performance



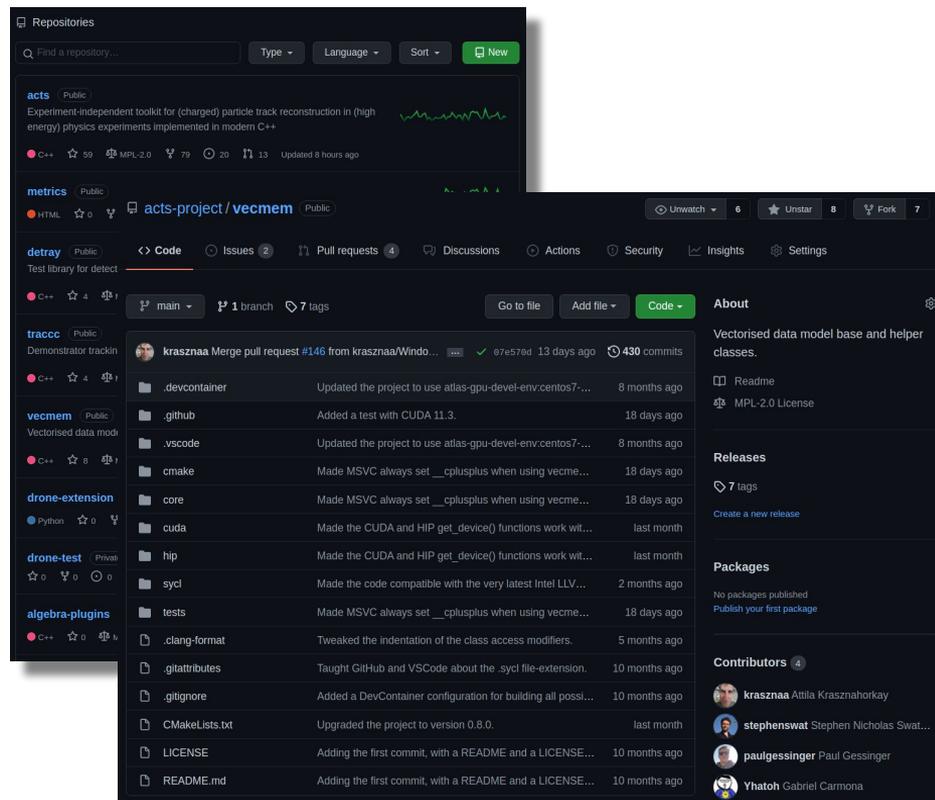
- Using the CUDA/HIP/SYCL “C API” is not only inconvenient, but it is also significantly slower than host memory allocation
- This is one of the reasons why we generally can't handle device memory with the same flexibility as we can host memory
 - Naive `push_back(...)`-s into a vector would be a performance killer



The VecMem Project



- In ATLAS we are currently working on a “full chain demonstrator” for running charged particle tracking on GPUs
 - The “top level” project being [acts-project/traccv](https://github.com/acts-project/traccv)
 - As the development started last year, it quickly became clear that we needed to break the issue down into separate domains
 - We ended up separating the code for [memory management \(VecMem\)](#), [algebra abstraction](#) and [detector geometry description](#) out of the “top level” repository



C++17 Memory Resources



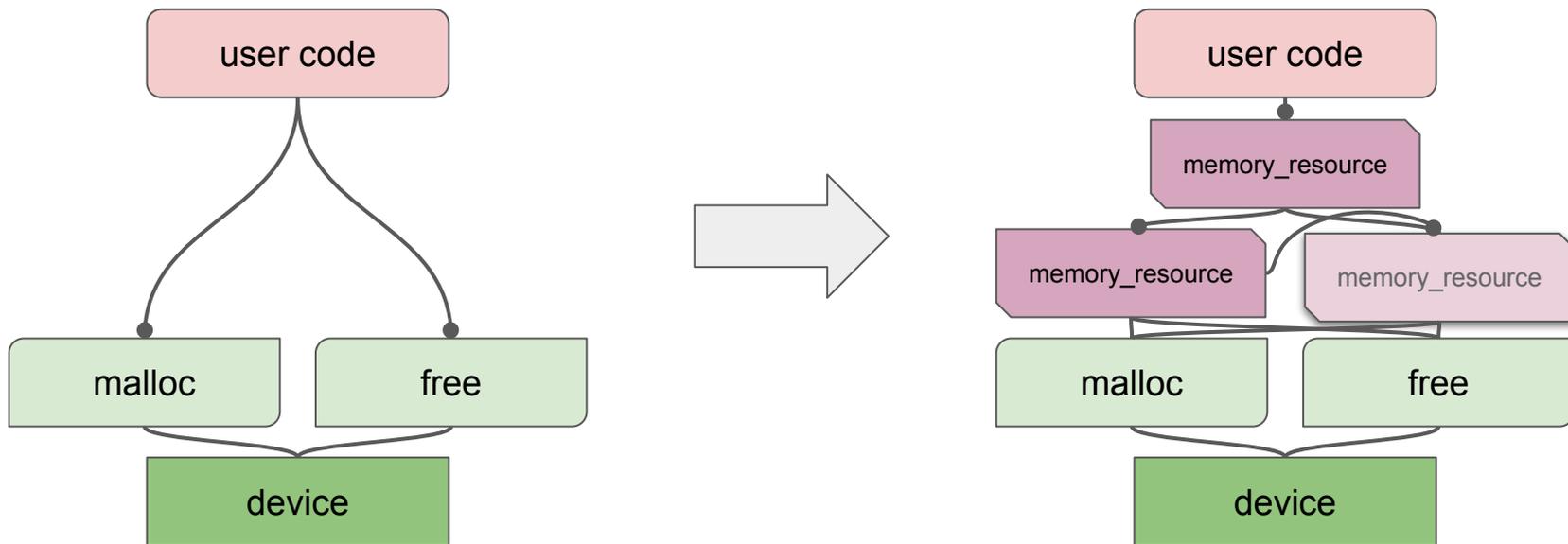
Memory resources	
Memory resources implement memory allocation strategies that can be used by <code>std::pmr::polymorphic_allocator</code>	
Defined in header <code><memory_resource></code> Defined in namespace <code>std::pmr</code>	
<code>memory_resource</code> (C++17)	an abstract interface for classes that encapsulate memory resources (class)
<code>new_delete_resource</code> (C++17)	returns a static program-wide <code>std::pmr::memory_resource</code> that uses the global operator <code>new</code> and operator <code>delete</code> to allocate and deallocate memory (function)
<code>null_memory_resource</code> (C++17)	returns a static <code>std::pmr::memory_resource</code> that performs no allocation (function)
<code>get_default_resource</code> (C++17)	gets the default <code>std::pmr::memory_resource</code> (function)
<code>set_default_resource</code> (C++17)	sets the default <code>std::pmr::memory_resource</code> (function)
<code>pool_options</code> (C++17)	a set of constructor options for pool resources (class)
<code>synchronized_pool_resource</code> (C++17)	a thread-safe <code>std::pmr::memory_resource</code> for managing allocations in pools of different block sizes (class)
<code>unsynchronized_pool_resource</code> (C++17)	
<code>monotonic_buffer_resource</code> (C++17)	

std::pmr::memory_resource	
Defined in header <code><memory_resource></code> <code>class memory_resource;</code> (since C++17)	
The class <code>std::pmr::memory_resource</code> is an abstract interface to an unbounded set of classes encapsulating memory resources.	
Member functions	
(constructor) (implicitly declared)	constructs a new <code>memory_resource</code> (public member function)
(destructor) [virtual]	destructs an <code>memory_resource</code> (virtual public member function)
<code>operator=</code> (implicitly declared)	implicitly declared copy assignment operator (public member function)
Public member functions	
<code>allocate</code>	allocates memory (public member function)
<code>deallocate</code>	deallocates memory (public member function)
<code>is_equal</code>	compare for equality with another <code>memory_resource</code> (public member function)
Private member functions	
<code>do_allocate</code> [virtual]	allocates memory (virtual private member function)
<code>do_deallocate</code> [virtual]	deallocates memory (virtual private member function)
<code>do_is_equal</code> [virtual]	compare for equality with another <code>memory_resource</code> (virtual private member function)

- C++17 introduced a set of new classes/interfaces in the `std::pmr` namespace for generalised memory management
 - With a design having its roots in [Thrust](#)
- Unfortunately it is taking some time to be implemented in `lib(std)c++`
 - [libstdc++](#) implemented full support in version 9.1
 - [libc++](#) still did not fully implement it
 - We are working around that in `VecMem`

Allocation Through Composition

- STL classes can perform all their memory allocations through the abstract [std::pmr::memory_resource](#) interface
 - But implementations of that interface do not need to do all the work by themselves. They can also delegate work to other resources through the same abstract interface.



- One can categorise memory resources into two groups. We use the following nomenclature for now.
- Upstream resources, which talk “directly” to some device/memory
 - We have implementations for all (“main”) types of memory that CUDA/HIP/SYCL have to offer
- Downstream resources, which only talk to other resources
 - We have multiple caching resources, some that can perform “logic operations”, and some with “side effects”
- Technically “hybrid” resources could also be made, but they should not be needed. Everything can be expressed with just the two types that we have at the moment.

Memory Resource Usage Examples

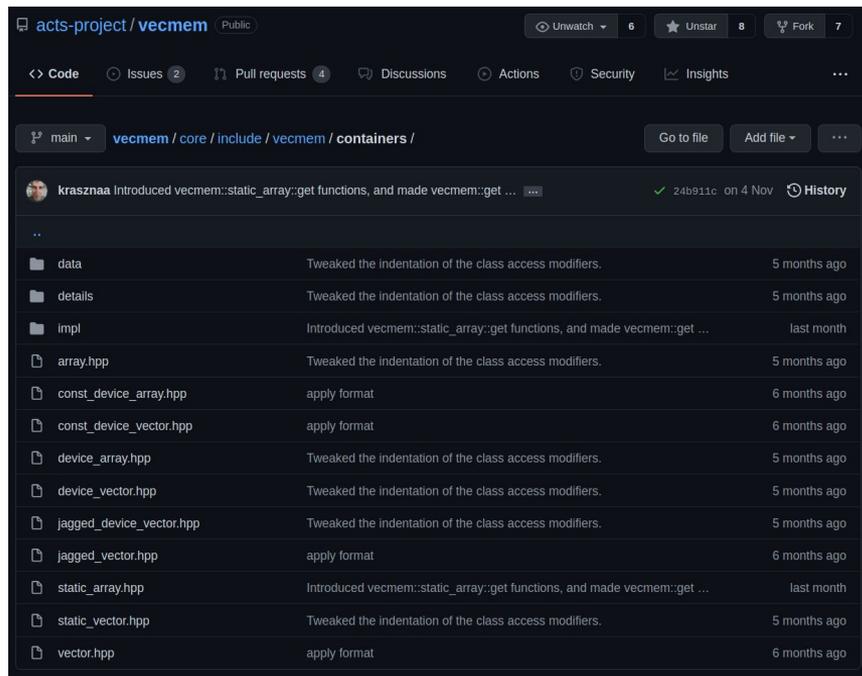


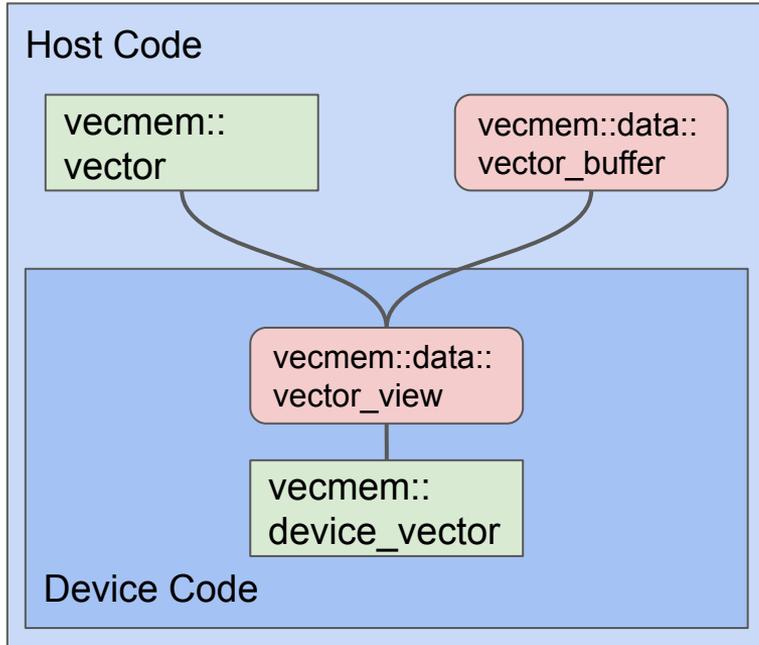
```
vecmem::cuda::managed_memory_resource cuda_mr;  
vecmem::arena_memory_resource arena_mr(cuda_mr);  
  
vecmem::vector<float> floats(&arena_mr);  
floats.resize(100, 5.f);  
  
superKernel<<<100, 1>>>(vecmem::get_data(floats));
```

```
vecmem::host_memory_resource host_mr;  
vecmem::cuda::device_memory_resource device_mr;  
  
vecmem::vector<int> ints(&host_mr);  
ints.resize(100, 5);  
  
vecmem::cuda::copy copy;  
vecmem::data::vector_buffer<int> device_ints =  
    copy.to(vecmem::get_data(ints), device_mr);  
  
hyperKernel<<<100, 1>>>(vecmem::get_data(device_ints));
```

Device Containers

- Managing device memory in host code using an STL syntax is just half of the story
- We also want to be able to interact with this memory in device code, with the same STL syntax
 - Since the (current) STL containers always assume that they manage their own memory, they can not be used in device code
 - We wrote a set of “device containers” that can interact with the data allocated by their “host container” counterparts





- We tried to formalise the exchange of information between host and device code using explicit data types
 - Types that are able to “point at” data owned by some other object
- The “device containers” can then be constructed on top of an appropriate data object
- “View type” objects are trivially constructible/copyable, and can be passed directly to device kernels
 - But we also introduced memory managing data types, which can only be instantiated in host code

Device Code Example



```
__global__
void linearTransformKernel(vecmem::vector_view<float> vec_data, float a, float b) {
    const std::size_t i = blockIdx.x * blockDim.x + threadIdx.x;
    if (i >= vec_data.size()) {
        return;
    }
    vecmem::device_vector<float> vec(vec_data);
    vec.at(i) = vec[i] * a + b;
}

void linearTransform(vecmem::vector<float>& vec) {
    linearTransformKernel<<<vec.size(), 1>>>(vecmem::get_data(vec), 2.5f, 3.5f);
}
```

On-Device Container Modifications



- Apart from allowing convenient access to the elements of fixed sized containers, we also allow for appending elements to resizable vectors
 - This only works with separate host/device memory management, as we can not change the memory under an [std::vector](#)'s feet
- In our current algorithms this comes in very handy in many situations

```
__global__  
void fillVector(  
    vecmem::data::vector_view<float> vec_data) {  
    ...  
    vecmem::device_vector<float> vec(vec_data);  
    vec.push_back(<foo>);  
    ...  
}  
  
{  
    vecmem::cuda::host_memory_resource host_mr;  
    vecmem::cuda::device_memory_resource device_mr;  
    vecmem::cuda::copy copy;  
  
    vecmem::data::vector_buffer<float>  
        device_buffer(MAX_SIZE, 0, device_mr);  
    copy.setup(device_buffer);  
  
    fillVector<<<...>>>(device_buffer);  
  
    vecmem::vector<float> host_vector;  
    copy(device_buffer, host_vector);  
}
```

Complex Types/Containers



```
55 /// Container describing objects in a given event
56 ///
57 /// This is the generic container of the code, holding all relevant
58 /// information about objects in a given event.
59 ///
60 /// It can be instantiated with different vector types, to be able to use
61 /// the same container type in both host and device code.
62 ///
63 /// It also can be instantiated with different edm types represented by
64 /// header and item type.
65 template <typename header_t, typename item_t,
66         template <typename> class vector_t,
67         template <typename> class jagged_vector_t>
68 class container {
69 public:
70     /// @name Type definitions
71     /// @{
72     291 // Convenience declaration for the container type to use in host code
73     292 template <typename header_t, typename item_t>
74     293 using host_container =
75     294     container<header_t, item_t, vecmem::vector, vecmem::jagged_vector>;
76     295
77     296 // Convenience declaration for the container type to use in device code
78     297 template <typename header_t, typename item_t>
79     298 using device_container = container<header_t, item_t, vecmem::device_vector,
80     299     vecmem::jagged_device_vector>;
81     300
82     301 // @name Types used to send data back and forth between host and device code
83     302 /// @{
84     303
85     304 // The header vector type
86     305 using header_vector = vector_t<header_t>;
87     306 // The item vector type
88     307 using item_vector = jagged_vector_t<item_t>;
89     308
90     309 // Structure holding (some of the) data about the container in host code
91     310 template <typename header_t, typename item_t>
92     311 struct container_data {
93     312     vecmem::data::vector_view<header_t> headers;
94     313     vecmem::data::jagged_vector_data<item_t> items;
95     314 };
96
97     315 // Structure holding (all of the) data about the container in host code
98     316 template <typename header_t, typename item_t>
99     317 struct container_buffer {
100     318     vecmem::data::vector_buffer<header_t> headers;
101     319     vecmem::data::jagged_vector_buffer<item_t> items;
102     320 };
103
104     321 // @brief The type name of the container
105     322 * methods in this class.
106     323 */
107     324 using element_view =
108     325     container_element<header_t,
109     326     container_data<header_t, item_t>,
110     327     container_buffer<header_t, item_t>>;
111
112     328 }
```

- In [detray](#) and [traccc](#) we define some types that we want to use both as “host” and as “device” objects / containers
 - Instead of duplicating definitions for the different vector types, we made these data types templated on the vector type that they should use
 - Allowing us to write even algorithmic code in a way that it could be used directly on top of host vectors, or in device code on device vectors

Miscellaneous Types

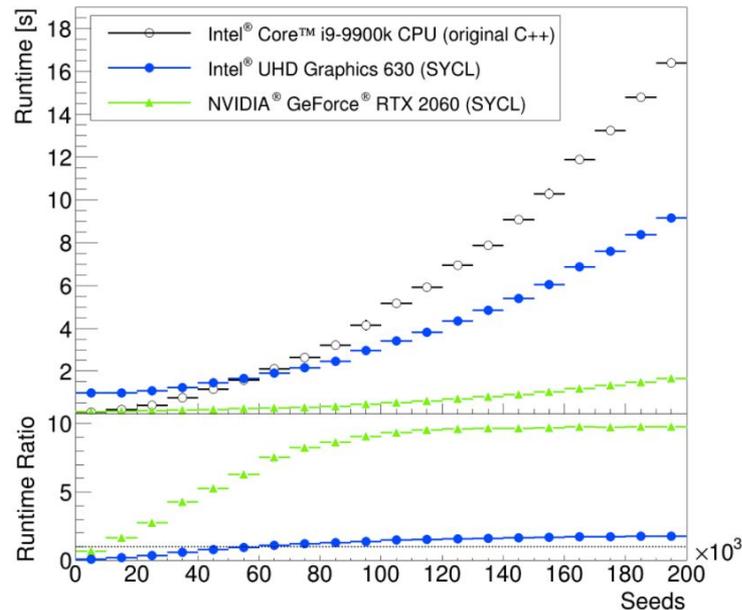
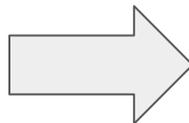
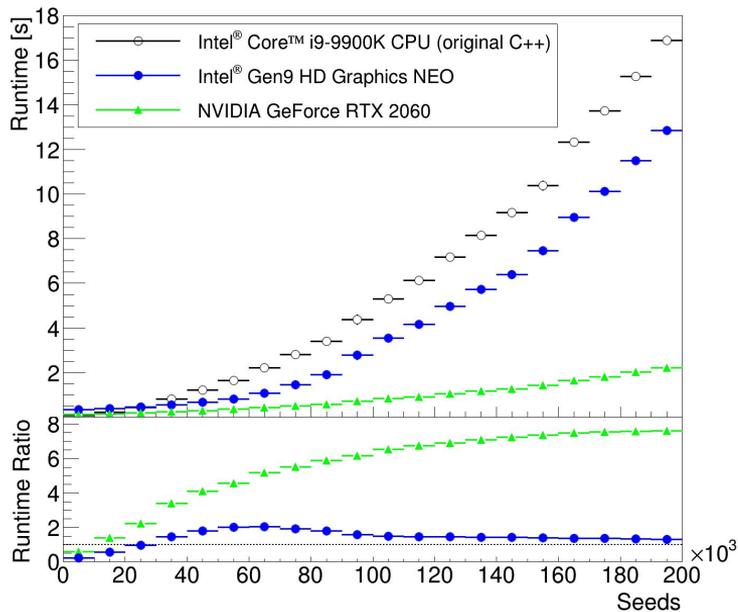
- We have also defined a number of other, hopefully useful types
 - Containers with a compile-time (maximum) size
 - [vecmem::allocator](#) for “non-container” memory allocations
 - [vecmem::atomic](#) for abstracting atomic operations for all supported languages
 - [vecmem::memory_monitor](#) for collecting simple information about an application’s behaviour
- Mainly driven by what we could think of, and what we ended up needing for the tracking demonstrator

```

18 namespace vecmem {
19
20 /// Class mimicking @c std::vector on top of a fixed sized array
21 ///
22 /// This can come in handy when needing vector arithmetics in device code,
23 /// without resorting to heap allocations.
24 ///
25 /// The type does come with a significant limitation over @c std::vector.
26 /// It has a maximal/fixed size that needs to be chosen at compile time.
27 ///
28 template <typename TYPE, std::size_t MAX_SIZE>
29 class static_vector {
30
31 public:
32     /// @name Type definitions
33     /// @{
34
35     /// Type of the array elem
36     typedef TYPE value_type;
37     /// Size type for the arra
38     typedef std::size_t size_t
39     /// Pointer difference typ
40     typedef std::ptrdiff_t dif
41
42     /// The maximal size of th
43     static constexpr size_type
44     /// The size of the vector
45     static constexpr size_type
46     /// Type of the array hold
47     typedef
48     typename details::arra
49     array_type;
50
51     /// Value reference type
52     typedef value_type& refere
53     /// Constant value referen
54     typedef const value_type&
55     /// Value pointer type
56     typedef value_type* pointe
57     /// Constant value pointer
58     typedef const value_type*
59
60     namespace vecmem {
61     /// Class collecting some basic set of memory allocation statistics
62     /// Objects of this class can be used together with
63     /// @c vecmem::instrumenting_memory_resource to easily access a common set of
64     /// useful performance metrics about an application.
65     /// Note that the lifetime of this object must be at least as long as the
66     /// lifetime of the connected memory resource!
67     class VECMEM_CORE_EXPORT memory_monitor {
68     public:
69         /// Constructor with a memory resource reference
70         memory_monitor(instrumenting_memory_resource& resource);
71
72         /// Get the total amount of allocations
73         std::size_t total_allocation() const;
74         /// Get the outstanding allocation left after all operations
75         std::size_t outstanding_allocation() const;
76         /// Get the average allocation size
77         std::size_t average_allocation() const;
78         /// Get the maximal concurrent allocation
79         std::size_t maximal_allocation() const;
80     };
81     };
82
83     };
84 };

```

Memory Caching in Action



- A lot of interesting features are coming in C++2X for GPU programming
 - We jumped on “memory resources” as something that is already available. But [std::atomic_ref](#), [std::mdspan](#) and [std::execution](#) will all be very relevant for HEP code. We should be mindful of these upcoming language features when writing our accelerated code already today.
- We are still very much doing R&D with this code in the context of the ATLAS GPU tracking demonstrator
 - With the goal of merging this kind of memory management into ATLAS’s offline software before Run-4
- If you are at all interested in giving the code a try, don’t be shy to get in contact with us
 - Hopefully next year we will have some news about interfacing [LLAMA](#) with VecMem as well 😊



<http://home.cern>