



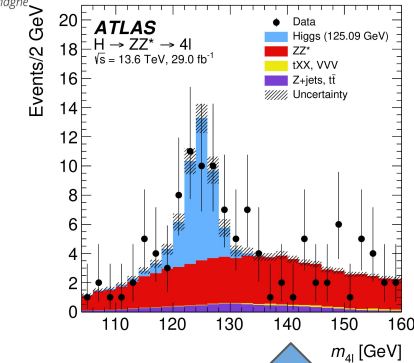
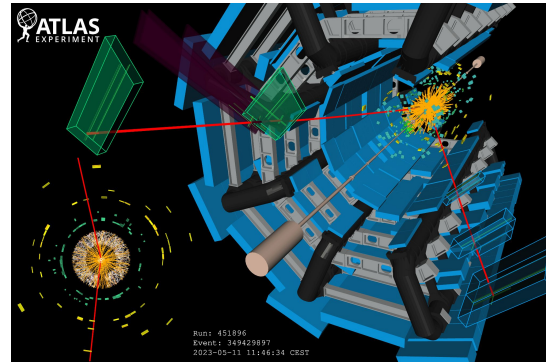
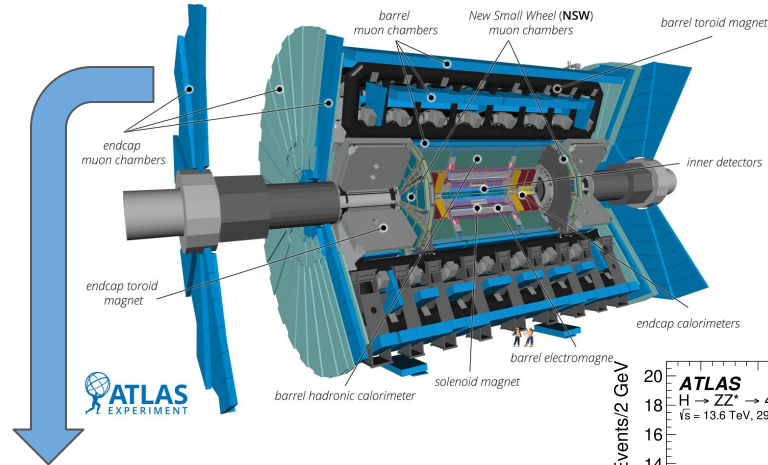
Accelerating HEP Software

Attila Krasznahorkay

trying to sample / expose the work of a lot of people...



- “Software” is used practically everywhere in HEP
 - But here I’ll only focus on software used in “data processing”. Software used for controlling the accelerator, detector hardware, etc. is a different beast entirely.
- The goal is to “do physics”
 - Fast and nicely written software is the means, not the actual goal of the endeavour.



Software ↔ Computing

The image shows a terminal window on the left and a Grafana dashboard on the right. The terminal window displays the setup of an Athena container and the output of the Application Manager, which is running on Celborn. The Grafana dashboard shows a stacked area chart titled 'Slots of Running jobs by Activity' for the last 30 days. The chart shows various activities such as MC Simulation Full, Group Production, MC Event Generation, MC Reconstruction, User Analysis, MC Simulation Fast, Group Analysis, Data Processing, and I/O processing. A table at the bottom of the dashboard provides summary statistics for each activity.

Activity	min	max	avg	current
MC Simulation Full	132 K	155 M	461 K	527 K
Group Production	7.85 K	235 K	95.8 K	18.6 K
MC Event Generation	20.1 K	181 K	86.5 K	25.5 K
MC Reconstruction	24.5 K	131 K	59.9 K	100 K
User Analysis	12.1 K	116 K	59.7 K	28.5 K
MC Simulation Fast	0	214 K	20.9 K	136 K
Group Analysis	0	485.4 K	12.9 K	35.2 K
Data Processing	0	570 K	10.3 K	94.7
I/O processing	0	38.7 K	1.64 K	582

- Even in data processing, we have 2 distinct areas
 - “Software” covers all the software written to perform the data transformations needed “for physics”.
 - “Computing” covers all the software created to allow the first category to execute on up to millions of CPU cores in parallel.
- Here I will only cover “software”
 - Even though using accelerators in job/network optimization, and even just making batch systems “accelerator aware” is a very important thing on our plates.

Development, Build and Deployment



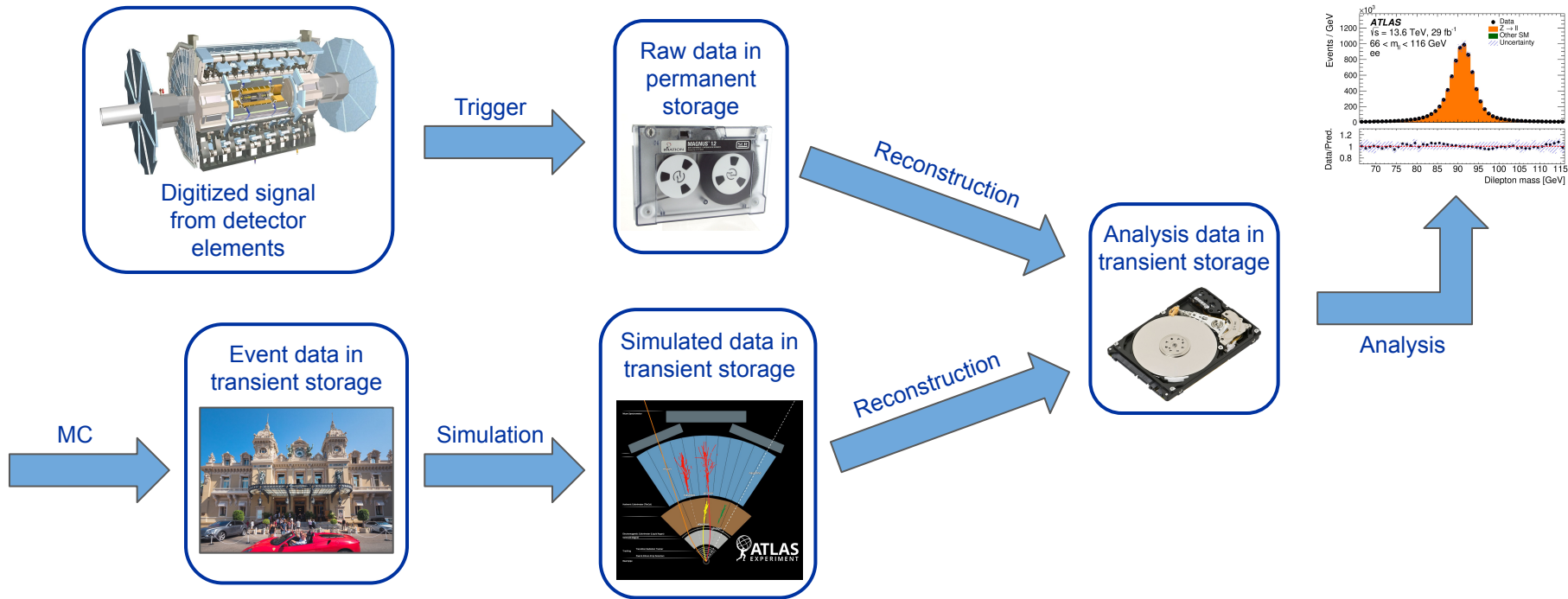
- Interpreting / processing the data coming from the LHC detectors takes many millions of C++ / Python code
 - We use very customized build systems to allow developers to deal with a small portion of the software at any given time, testing its integration with the rest of the software.
 - This requires highly curated, many GB development environments to be made available around the world.
- The distribution uses CVMFS
 - Which needs to be able to house the full SDKs (making licensing a challenge...)

The top screenshot shows the GitHub repository for 'athena', the ATLAS Experiment's main offline software repository. It displays 127,815 commits, 74 branches, 2,752 tags, and 441 releases. A recent commit by Edward Mysse is highlighted.

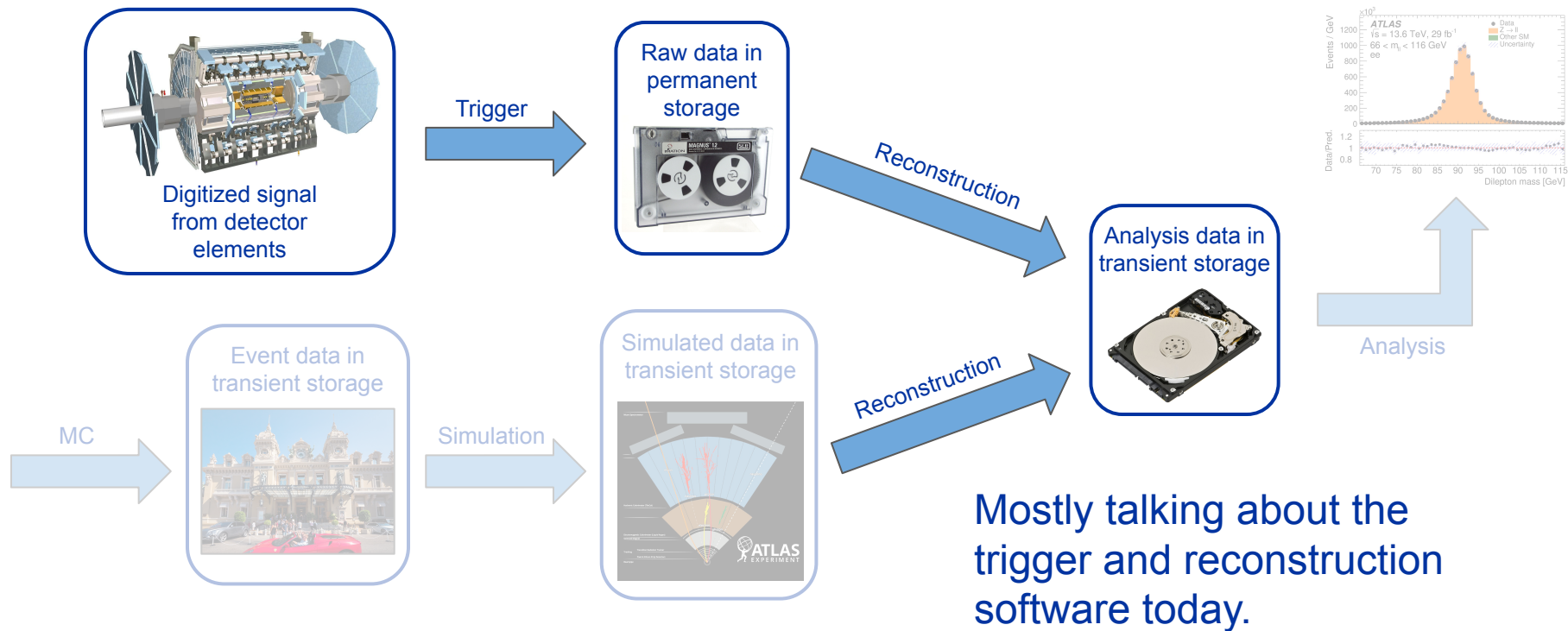
The bottom screenshot shows the PanDA monitor dashboard for ATLAS Nightlies and CI Global Page. It displays a table of build results for various branches and versions.

Nightly Group	Branch	Recent Release	Build time	Compilation errors (warnings)	C Test (or CI) test status (warnings)	ART LOCAL	ART GRID	CVMFS (on client)
CI	MR-CI-build	MR-70102-2024-03-25-19-11	25-MAR-17:56	0 (0)	0 (0)	N/A	N/A	N/A
PRODUCTION	24_0_Athena_n98_64-68-gcc13-dbg	2024-03-23T0301	23-MAR-05:52	0 (0)	13 (13)	N/A	N/A	23-MAR-09:51
PRODUCTION	24_0_Athena_n98_64-68-gcc13-opt	2024-03-24T2101	24-MAR-23:30	0 (0)	0 (0)	7.6.6	0.198.93.9	25-MAR-01:11
PRODUCTION	24_0_AthDemulation_n98_64-68-gcc13-opt	2024-03-22T2101	23-MAR-22:49	0 (0)	0 (0)	3.0.0	0.35.0.1	22-MAR-23:21
PRODUCTION	24_0_DeCommon_n98_64-68-gcc13-opt	2024-03-24T2101	24-MAR-21:04	0 (0)	0 (0)	N/A	N/A	24-MAR-21:22
PRODUCTION	main_AnalysisBase_n98_64-centos7-gcc11-opt	2024-03-19T0220	19-MAR-05:57	0 (0)	0 (0)	N/A	0.5.2.0	19-MAR-04:42
PRODUCTION	main_AnalysisBase_n98_64-68-gcc13-opt	2024-03-25T0220	25-MAR-02:58	0 (0)	0 (0)	N/A	0.6.2.0	25-MAR-04:52
PRODUCTION	main_AthAnalysis_n98_64-centos7-gcc11-opt	2024-03-19T0001	19-MAR-00:38	0 (0)	0 (0)	N/A	0.0.0.0	19-MAR-01:21
PRODUCTION	main_AthAnalysis_n98_64-68-gcc13-opt	2024-03-25T0001	25-MAR-00:29	0 (0)	0 (0)	N/A	0.6.0.0	25-MAR-01:11

HEP Data Processing (1)

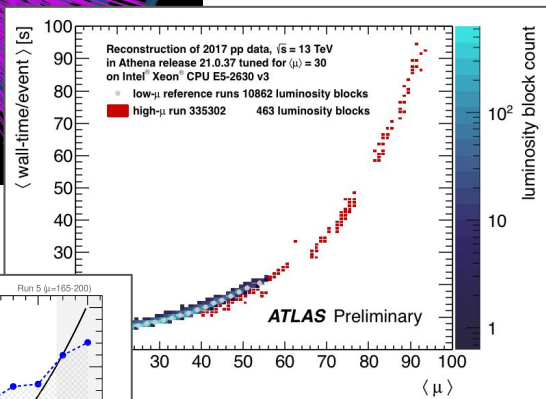
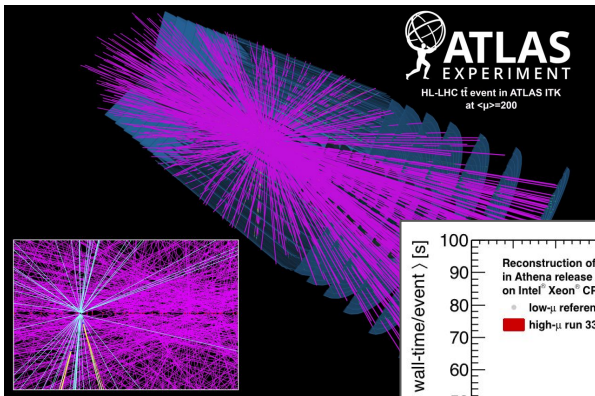


HEP Data Processing (2)

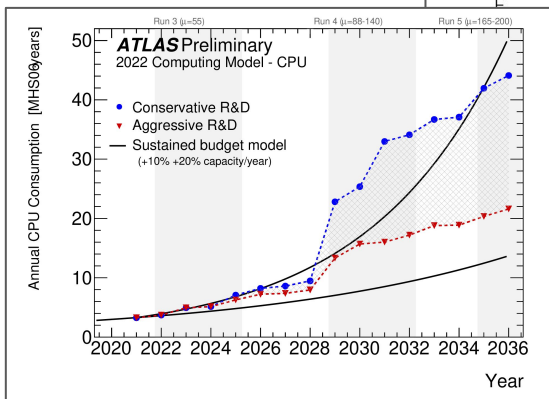


Mostly talking about the trigger and reconstruction software today.

Trigger / Reconstruction



- The goal of this software is to figure out what particles, with what exact properties, left signals in the detector
 - A little similar to “finding the cats” in a CCD camera’s recorded data 🐱
 - In trillions of >100 megapixel, “3D camera” pictures
- “Classical” algorithms do this very well in the ongoing LHC data taking. (Though ATLAS is the only one **only** using classical algorithms at this point.)
 - However during the High-Luminosity LHC era these current algorithms would need too much CPU power.



Things To Do Right: Event Data Model



- Our software uses “rich” data
 - Many, well defined properties, in **very jagged arrays**
- Our 20+ year old existing software tends to manage these in ways not suited for accelerators (AoS vs. SoA)
- We need to move this data with a high frequency between $H \leftrightarrow D$
- Must be manageable to implement/maintain the code that describes all our data types
 - Note that a single “type” usually requires multiple C++ classes to describe it

```
1 #ifndef DataFormats_PortableTestObjects_Interface_TestSoA_h
2 #define DataFormats_PortableTestObjects_Interface_TestSoA_h
3
4 #include <array>
5
6 #include <Eigen/Core>
7 #include <Eigen/Dense>
8
9 #include "DataFormats/Common/interface/StdArray.h"
10 #include "DataFormats/SoATemplate/interface/SoACommon.h"
11 #include "DataFormats/SoATemplate/interface/SoALayout.h"
12 #include "DataFormats/SoATemplate/interface/SoAView.h"
13
14 namespace portabletest {
15
16 // the type aliases are needed because commas confuse macros
17 using Array = eds::StdArray<short, 4>;
18 //using Array = std::array<short, 4>;
19 using Matrix = Eigen::Matrix<double, 3, 6>;
20
21 // SoA layout with x, y, z, id fields
22 GENERATE_SOA_LAYOUT(TestSoALayout,
23 // columns: one value per element
24 SOA_COLUMN<double, x>,
25 SOA_COLUMN<double, y>,
26 SOA_COLUMN<double, z>,
27 SOA_COLUMN<int32_t, id>,
28 // scalars: one value for the whole structure
29 SOA_SCALAR<double, r>,
30 // column of arrays: one fixed-size array per element
31 SOA_COLUMN<Array, flags>,
32 // Eigen columns: each matrix element is stored in a separate column
33 SOA_EIGEN_COLUMN<Matrix, m>);
34
35 using TestSoA = TestSoALayout<>;
```

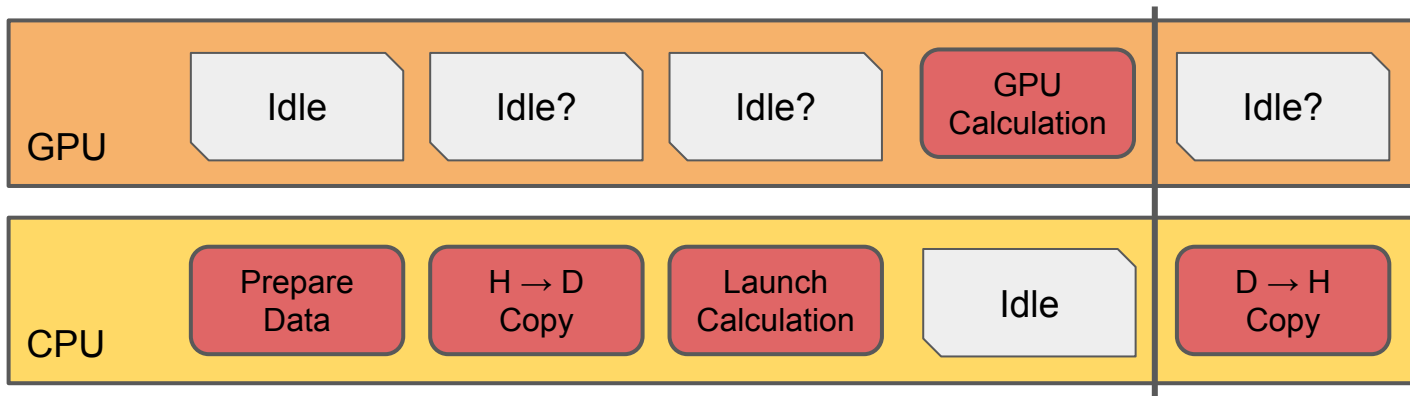
```
45 // Global "average" of something (non-const)
46 VECEM_HOST_AND_DEVICE
47 auto& average() { return BASE::template get<3>(); }
48 // Global "average" of something (const)
49 VECEM_HOST_AND_DEVICE
50 const auto& average() const { return BASE::template get<3>(); }
51
52 // "Indices" of something (non-const)
53 VECEM_HOST_AND_DEVICE
54 auto& indices() { return BASE::template get<4>(); }
55 // "Indices" of something (const)
56 VECEM_HOST_AND_DEVICE
57 const auto& indices() const { return BASE::template get<4>(); }
58
59 // "Index" of something (non-const)
60 VECEM_HOST_AND_DEVICE
61 auto& index() { return BASE::template get<>(); }
62 // "Index" of something (const)
63 VECEM_HOST_AND_DEVICE
64 const auto& index() const { return BASE::template get<>(); }
65
66 }; // class jagged_interface
67
68 // "Jagged" container for the tests
69 using jagged_soa_container =
70 eds::container<jagged_soa_interface, eds::type::scalar<int>,
71 eds::type::vector<float>, eds::type::jagged_vector<double>,
72 eds::type::scalar<float>, eds::type::jagged_vector<int>,
73 eds::type::vector<int>>;
74
75 } // namespace testing
76 } // namespace vecmem
```

C++ data layout portability
Thursday 30 Nov 2023, 16:00 → 17:30 Europe/Zurich
32/1-A24 (CERN)

Videoconference C++ data layout portability workshop

Time	Topic	Speaker	Duration
16:00 → 16:30	Data layout: our reality	Eric Cano (CERN)	30m
16:30 → 16:50	How reflection could help	Bernhard Manfred Gruber (Technische Universität Dresden (DE))	20m
16:50 → 17:30	Discussion		40m

Things To Do Right: Task Scheduling



- Not all aspects of our data processing are suited for accelerators
 - Any calculation not using an accelerator must be using the available CPU cores efficiently in parallel with whatever the GPU is doing, maximising the usage of all components
- A number of different solutions are used by the experiments currently
 - Built around [TBB](#) (for ATLAS and CMS) and [ZeroMQ](#) (for ALICE and LHCb)
- But new, hopefully better solutions seem to be on the horizon
 - Using [Boost::fiber](#) and/or [C++ coroutines](#)

Things To Do Right: Algorithms

```

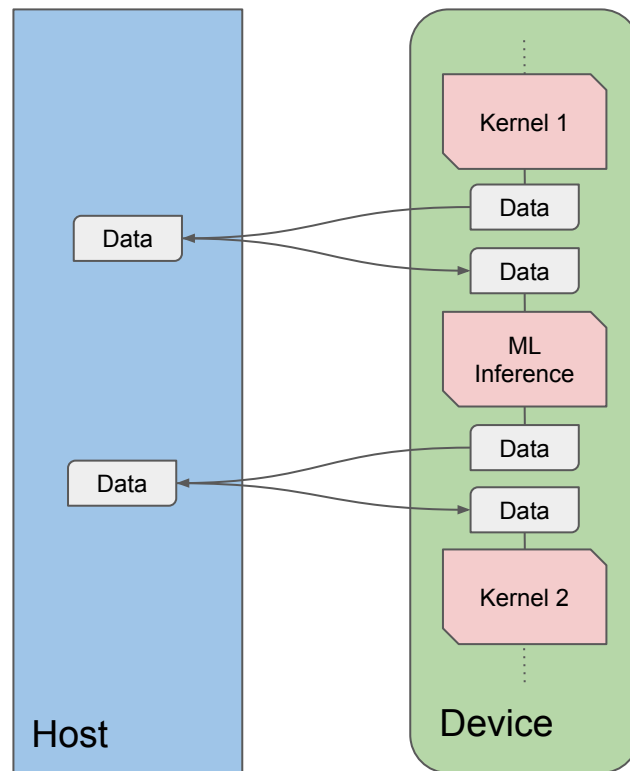
13 namespace traccs::device {
14
15 /// Implementation of a FastSV algorithm with the following steps:
16 /// 1) mix of stochastic and aggressive hooking
17 /// 2) shortcutting
18 ///
19 /// The implementation corresponds to an adapted version of Algorithm 3 of
20 /// the following paper:
21 /// https://www.sciencedirect.com/science/article/pii/S07433 39 _global_
22 /// 40 // __launch_bounds__(256,4)
23 /// 41 void
24 /// This array only gets updated at the 42 findclus(uint16_t const* restrict_id, // module id of each pixel
25 /// to prevent race conditions. 43 uint16_t const* restrict_x, // local coordinates of each pixel
26 /// The number of adjacent cells 44 uint16_t const* restrict_y, //
27 /// @param[in] adjc Vector of adjacent cells 45 uint32_t const* restrict_modulestart, // index of the first pixel of each module
28 /// @param[in] tid The thread index 46 uint32_t* restrict_nclustersModule, // output: number of clusters found in each module
29 /// @param[in] blkDim The block size 47 uint32_t* restrict_moduleId, // output: module id of each module
30 /// @param[inout] f array holding the parent cell ID for 48 int32_t* restrict_clusterId, // output: cluster id of each pixel
31 /// iteration. 49 int numElements) {
32 /// @param[inout] gf array holding grandparent cell ID fr 50 if (blockIdx.x >= modulestart[0])
33 /// iteration. 51 return;
34 /// @param[in] barrier A generic object for block-wide sync 52
35 /// 53 auto firstPixel = modulestart[1 + blockIdx.x];
36 /// 54 auto thisModuleId = id[firstPixel];
37 template <typename barrier_t> 55 assert(thisModuleId < MAXnumModules);
38 TRACCC_DEVICE void fast_sv_1(index_t* f, index_t* gf, 56
39 unsigned char adjc[MAX_CELLS_PE] #ifdef GNU_DEVUS
40 index_t adjv[MAX_CELLS_PER_THRE 57 if (thisModuleId % 100 == 1)
41 const index_t tid, const index 58 if (threadIdx.x == 0)
42 barrier_t& barrier) { 59 printf("start clusterizer for module %d in block %d\n", thisModuleId, blockIdx.x);
43 60 #endif
44 /* 61
45 * The algorithm finishes if an iteration leaves the arr 62
46 * This variable will be set if a change is made, and djc 63
47 * loop is necessary. 64
48 */ 65 // find the index of the first pixel not belonging to this module (or invalid)
49 bool gf_changed; 66 _shared int msize;
50 67 msize = numElements;
51 68 _syncthreads();
52 do { 69
53 /* 70 // skip threads not associated to an existing pixel
54 * Reset the end-parameter to false, so we can set i 71 for (int i = first; i < numElements; i += blockDim.x) {
55 * make a change to the gf array. 72 if (id[i] == -1) // skip invalid pixels
56 */ 73 continue;
57 gf_changed = false; 74 if (id[i] != thisModuleId) { // find the first pixel in a different module
58 75 atomicMin(&msize, i);
59 76 break;
60 77 }
61 /* 78 }
62 * The algorithm executes in a loop of three distinct 79
63 * stages. In this first one, a mix of stochastic and aggressive 80
64 * hooking, we examine adjacent cells and copy their grand parents 81
65 * cluster ID if it is lower than ours, essentially merging the two 82
66 * together. 83
67 */ 84 }
68 for (index_t tst = 0; tst < MAX_CELLS_PER_THREAD; ++tst) { 85
69 const index_t cid = tst * blkDim + tid; 86

```

- We have a wealth of experience with writing efficient data processing code in a single thread
 - Even our multi-threaded jobs are built from single threaded tasks that do independent, well defined operations
- We have been developing novel algorithms for many years now, that could take advantage of hundreds of thousands of threads
 - But I believe we are still only at the beginning of doing this correctly 🤔

Things To Do Right: ML Inference

- HEP software has used machine learning since a long-long time
 - At the same time, the innovations in industry do present very new opportunities for us
- ML/AI techniques already exist for doing very complicated operations in the trigger / reconstruction software of the experiments
- What is missing are efficient algorithm chains that would “stay on a device”, interleaving ML/AI inference and hand-written algorithmic steps



Next Generation Triggers



Next Generation HEP Triggers Proposal

CERN - European Organization for Nuclear Research

V 1.2_dist (20231212)

Executive Summary

The High Energy Physics (HEP) program at CERN has achieved major breakthroughs in particle physics, technology, and algorithms, including the discovery of the Higgs boson in 2012. This allowed the validation of compatibility of the theoretical construction behind the Standard Model (SM) of particle physics with the data, but the existing uncertainties leave room for models beyond the SM. With the experimental collider framework in place, scientific exploration continues to answer questions around the origin of dark matter, the disproportionately low abundance of antimatter and the nature of the discovered Higgs boson. Hard physics problems aside, much can be gained from improvements to the data acquisition pipeline allowing for capturing a richer set of collision events, furthering scientific understanding.

The Large Hadron Collider (LHC) consists of a 27 km tunnel where superconducting magnets guide bunches of protons, circulating in opposite directions, which are then caused to collide at experimental sites (e.g. ATLAS and CMS) at a rate of 40 million times per second. The collision events emit various particles, which are tracked through a multitude of radiation-hardened detectors and fed into the L1 trigger system, which needs to reject >99% of the events within 10 microseconds due to detector cache constraints and available network capacity.

This data is further reduced by >99% in the High-Level Trigger (HLT) to conform to the current event analysis and simulation capacity. HEP experimentation is fundamentally stochastic, so without changing other factors, an increase in data collection throughput would allow for higher confidence in current results while increasing the likelihood of detecting novel particles in the current LHC setup. Furthermore, this capacity increase is absolutely needed for future LHC upgrades where each collision will have many more interesting events.

The interpretation of the LHC data relies on theoretical simulations of particle interactions in the Standard Model (SM) and in scenarios of new physics beyond the SM (BSM). The full exploitation of the immense HL-LHC datasets, and in perspective of the data from Future Colliders, will require radical improvements in the computing strategies of theory calculations, to increase their accuracy while keeping affordable computing times. A multitude of theoretical tools must be addressed, in a coordinated effort, to preserve their interoperability and harmonize the overall precision. In addition to the several ingredients needed to describe the final states of proton collisions, the infrastructure developed for the triggers, e.g. the GPU cluster, also supports the advancement of software and algorithms for lattice Quantum Field Theory (LQFT) calculations, as a unique approach to control relevant non-perturbative ingredients. The engagement of LQFT experts would also bring to the trigger

- All of the previously highlighted areas are pursued by the “Next Generation Triggers” project
 - With external funding, CERN is starting a 5-year R&D project that will hire a large number of people, and get a large amount of accelerator hardware
- The goal is to find new ways for using the latest types of hardware in processing data from the HL-LHC
- If you’re a graduate/post-doc, and interested in such developments, consider applying! 😊

- HEP experiments worldwide (and especially at CERN) are implementing support for accelerators in their software
 - With industry moving to providing computing power primarily in such a form in the future, this is a must
- CERN (experiments) have a rich development program for the following years to put accelerator usage into production
 - Both on our custom-built clusters, and the Worldwide LHC Computing Grid
- Some of this development is difficult to get help with from non-physicists
 - But with most of it we can certainly take advantage of industry support! Partnerships through CERN OpenLab can be of great help for the HL-LHC data taking / analysis.



<http://home.cern>