



# **traccc** - A (Close To) Single-Source Tracking Demonstrator on CPUs/GPUs

#### Attila Krasznahorkay on behalf of the Acts Parallelization R&D Team



## The Acts Project

C Edit on GitHub



Search docs

Core library

Plugins

Examples

File formats

Authors

License

The <u>Acts project</u> aims to provide HEP / NP experiments with a toolbox for charged particle reconstruction

- It allows experiment software to have its 0 own specific event data types, which can be used with zero copy in the Acts components
- Its development is heavily influenced by the ATLAS Experiment at this stage, but it is truly meant to be a collaborative project across experiments
  - It is on track to be used by ATLAS in all of 0 its track reconstruction by the LHC's Run-4

# The HL-LHC Computing Challenge

- The LHC experiments will collect events of much higher complexity with much higher rate than ever before in the <u>High Luminosity LHC</u> era
- The computing requirements of "classical reconstruction algorithms" increase non-linearly with event complexity in many cases
  - Charged particle reconstruction being a dominant part of this
- In order to tackle this challenge new types of event reconstruction methods need to be tried, including using devices beyond classical CPUs



# Track Reconstruction at High Pileup





Currently used algorithms were primarily designed for relatively low complexity collision events

## Track Reconstruction at High Pileup



here for the ATLAS ITk geometry)

## **Track Reconstruction in Acts**



CKF chain



## Track Reconstruction in Acts



ERI

## The Acts Parallelization R&D



- The development happens independently from the main Acts repository to make developments quicker
- Broken up into multiple, task specific projects
  - vecmem: Common memory management
  - <u>covfie</u>: Generic vector field handling
  - <u>algebra-plugins</u>: Small matrix linear algebra abstractions
  - <u>detray</u>: Tracking geometry handling in device code
  - <u>traccc</u>: The main repository of the R&D effort, holding most algorithmic code

# **Implementing Performance Portability**



- While the CPU and GPU algorithms themselves are implemented separately, they do share a lot of functions
- These all need to be implemented **inline** to:
  - Allow the host and device compilers to generate code from them as they need it;
  - Allow the same function to be compiled into multiple object files (with different compilers/flags) during the build.

```
namespace traccc {
18
19
    /// helper functions (for both cpu and gpu) to perform conformal transformation
20
21
    /// @param x is the x value
22
    /// Oparam y is the y value
24
    /// @return is the conformal transformation result
    inline TRACCC_HOST_DEVICE vector2 uv_transform(const scalar& x,
26
                                                     const scalar& y) {
27
        vector2 uv:
28
        scalar denominator = x * x + y * y;
29
        uv[0] = x / denominator;
30
        uv[1] = y / denominator;
31
        return uv;
32
33
34
    /// helper functions (for both cpu and gpu) to calculate bound track parameter
35
    /// at the bottom spacepoint
36
37
    /// Oparam seed is the input seed
38
    /// Oparam bfield is the magnetic field
39
    /// Oparam mass is the mass of particle
40
    template <typename spacepoint_collection_t>
41
    inline TRACCC_HOST_DEVICE bound_vector seed_to_bound_vector(
42
        const spacepoint_collection_t& sp_collection, const seed& seed,
43
        const vector3& bfield, const scalar mass) {
44
45
        bound_vector params;
46
47
        const auto& spB = sp collection.at(seed.spB link);
48
        const auto& spM = sp_collection.at(seed.spM_link);
49
        const auto& spT = sp collection.at(seed.spT link);
50
51
        darray<vector3, 3> sp_global_positions;
52
        sp_global_positions[0] = spB.global;
53
        sp_global_positions[1] = spM.global;
         sp_global_positions[2] = spT.global;
```

# Implementing Performance Portability



/// CUDA kernel for running @c traccc::device::find\_doublets 52 \_\_global\_\_ void find\_doublets( 53 54 seedfinder\_config config, sp\_grid\_const\_view sp\_grid, 55 device::doublet\_counter\_collection\_types::const\_view doublet\_counter, 56 device::device doublet collection types::view mb doublets, 57 device::device\_doublet\_collection\_types::view mt\_doublets) { 58 59 device::find\_doublets(threadIdx.x + blockIdx.x \* blockDim.x, config, sp\_grid, doublet\_counter, mb\_doublets, mt\_doublets); 60 61

150	// Find all of the spacepoint doublets.				
151	<pre>device::device_doublet_collection_types::view mb_view = doublet_buffer_mb;</pre>				
152	<pre>device::device_doublet_collection_types::view mt_view = doublet_buffer_mt;</pre>				
153	<pre>auto find_doublets_kernel =</pre>				
154	<pre>details::get_queue(m_queue).submit([&amp;](::sycl::handler&amp; h) {</pre>				
155	h.parallel_for <kernels::find_doublets>(</kernels::find_doublets>				
156	doubletFindRange,				
157	<pre>[config = m_seedfinder_config, g2_view, doublet_counter_view,</pre>				
158	<pre>mb_view, mt_view](::sycl::nd_item&lt;1&gt; item) {</pre>				
159	<pre>device::find_doublets(item.get_global_linear_id(), config,</pre>				
160	g2_view, doublet_counter_view,				
161	<pre>mb_view, mt_view);</pre>				
162	<pre>});</pre>				
163	<pre>});</pre>				

 Language specific kernel launches call on shared functions for the heavy lifting

21	<pre>namespace traccc::device {</pre>		
22			
23	/// Function finding all of the spacepoint doublets		
24	///		
25	/// Based on the information collected by @c traccc::device::count_doublets it		
26	/// can fill collection with the specific doublet pairs that exist in the event.		
27	111		
28	/// @param[in] globalIndex The index of the current thread		
29	/// @param[in] config Seedfinder configuration		
30	/// @param[in] sp_view The spacepoint grid to find doublets on		
31	/// @param[in] dc_view Collection with the number of doublets to find		
32	/// @param[out] mb_doublets_view Collection of middle-bottom doublets		
33	/// @param[out] mt_doublets_view Collection of middle-top doublets		
34	///		
35	TRACCC_HOST_DEVICE		
36	inline void find_doublets(		
37	<pre>std::size_t globalIndex, const seedfinder_config&amp; config,</pre>		
38	<pre>const sp_grid_const_view&amp; sp_view,</pre>		
39	<pre>const doublet_counter_collection_types::const_view&amp; dc_view,</pre>		
40	<pre>device_doublet_collection_types::view mb_doublets_view,</pre>		
41	<pre>device_doublet_collection_types::view mt_doublets_view);</pre>		
42			
43	} // namespace traccc::device		

# Floating Point Precision



Running ./bin/traccc seq example sycl tml detector/trackml-detector.csv tml full/ttbar mu140/ 1 Running Seeding on device: NVIDIA GeForce RTX 3080 ===>>> Event 0 <<<=== Number of spacepoints: 67442 (host), 67442 (device) Matching rate(s): - 97.9256% at 0.01% uncertainty - 99.3046% at 0.1% uncertainty - 99.3076% at 1% uncertainty - 99.3076% at 5% uncertainty Number of seeds: 11514 (host), 11514 (device) Matching rate(s): - 86.5902% at 0.01% uncertainty - 98.7233% at 0.1% uncertainty - 99.2705% at 1% uncertainty - 99.3399% at 5% uncertainty Number of track parameters: 11514 (host), 11514 (device) Matching rate(s): - 98.3933% at 0.01% uncertainty **FP32** - 99.8002% at 0.1% uncertainty - 99.835% at 1% uncertainty - 99.835% at 5% uncertainty Running ./bin/traccc seq example sycl tml detector/trackml-detector.csv tml full/ttbar mul40/ Running Seeding on device: NVIDIA GeForce RTX 3080 ===>>> Event 0 <<<=== Number of spacepoints: 67442 (host), 67442 (device) Matching rate(s): - 100% at 0.01% uncertainty - 100% at 0.1% uncertainty - 100% at 1% uncertainty - 100% at 5% uncertainty Number of seeds: 11515 (host), 11515 (device) Matching rate(s): - 100% at 0.01% uncertainty - 100% at 0.1% uncertainty - 100% at 1% uncertainty - 100% at 5% uncertainty Number of track parameters: 11515 (host), 11515 (device) Matching rate(s): - 100% at 0.01% uncertainty FP64 - 100% at 0.1% uncertainty - 100% at 1% uncertainty - 100% at 5% uncertainty

- With the minimal vectorization that HEP code typically has, modern x86\_64 CPUs have (virtually) the same performance for FP32 and FP64 operations
  - However accelerators usually do not
- On the other hand FP64 operations generally provide much better agreement between CPU and GPU algorithms (5)
- All projects are set up to use a user-defined floating point type
  - For the moment higher level projects select the type on a project level, but it shall be possible to select the precision algorithm-by-algorithm in the future

# Status of the Project



- Many programming techniques tried for GPUs during development
  - But only seriously using <u>CUDA</u> and <u>SYCL</u> for now
- The most complicated part, CKF, is under heavy development at the moment
  - It will be the final word on whether the project would fully succeed
- Track fitting works, but is not integrated into the "full chain" of algorithms at this point

Category	Algorithms	CPU	CUDA	SYCL	Futhark
Clusterization	CCL	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
	Measurement creation	V	$\checkmark$	$\checkmark$	$\checkmark$
	Spacepoint formation	V	V	V	
Track finding	Spacepoint binning	$\checkmark$	$\checkmark$	$\checkmark$	
	Seed finding	V	$\checkmark$	$\checkmark$	
	Track param estimation	V	$\checkmark$	$\checkmark$	
	Combinatorial KF	•			
Track fitting	KF	$\checkmark$	$\checkmark$	$\checkmark$	

# Measuring Performance



79	full_chain_algorithm	::output_type full_chain_algorithm::operator()(					
80	<pre>const cell_collection_types::host&amp; cells,</pre>						
81	<pre>const cell_module_collection_types::host&amp; modules) const {</pre>						
82							
83	// Create device copy of input collections						
84	cell_collection_	<pre>cell_collection_types::buffer cells_buffer(cells.size(),</pre>					
85		<pre>*m_cached_device_mr);</pre>					
86	m_copy(vecmem::g	::get_data(cells), cells_buffer);					
87	cell_module_coll	cell_module_collection_types::buffer modules_buffer(modules.size(),					
88	*m_cached_device_mr);						
89	m_copy(vecmem::g	et_data(modules), modules_buffer);					
90							
91	// Run the clust	erization (asynchronously).					
92	const clusteriza						
93	m_clusteriza	ulti-threaded CUDA GPU throughput tests					
94	const track_para	>>> Throughput options <<<					
95	m_track_para	Input data format : csv					
96		Input directory : tml_full/ttbar_mu140/					
97		Detector geometry : tml_detector/trackml-detector.csv					
98	// Get the final	Target cells per partition : 1024					
99	bound_track_para	Loaded event(s) : 10					
100	m_copy(track_par	Cold run event(s) : 10					
101	m_stream.synchro	Processed event(s) : 10000					
102		>>> Multi-threading options <<<					
103	// Return the ho	CPU threads: 4					
104	return result;						
105	}	Using CUDA device: NVIDIA GEForce RTX 3080 [id: 0, bus: 1, device: 0]					
		Using CUDA device: NVIDIA GeForce RTX 3080 [id: 0, bus: 1, device: 0]					
_		Using CUDA device: NVIDIA GeForce RTX 3080 [id: 0, bus: 1, device: 0]					
		Using CUDA device: NVIDIA GeForce RTX 3080 [id: 0, bus: 1, device: 0]					
		Reconstructed track parameters: 104221093					
		Time totals:					
		File reading 8546 ms					
		Warm-up processing 18 ms					
		Throughput:					
		Warm-up processing 1.89721 ms/event, 527.091 events/s					
		Event processing 1.07808 ms/event, 927.578 events/s					

- On top of the many tests for the correctness of the output of the algorithms, we also build executables testing the throughput of the algorithms
  - Pre-load "cell data" for a set number of events into host memory;
  - Process the data into track parameters on possibly multiple CPU threads (using <u>TBB</u>);
  - Copy the results back into host memory if needed.

## Measuring Performance (Precision)



With FP64 operations one needs (with our current code) a very high-end GPU to compete with a high-end CPU. With FP32 operations lower-end GPUs perform much better.

# Measuring Performance (Language)



With our current code CUDA performs better at low  $\mu$ . (At high  $\mu$  the difference is insignificant.) Though we know about current inefficiencies in both implementations.

# Measuring Performance (Summary)



GPUs become competitive at high pile-up. Highest performance observed on NVIDIA<sup>®</sup> workstation GPUs so far.





- The Acts Parallelization R&D is a significant effort, providing one of the demonstrators for the ATLAS HLT upgrade for the HL-LHC
  - Developments will start soon on exercising the code with the simulations / geometry of the ATLAS ITk
- The performance of the code is promising, providing a higher throughput with attainable GPUs than with the fastest CPUs that we could run tests on
  - What makes sense to use of course also very much depends on price and power usage, which we are not testing for / including in our results at the moment
- One significant development step still in the works: CKF



#### Performance





mu40



#### Performance

mu300





#### Performance

Peak performance across different hardware



ttbar pile-up



http://home.cern