Matteo Bauce

# GPU integration in High Energy Physics experiment online event selection systems

Perspective of GPU computing in Science, 26-28/09/16, Roma

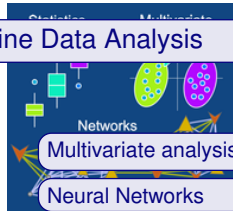Run: 286665
Event: 419161
2015-11-25 11:12:50 CEST

Many applications of GPU in High Energy Physics

Monte Carlo simulation

Lattice QCD calculations

Offline Data Analysis
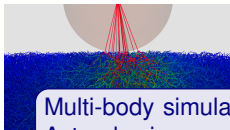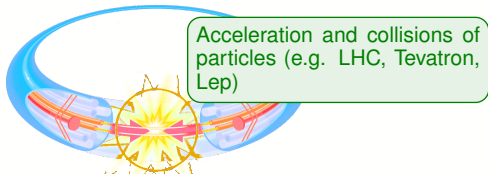
Multivariate analysis

Neural Networks

Higher Levels

Realtime event selection

Lower Levels

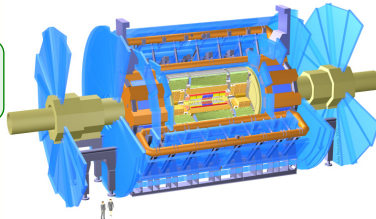Multi-body simulation in Astrophysics

Acceleration and collisions of particles (e.g. LHC, Tevatron, Lep)



proton - (anti)proton cross sections

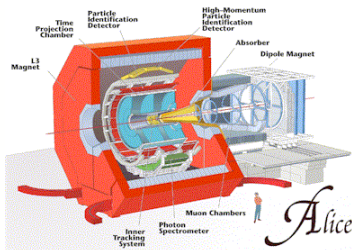Center-of-mass energy up to 13 TeV *pp* collision, rates up to 40 MHz

Multipurpose detector to reconstruct most of the collision information
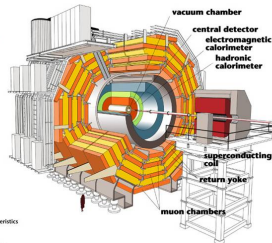
Interesting events are only $1/10^{7-12}$: need to reject most of the others.

**Realtime selection plays a fundamental role**

Event size: 100 MB, 300 Hz: 30 GB/s



Event size: 1-2 MB, 40 MHz: 40 TB/s



Event size: 1.5 MB, 40 MHz: 30 TB/s



Event size: 100 kB, 40 MHz: 4 TB/s

|  | Run 1 | Run 2 | Run 3 |  |
|---|---|---|---|---|
| Energy ($\sqrt{s}$) | 7/8 TeV | 13 TeV | 14 TeV |  |
| Peak Luminosity ($cm^{-2}s^{-1}$) | $10^{34}$ | $1.5 \cdot 10^{34}$ | $2\text{-}3 \cdot 10^{34}$ |  |
| Interactions/bunch crossing | 30 | 23 | 55-80 | ◄ pileup |
| Bunch crossing rate | 20 MHz | 40 MHz | 40 MHz |  |
| Offline Storage rate | 600 Hz | 1000 Hz | 1000 Hz |  |
| Bunch spacing | 50 ns | 25 ns | 25 ns |  |

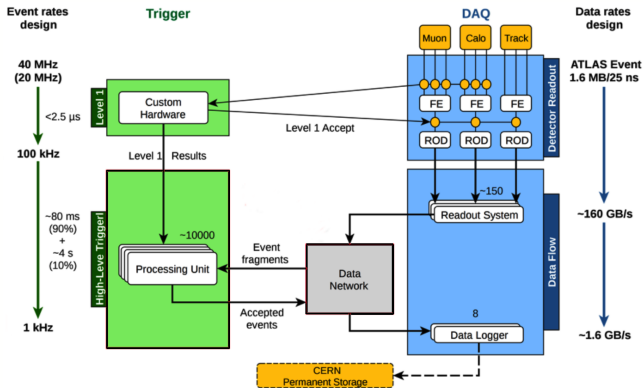Data-taking conditions will be more and more demanding in the upcoming years

- Higher collision rates
- Higer number of multiple overlapping events (pileup)
- Detector upgrades might increase event size
► Processing latencies should remain almost the same $\mathcal{O}(100 \; ms)$

- Multi-stage system based on hardware (LLT) and software (HLT)
- CPU computing power reaching saturation: change of paradigma toward parallel computing

- Multi-stage system based on hardware (LLT) and software (HLT)
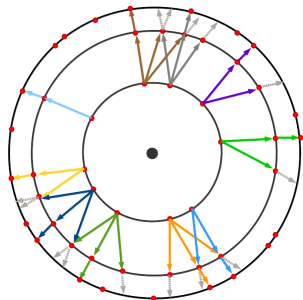- CPU computing power reaching saturation: change of paradigma toward parallel computing
- Try to include GPUs in trigger systems

- *Pattern-recognition* algorithms suitable for parallelization (SIMD)

e.g. Different color $\rightarrow$ Different core

- *Pattern-recognition* algorithms suitable for parallelization (SIMD)
- memory usage is a limitation: small amount available, overhead for data cross-reading algorithm.

e.g. Different color → Different core

- *Pattern-recognition* algorithms suitable for parallelization (SIMD)
- memory usage is a limitation: small amount available, overhead for data cross-reading algorithm.
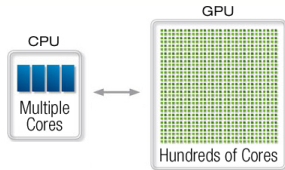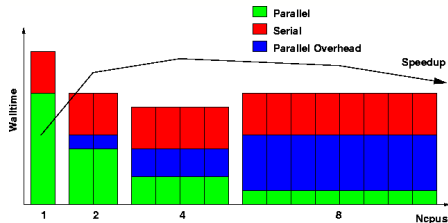- Multi-event parallelization is a BONUS!

e.g. Different color → Different core

► Main questions that need an answer:

1. How to integrate a GPU in a pre-existing data-taking software?
   - ► Need to redesign software from scratch?

2. How fast can a GPU be within time constraints from the DAQ system?
   - ► i.e. how low can you go in the trigger levels?

3. What algorithms get the best from parallelization on GPU?
   - ► Existing ones are suitable for parallelization?
   - ► How innovative ones compare in terms of efficiency?

- Aim at the evaluation of benefit and disadvantages
  - ► Need to suppress increase in CPU time due to pileup
  - ► Limit on HLT farm size from cooling and power
  - ► Evaluate processing time/event per unit cost
- Investigation on trigger algorithm for Inner Detector (tracking), Calorimeter clustering, Muon segment reconstruction

Server with NVidia Tesla K80

- 2 chips in each card
- 2 GB RAM
- 13 multi-processor
- 192 cores per multiprocessor

- 2496 CUDA cores
- 824 MHz GPU, 2505 MHz memory clock

# Flexible client-server architecture

**Client side** | **Server side**



**Client:**

- One HLT processing unit per core
- Offline & Online framework (Athena)
  - ▶ manage data
  - ▶ execute chains of algorithms
  - ▶ monitors data-processing

**Server:**

- Independent from Client framework
- Flexible hardware resources management (multi-devices)
- Preallocate memory for data and store constants

Tracking is the most time consuming algorithm

Bytestream decoding → Hit clustering → Track seeding → Track Following → Clone Removal

- Sequential steps: silicon hit clustering, seeds creation, track following

Tracking is the most time consuming algorithm



- Sequential steps: silicon hit clustering, seeds creation, track following
- Parallelization on GPU of **track-seeding**
- Huge data multiplicity for a full-detector scan tracking: a GPU makes it feasible



$10^5$ spacepoints → $10^9$ triplets → $10^4$ seeds → $10^3$ tracks

**Pair formation**: 2D thread array checking for pairing conditions

**Triplet formation** through 2D thread block

Tracking is the most time consuming algorithm



Track following also may benefit

- Sequential steps: silicon hit clustering, seeds creation, track following
- Parallelization on GPU of **track-seeding**
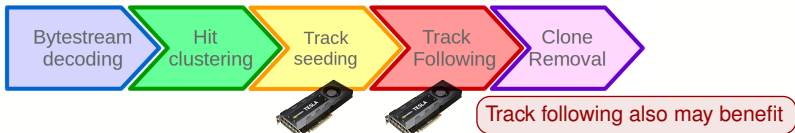- Huge data multiplicity for a full-detector scan tracking: a GPU makes it feasible



$10^5$ spacepoints    $10^9$ triplets    $10^4$ seeds    $10^3$ tracks

**Pair formation**: 2D thread array checking for pairing conditions

**Triplet formation** through 2D thread block

GPU algorithm has same efficiency and resolution as CPU one

- Algorithm execution time reduced by a factor $\sim$5
- Small data transfer overhead: $\sim$0.6%

Algorithm needed to associate energy deposits from the same shower of particles (hadronization)
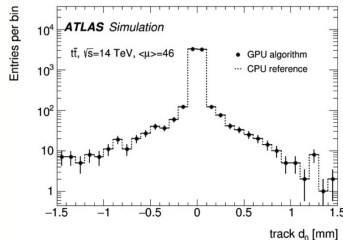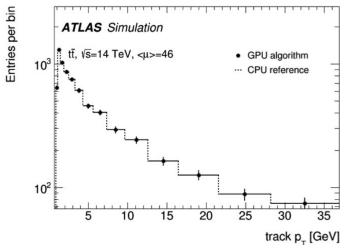
- Topological Calorimeter Cell Clusters reconstruction on CPU: $\sim 8\%$ of total time
  - Cells are grouped according to their signal-to-noise ratio



| | |
|---|---|
| Seed | |
| Growing | |
| Terminal | |
| Not enough S/N | |
| Not evaluated | |

- Topo-Automaton Clustering on GPU to maximize parallelism
  - Propagation of a flag through a grid of elements (cell pairs)
  - Cells get the largest flag and continue until no flag changes

# Calorimeter



*ATLAS* Simulation
Monte Carlo 14 TeV di-jet events <μ>=40

— CPU clusters $E_T$
  <$E_T$> = 0.907 GeV
- - GPU clusters $E_T$
  <$E_T$> = 0.912 GeV

◄ Energy difference <5% in most clusters

► no significant effect on jet reconstruction



*ATLAS* Simulation
Monte Carlo 14 TeV di-jet events <μ>=40

— CPU number of jets
  <number Jets> = 32
- - GPU number of jets
  <number Jets> = 32

- 30% reduction for di-jet events with 40 interactions/bunch-crossing($\mu$), $\times 3$ reduction for $t\bar{t}$ with $\mu$=138
- Data-format conversion reduce the benefits
- Potential larger gain from parallelization of following clusterization steps



Calorimeter Clustering on CPU

Athena 91.8%
CPU Clustering 84ms 8.2%

Time per event 1.02 s



Calorimeter Clustering on GPU

Athena 95.9%
GPU Clustering 44ms 4.1%

Conversion CPU→GPU 0.2%
Conversion GPU→CPU 1.0%
Data Transfer 0.1%
IPC 0.7%
Classification 0.05%
Growing 1.8%
Tagging 0.3%

Time per event 1.06 s

▶ Muon segment reconstruction through Hough Transform

- algorithm translates track finding to maxima finding
- Filter hits and fill Hough parameter space
- Select maxima above a given threshold and reconstruct track parameters



Development ongoing - public results expected soon

► Testing E5-2695 v3 14-core vs. 1(/2) NVidia K80 GPU

- 20-40% gain in throughput, depending on the number of processes running
- 1 GPU saturation when serving 14 clients (no performance loss)
- Slight benefit from the additional GPU

1s → 300 ms reduction by deploying GPU in TPC tracking system

Considering improvements and modification to the trigger scheme for the experiment upgrades.

more info in talk from D. Rohr

**40 MHz bunch crossing rate**

**L0 Hardware Trigger : 1 MHz readout, high $E_T/P_T$ signatures**

| 450 kHz $h^\pm$ | 400 kHz $\mu/\mu\mu$ | 150 kHz $e/\gamma$ |

**Software High Level Trigger**

Partial event reconstruction, select displaced tracks/vertices and dimuons

Buffer events to disk, perform online detector calibration and alignment

Full offline-like event selection, mixture of inclusive and exclusive triggers

**5 kHZ Rate to storage**

- For HL-LHC aim at a triggerless scheme (no hardware)
- GPU deployment can boost software trigger level
- Evaluation in progress to minimize communication latencies and throughput

**30 MHz inelastic event rate and full event rate building**

**LLT : 15-30 MHz output rate, select high $E_T/P_T$ ($h^\pm/\mu/e/\gamma$)**

**Software High Level Trigger**

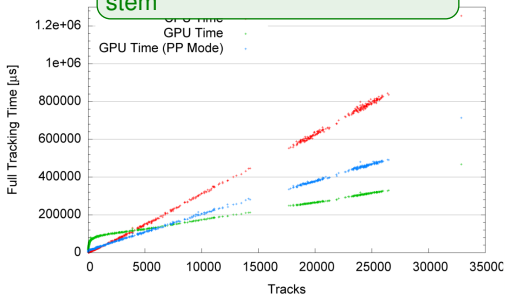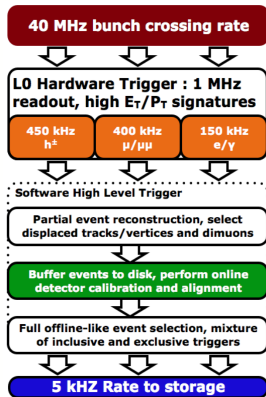Full event reconstruction, inclusive and exclusive kinematic/geometric selections

Run-by-run detector calibration

Add offline precision particle identification and track quality information to selections

**2-10 GB/s rate to storage**

Focusing on vertex reconstruction and tracking algorithm ▶

- VErtex LOcator detector fundamental for displaced vertices detection
- Tracking in muon detectors

only ∼50 $\mu$s overhead ▶

Mirror Mosaic (17 m focal length)

17 m

Beam

2024 TDC channels, 4 TEL62

TEL62

GPU

**Vessel diameter 4→3.4 m**
Volume ~ 200 m³

Beam Pipe

2 × ~1000 PM

External Power

Gbit Ethernet

SO-DIMM DDR3

Programmable Device

PCI-e connector

Gbit Ethernet
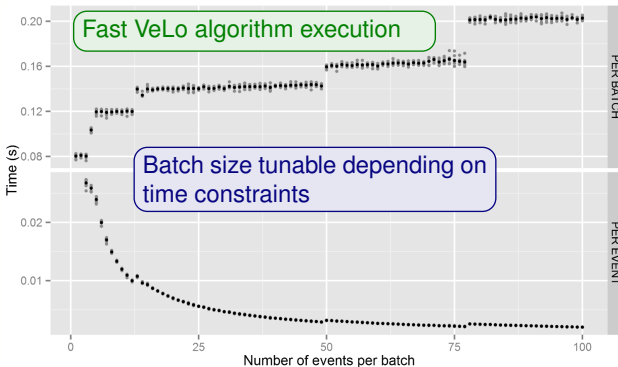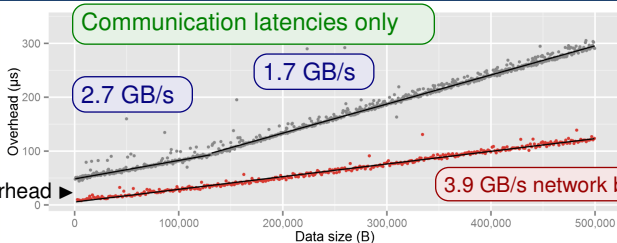
mini-USB

QSFP+ Connectors

Improved RDMA scheme allowed to increase throughput and deploy GPU in **low-level trigger** (smaller event size)



- NaNet-10 moves Data to CPU memory
- NaNet-10 moves Data to GPU memory (GPUDirect v2)
- NaNet-10 moves Data to GPU memory (GPUDirect RDMA)
- NaNet-1 moves Data to GPU memory (GPUDirect v2)

Bandwidth (MB/s) axis: 0, 200, 400, 600, 800, 1000, 1200, 1400

Message size (Byte) axis: 16, 32, 64, 128, 256, 512, 1K, 2K, 4K, 8K

more info in talk from A. Biagioni

▶ Parallelism in realtime selection system is a must: GPUs deployment is crucial

- GPU integration can be achieved in a transparent way
    - ▶ Client-Server architecture: the most flexible solution for DAQ existing frameworks
    - ▶ New experiment can deploy different scheme, no constraints
    - ▶ Careful design of EDM needed

- Communication overhead latencies define the feasibility domains

    - ▶ High-level trigger applications accessible for typical HEP experiment sizes (latencies $\mathcal{O}(100 \ \mu s - 100 ms)$)
    - ▶ From the detector to the GPU in Low-level trigger application, achieved thanks to dedicated interface cards

- Algorithm optimization can add gain in parallelization
    - ▶ Several developed for pattern recognition algorithms (Hough Transform, Cellular Automaton, ... )
    - ▶ Neural Networks (and MVA) might come into the game in the future