# Use of Hardware Accelerators for ATLAS Computing
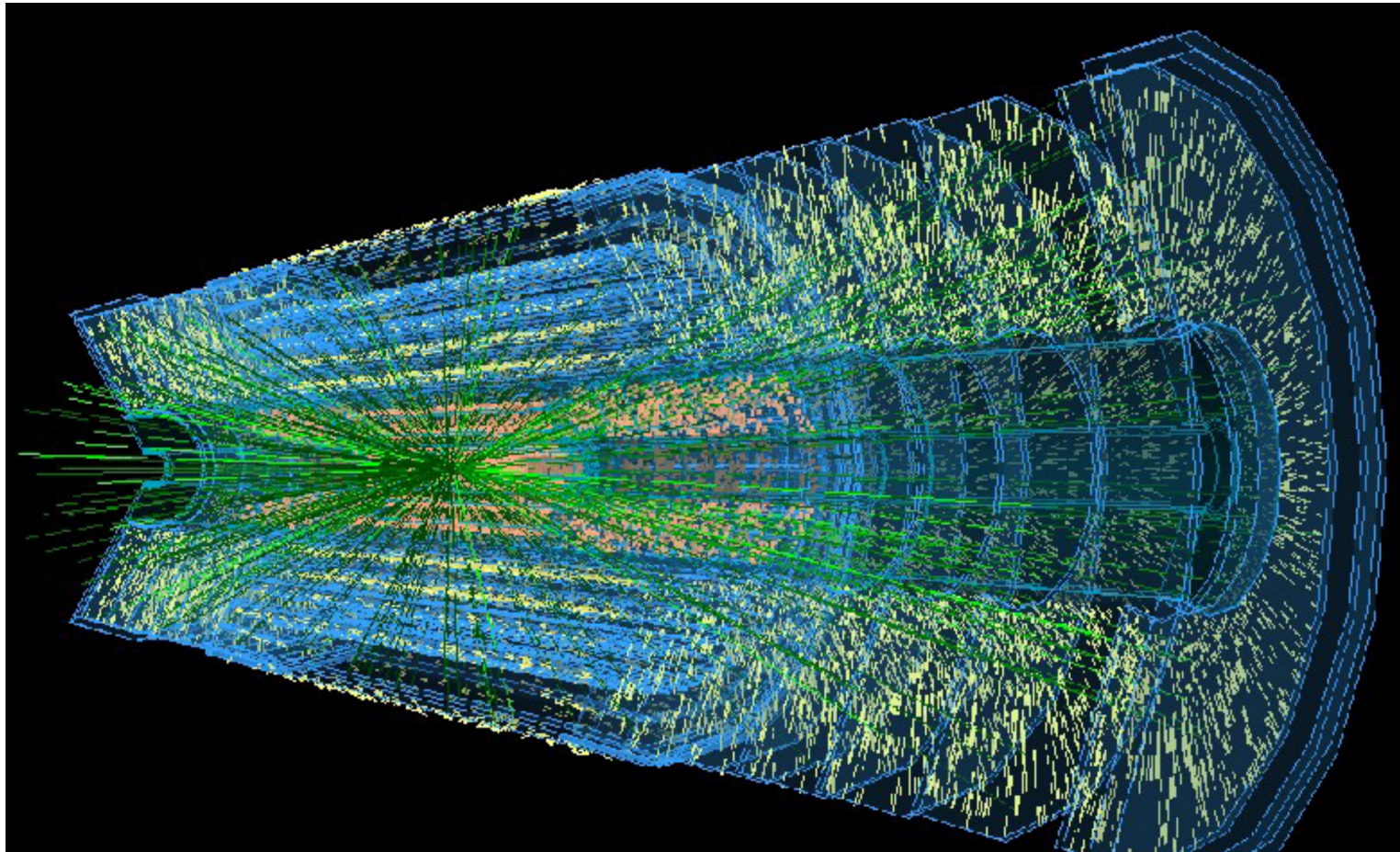
## Maik Dankel
on behalf of ATLAS Collaboration

GPU Computing in High Energy Physics
Pisa, Italy, September 10-12, 2014

# Outline

➤ Motivation

➤ Software Framework integration

➤ Online studies

  – GPU deployment in Muon Trigger algorithms

  – GPU Tracking Algorithms in the ATLAS Trigger
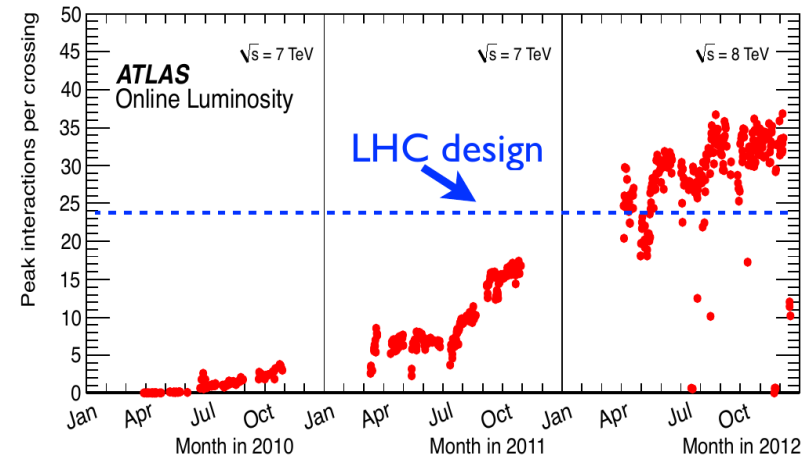
➤ Offline studies

  – GPU based (reference) Kalman-Filter

➤ Summary and Outlook

# Pile-up in Future

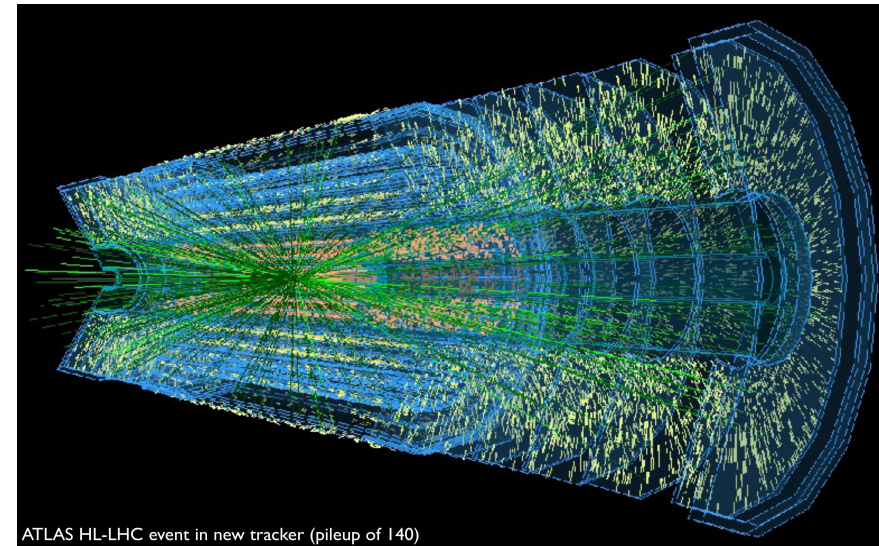- Pile-up (number of instantaneous collisions) exceeded design already in 2012

  - avg. pile-up above 35 interactions per bunch crossing

    → ≈ 1.200 tracks per b.c.



- Expectations for Run-2

  - luminosity up to 2-3 × $10^{34}$ cm$^{-2}$s$^{-1}$

    → pile-up of 40 - 80

  - ATLAS will record events at about 1kHz rate for offline processing

- Run-3 will be even higher



ATLAS HL-LHC event in new tracker (pileup of 140)

# More efficient use of resources

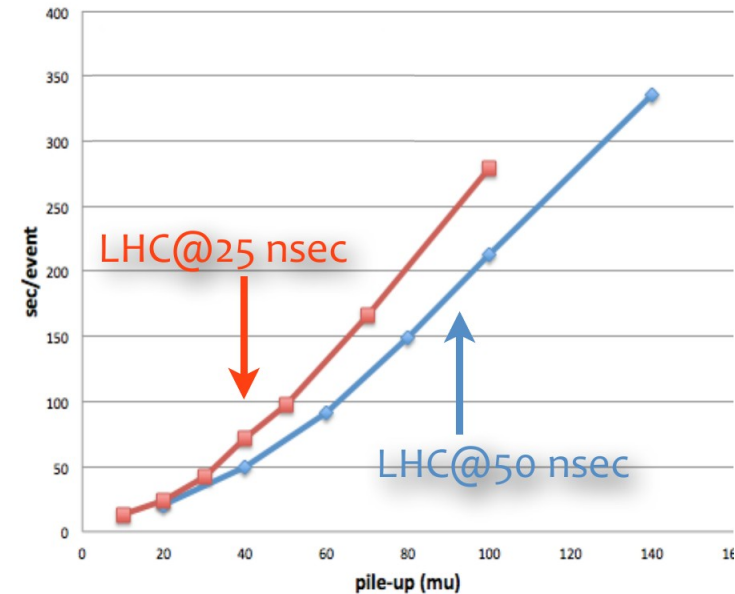> Track finding is a combinatorial problem

- processing time increases highly non-linear with pile-up

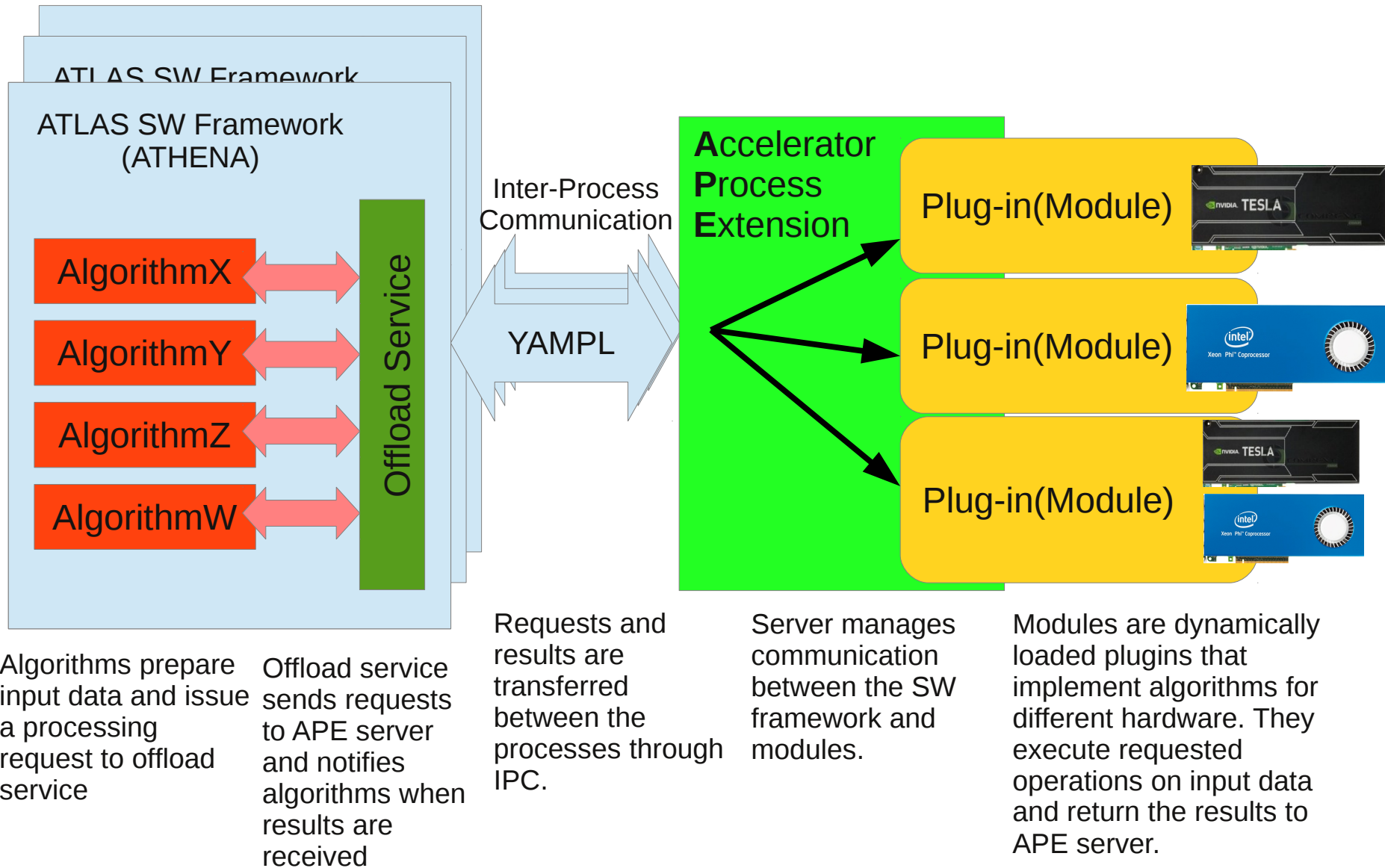> Flat-Budgets → Limited increase in CPU farms

- frequency scaling is halted since advent of manycore chips

- memory and power constraints necessitate parallel processing

- coprocessors promise high compute power at a cheaper price

> The use of GPUs could help to handle upcoming processing challenges

**RAW-> ESD Reconstruction time @ 14 TeV**

LHC@25 nsec

LHC@50 nsec

# Integration to the ATLAS SW framework



**ATLAS SW Framework**

**ATLAS SW Framework (ATHENA)**

- AlgorithmX
- AlgorithmY
- AlgorithmZ
- AlgorithmW

Offload Service

Inter-Process Communication

YAMPL

**A**ccelerator **P**rocess **E**xtension

Plug-in(Module)

Plug-in(Module)

Plug-in(Module)

Algorithms prepare input data and issue a processing request to offload service

Offload service sends requests to APE server and notifies algorithms when results are received

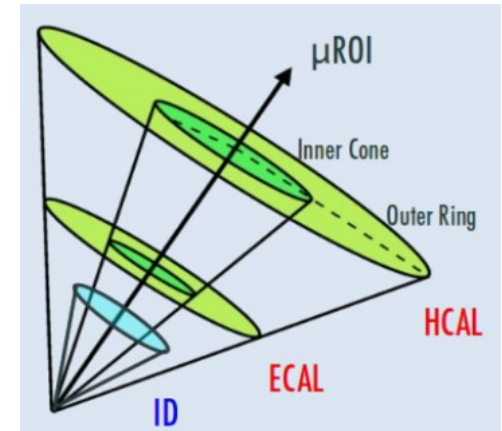Requests and results are transferred between the processes through IPC.

Server manages communication between the SW framework and modules.

Modules are dynamically loaded plugins that implement algorithms for different hardware. They execute requested operations on input data and return the results to APE server.
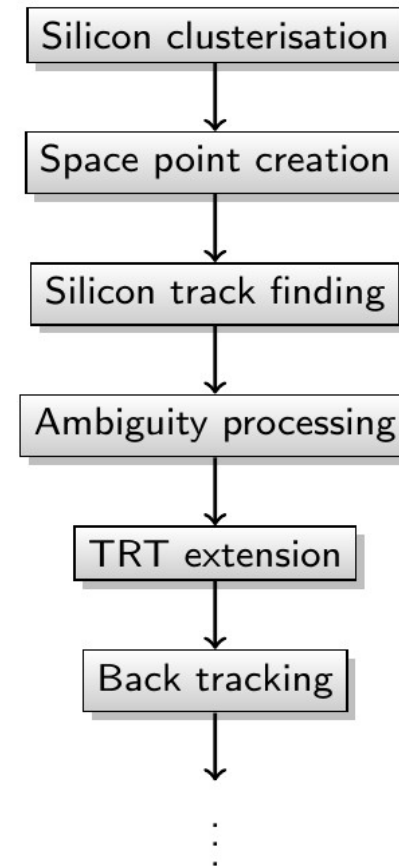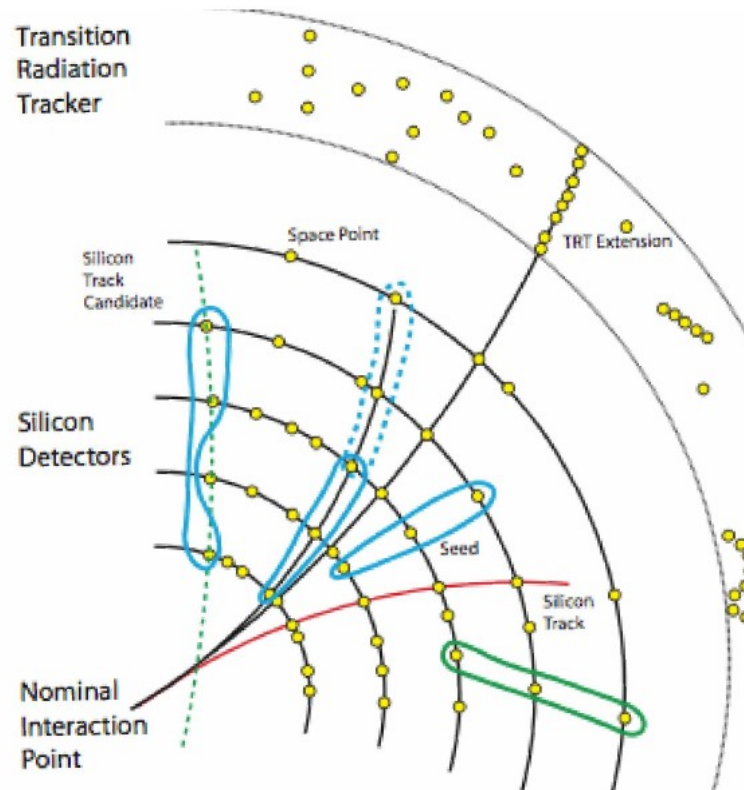
# Integration to the ATLAS SW framework

- APE plugin structure allows different parallel languages
  - Modules can be written in CUDA, OpenCL, OpenMP, …

- APE Server handles multiple ATHENA processes at once

- ATHENA algorithms are indifferent to loaded module

  - Different modules implement algorithms for different hardware

- Yampl abstracts different IPC technologies

  - APE server can run in same host or a dedicated server host

- Standalone clients ease algorithm development and optimization

- Initial tests show negligable overhead

# Using GPUs for Muon Triggers



➤ First Test with simple algorithm for muon isolation

➤ Test GPU - ATLAS framework interaction

➤ Future plans on GPU implementation of a Neural Network for particle identification at Trigger level

➤ See talk from yesterday morning
The GAP Project: GPU applications for high level trigger and medical imaging

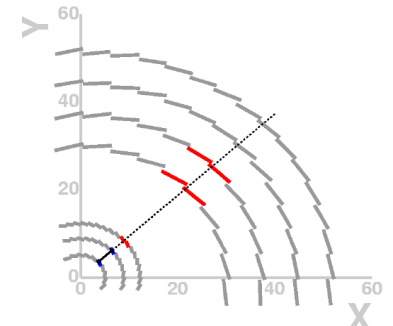11/09/2014

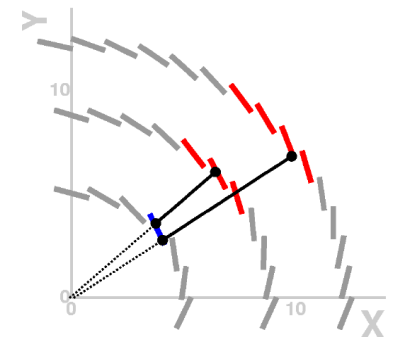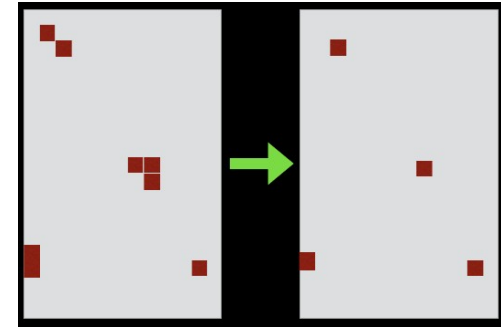# Track reconstruction in ATLAS



Track multiplicity and combinatronic problem makes tracking a natural candidate for parallelization

# Using GPUs for track finding in Triggers

>Containing four main steps:

- Data decoding

  - Decode each data word in parallel on a GPU thread

- Clusterization

  - Use cellular automaton to do parallel clusterization in GPU

- Track Formation

  - Parallel spacepoing creation and seed finding on GPU

- Clone removal
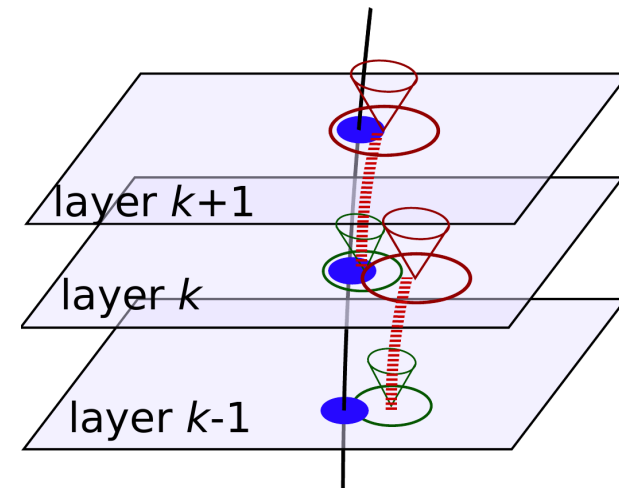
  - Track pair ranges are processed by GPU threads

>See previous talk

An evaluation of the potentials of GPUs to accelerate tracking algorithms for the ATLAS trigger
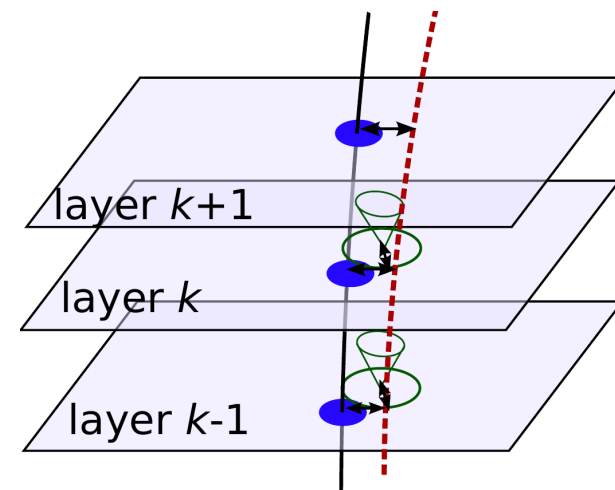
11/09/2014

# GPU-based (Reference) Kalman-Filter

➤ Default Kalman-Filter implementation in ATLAS: Extended KF

– Measurement updates alternate with extrapolation



➤ Reference Kalman-Filter

– Uses a precalculated reference track

– Reference extrapolated through whole volume

– Fitter runs only on differences between measurements and reference trajectory

– More stable in case of outliers

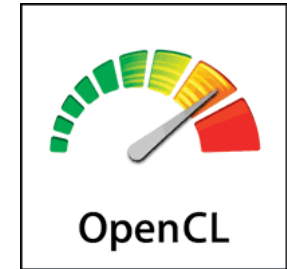# GPU-based (Reference) Kalman-Filter

➤ Standalone (non-ATHENA) code
  – Using .root file as input

➤ Four different implementations
  – OpenCL
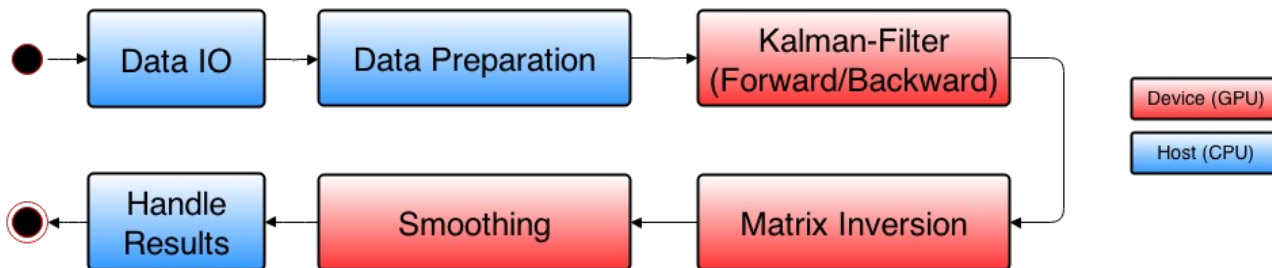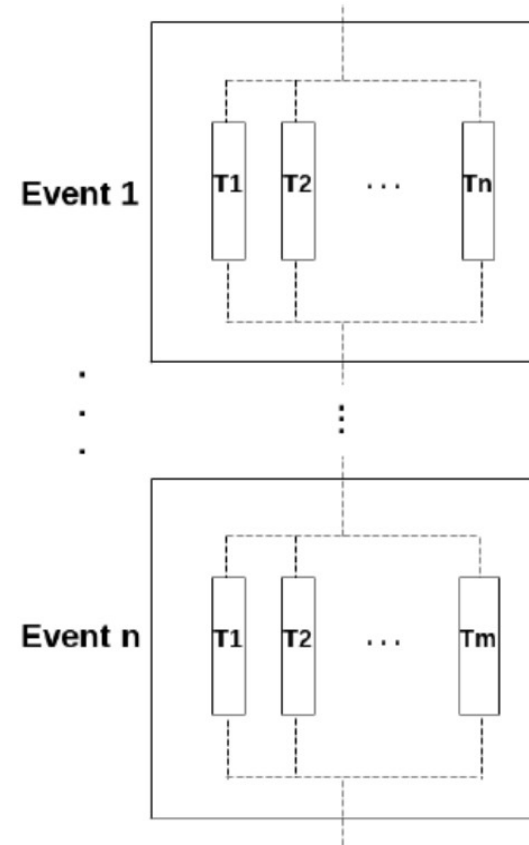  – CUDA
  – OpenMP
  – serial C++

➤ Using the same flat data structures

➤ Producing the same results

# GPU-based (Reference) Kalman-Filter

➤ Whole Kalman-Filter chain is processed on the GPU

– forward/backward filtering

– smoothing

➤ All tracks in an event are processed in parallel

– 5x5 GPU threads per track

➤ Inversion of up to 224 matrices in parallel

# GPU-based (Reference) Kalman-Filter

➤ Runtime of standalone code compared to OpenCL and CUDA on NVIDIA GPU

- GPU Code runs on Nvidia GeForce GTX TITAN
- CPU versions run on Intel Xeon E5-1620 @ 3.60GHz
- OpenMP multithreaded with 8 threads



→ Significant speed-up observed (up to 15× compared to serial)

# Summary and Outlook

➤ Frequency scaling is gone but data rates are increasing

- Need to reduce cost of processing

- Use of hardware accelerators looks promising

➤ Already several encouraging results with significant speedup

- Up to 15/26x compared to serial implementation

➤ Framework integration for short-to-medium term is there

➤ Various ongoing studies in ATLAS

- Evaluating step by step possible parallelizable problems

- Combinatorial nature of track reconstruction seems to be made for parallelization

➤ Porting CPU algorithms to GPU is not an easy task

- still a lot of work to do in the future

# Thank you for your attention

# The ATLAS Experiment

- One of the two biggest general purpose detectors on the LHC

- Cylindrical detectors with endcaps
    - Trackers in center
    - Calorimeters around trackers
    - Muon detectors in outermost shell

- Solenoid and Toroidal magnets for magnetic fields

M. Dankel, Use of Hardware Accelerators for ATLAS Computing, GPU Computing at HEP
11/09/2014

# Inner Detector

▶ Pixel Detector

- 1.744 modules in three concentric layers (4 in next running period)

- 46.080 pixels per module sensitive to traversing charged particles

  → ≈ 80 million readout channels

▶ Semiconductor Tracker (SCT)

- 8.176 modules in eight layers and over 6 million implanted readout strips
  → ≈ 6 million readout channels

- modules are built from doulbe-sided strip sensors with a small relative stereo angle

▶ Transition Radiation Tracker (TRT)

- Straw tube which gives additional information on the traversed particle type

- ≈ 350.000 readout channels