

# CERN openlab III

Short GPU/multi-  
core “slide-show”  
from CHEP-2010



Sverre Jarp, CERN IT  
2 November 2010

- Plenary:
  - The experiment offline systems after one year (R.Jones)
  - How to harness the performance potential of current **Multi-Core CPUs and GPUs** (S.Jarp)
  - Computing **Paths to the Future** (R.Goff/DELL)
- Parallel:
  - **Multicore-aware** Applications in CMS
  - Parallelizing Atlas Reconstruction and Simulation: Issues and Optimization Solutions for Scaling on Multi- and **Many-CPU** Platforms
  - **Multi-threaded** Event Reconstruction with JANA
  - Track Finding in a High-Rate Time Projection Chamber Using **GPUs**
  - **Fast Parallel** Tracking Algorithm for the Muon System and Transition Radiation Detector of the CBM Experiment at FAIR
  - Real Time Pixel Data Reduction with GPUs And Other HEP GPU Applications
  - Algorithm Acceleration from **GPGPUs** for the ATLAS Upgrade
  - Maximum Likelihood Fits on **Graphics Processing Units**
  - Partial Wave Analysis **on Graphics Processing Units**
  - **Many-Core** Scalability of the Online Event Reconstruction in the CBM Experiment
- BOF 3:
  - **GPUs**: High Performance Co-Processors



# Plenary on Monday by Roger Jones(ATLAS)

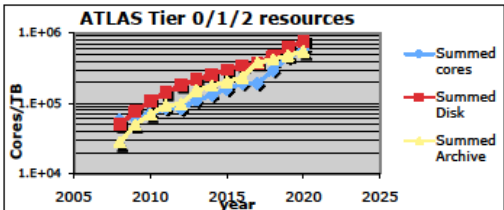
CERN  
openlab

■ “The experiment offline systems after one year”

LANCASTER UNIVERSITY 10/18/10 RWL Jones CHEP2010

## Future Challenges


- We assume we can use growth in CPU
  - But this implies changing architectures
  - And handle the data throughput



**ATLAS Tier 0/1/2 resources**

Year	Summed cores (1.E+06)	Summed Disk (1.E+05)	Summed Archive (1.E+04)
2005	~1.5E+05	~1.5E+04	~1.5E+04
2010	~3.0E+05	~3.0E+04	~3.0E+04
2015	~6.0E+05	~6.0E+04	~6.0E+04
2020	~1.2E+06	~1.2E+05	~1.2E+05
2025	~2.4E+06	~2.4E+05	~2.4E+05



- Experiments already working to deal with multi cores
  - Many cores and GPGUs are down the line
- We need to use them or be very clear why we cannot



LANCASTER UNIVERSITY 10/18/10 RWL Jones CHEP2010 26

## Related developments

- IO challenges being (partly) addressed by fast merging
- Re-write of Gaudi with stronger memory model planned
- Down the line, we may need to parallize the code
  - This could be either for many-core processors or for Graphical Processing Units – but the development might address both
    - GPUs having big success & cost savings in other fields
    - Harder for us to use, but funders will continue to ask
    - We need the R&D to know which path to take
  - Developments require O(3 years) to implement
  - This includes Geant4 – architectural review this year



CERN  
openlab

# Plenary on Monday by S.Jarp

- “How to harness the performance potential of current Multi-Core CPUs and GPUs”

Today:

Seven dimensions of multiplicative performance

## ■ First three dimensions:

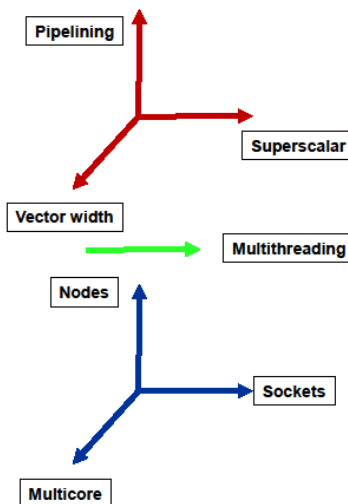
- Pipelined execution units
- Large superscalar design
- Wide vector width (SIMD)

## ■ Next dimension is a “pseudo” dimension:

- Hardware multithreading

## ■ Last three dimensions:

- Multiple cores
- Multiple sockets
- Multiple compute nodes



SIMD = Single Instruction Multiple Data

Sverre Jarp - CERN



## What are the multi-core options?

- There is a discussion in the community about the best way(s) forward:

- 1) Stay with event-level parallelism (and entirely independent processes)
  - Assume that the necessary memory remains affordable
  - Or rely on tools, such as KSM, to help share pages
- 2) Rely on forking:
  - Start the first process; Run through the first “event”
  - Fork N other processes
  - Rely on the OS to do “copy on write”, in case pages are modified
- 3) Move to a fully multi-threaded paradigm
  - Still using coarse-grained (event-level) parallelism
    - But, watch out for increased complexity

17

Sverre Jarp - CERN



- “How to harness the performance potential of current Multi-Core CPUs and GPUs”

Today:

Seven dimensions of multiplicative performance

■ **First three dimensions:**

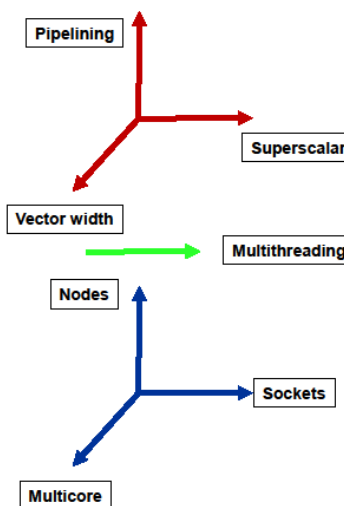
- Pipelined execution units
- Large superscalar design
- Wide vector width (SIMD)

■ **Next dimension is a “pseudo” dimension:**

- Hardware multithreading

■ **Last three dimensions:**

- Multiple cores
- Multiple sockets
- Multiple compute nodes



SIMD = Single Instruction Multiple Data

Sverre Jarp - CERN



## Shortlist

- 1) Broad Programming Talent
- 2) Holistic View with a clear split:  
Prepare to compute – Compute
- 3) Controlled Memory Usage
- 4) C++ for Performance
- 5) Best-of-breed Tools



# Plenary on Tuesday by Roger Goff/DELL

**CERN**  
openlab


- Plenary “vendor” session:

## Computing Paths to the Future



Roger Goff  
Dell Global CERN/LHC Technologist  
+1 970 672 1252 | [Roger\\_Goff@dell.com](mailto:Roger_Goff@dell.com)

## Final Takeaways

1. CPU cores are not getting faster.
2. Co-processors are here to stay.
3. Heterogeneous processors are inevitable.
4. Preparing applications for extreme parallelism will enable users to get the most out of future systems. 

8

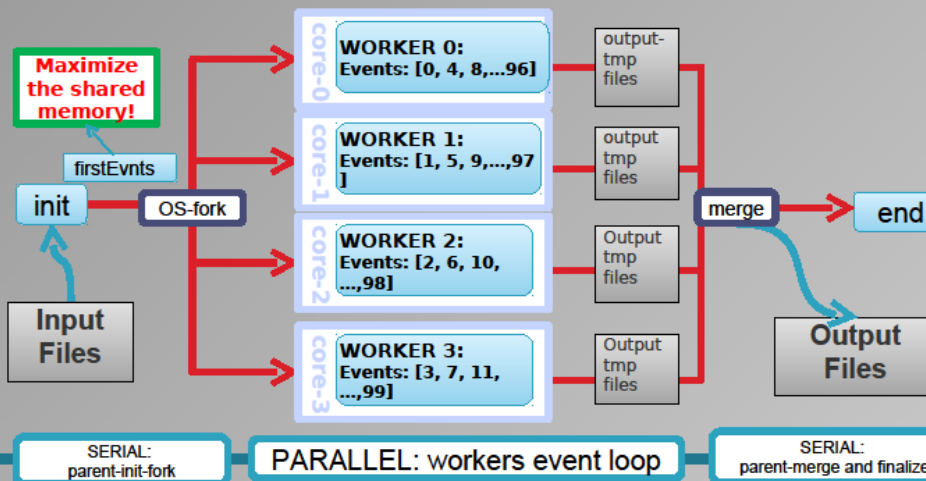
Dell CERN/LHC Program



- “Parallelizing Atlas reconstruction and simulation on multi-core platforms”

## Event Level Parallelism with AthenaMP

```
> Athena.py --nprocs=4 -c EvtMax=100 Jobo.py
```

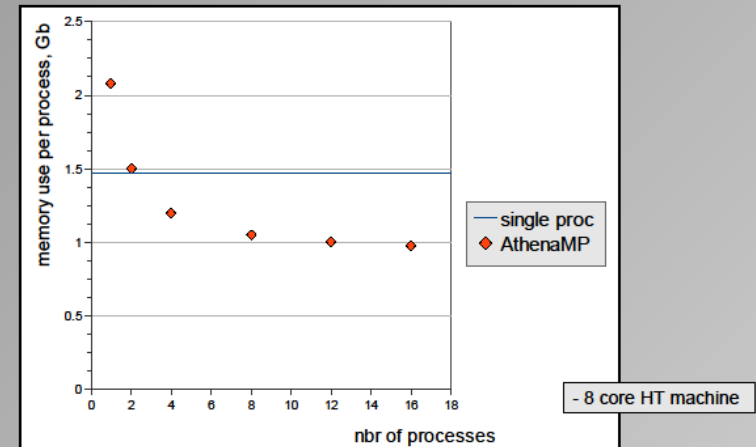


5/23

CHEP 2010

## Memory Usage

- Single reco process uses ~1.5 Gb
- We can save significant memory through OS level sharing



**AthenaMP ~0.5 Gb physical memory saved per process**

6/23

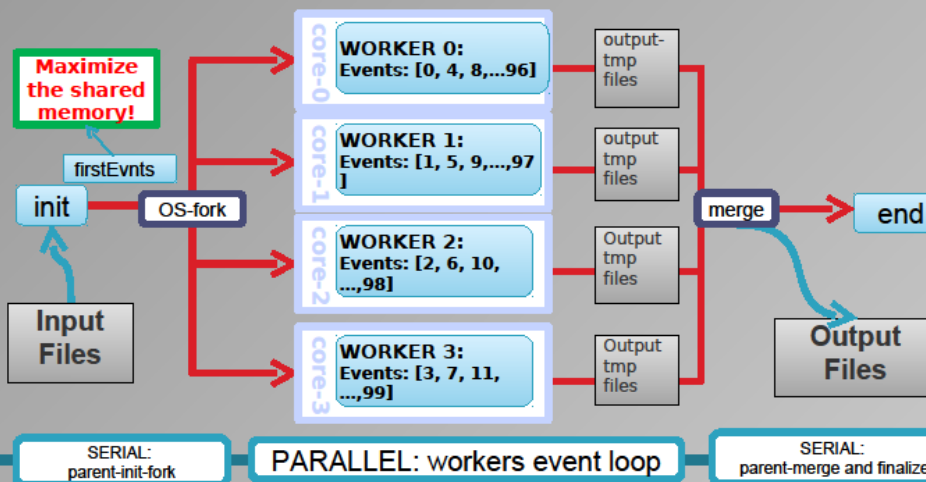
CHEP 2010

10/19/10

- “Parallelizing Atlas reconstruction and simulation on multi-core platforms”

## Event Level Parallelism with AthenaMP

```
> Athena.py --nprocs=4 -c EvtMax=100 Jobo.py
```



5/23

CHEP 2010

## Issues with Large OOP Code Bases

- Function calls result in added instructions
  - Call and return
  - Runtime address resolution (trampolines) required for position independent code/ shared object cross invocations
    - Indirect branches can be more costly
  - Freeing & restoring registers for local use
  - Setting and reading function arguments
- Virtual function calls (function pointers) increase indirect call instructions and associated pointer loads
  - Virtual functions can't be inlined!
- Atlas code has 2500 shared libraries!

19/23

CHEP 2010

10/19/10





# “Multi-core aware Applications in CMS”

## Why Bother?



HEP processing is naturally parallelizable

We have billions of events

Each event can be processed independently

Memory is becoming a limitation

Historically GB/US\$ increases at the same rate as number of transistors in a CPU

<http://www.icmit.com/memoryprice.htm>

Funding levels are not guaranteed to stay this high

We can afford 2GB/core now but may not in the future

Opportunistic use of grid sites improves if we lower our memory requirements

Not all grid sites have 2GB/core

Technical limitations on connecting many cores to shared system memory

<http://www.intel.com/technology/itj/2007/v11i3/3-bandwidth/7-conclusion.htm>

Multi-core aware applications can improve memory sharing

Threading

All threads share the same address space but have to worry about concurrent usage

Forking

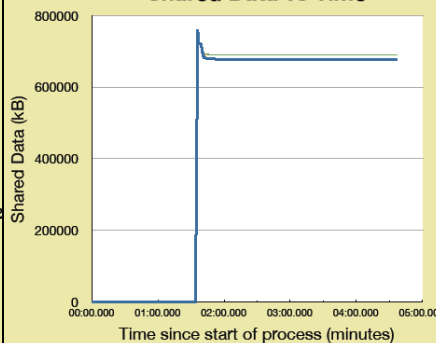
Each child process gets its own address space

Untouched memory setup by the parent is shared between the child processes

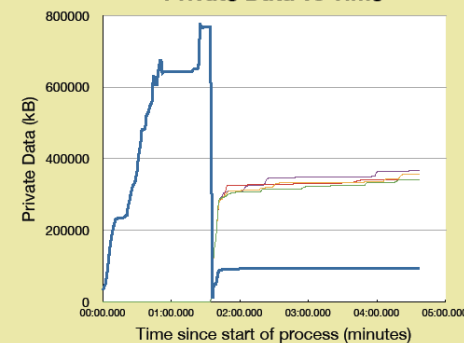
## Memory Sharing



Shared Data vs Time



Private Data vs Time



Measurements done using reconstruction with 64bit software on 4 CPU, 8 core/CPU 2GHz AMD Opteron(tm) Processor 6128

Shared memory per child: ~700MB

Private memory per child: ~375MB

Total memory used by 32 children: 13GB

Total memory used by 32 separate jobs: 34 GB

Saved 62% of memory





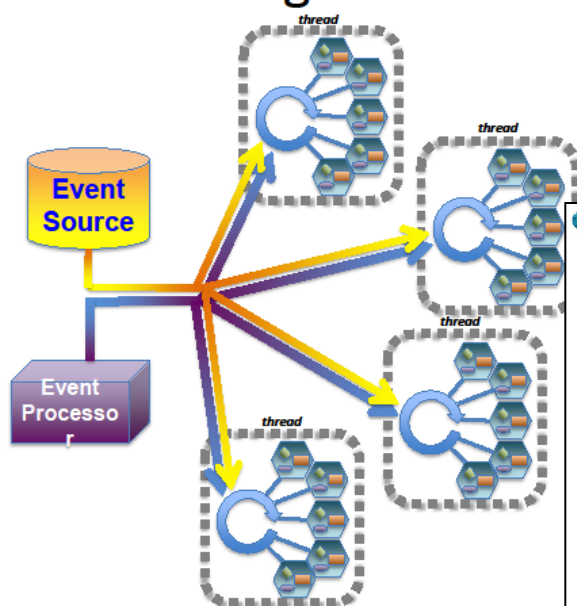
## “Multithreaded event reconstruction with JANA”

### Multi-threading

Each thread has a complete set of factories making it capable of completely reconstructing a single event

Factories only work with other factories in the same thread eliminating the need for expensive mutex locking within the factories

All events are seen by all Event Processors (multiple processors can exist in a program)

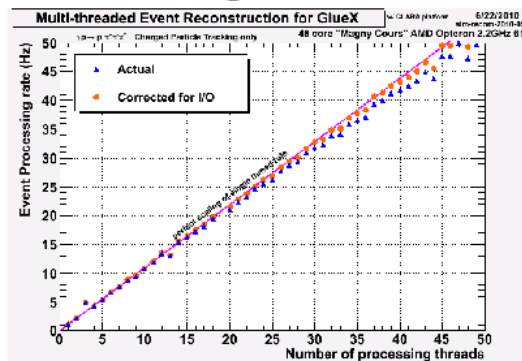


5/25/10

JANA - Lawrence - CLAS12 Software Workshop

6

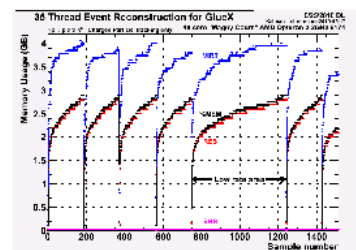
### Testing on a 48-core “Magny Cours”



Event reconstruction using 48 processing threads on a CPU with 48 cores generally scales quite well.

Eventually, an I/O limit will be encountered

Memory Usage vs. time while repeatedly running the 35 thread test. The marked area indicates one test where the program ran slower.



- Occasionally some problems with inexplicably lower rates.
- Program appears to simply run slower while not operating any differently.
- Unclear if this is due to hardware or Linux kernel

Oct. 19, 2010

CHEP10, Taipei – David Lawrence, JLab

9

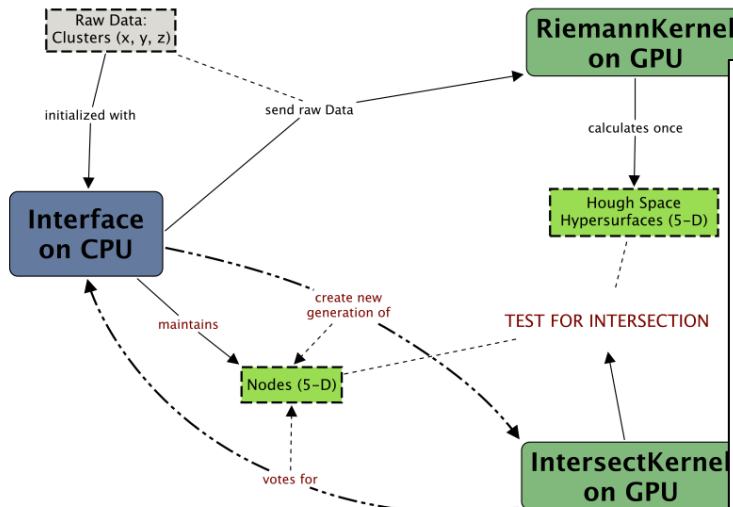
## “Track finding in a high-rate time projection Chamber using GPUs”

A TPC for PANDA 5D Hough Transform for Helical Track Detection Application in our Testbeam at CERN



Implementation on a GPU using CUDA™

TUM  
Technische Universität München



A TPC for PANDA 5D Hough Transform for Helical Track Detection Application in our Testbeam at CERN

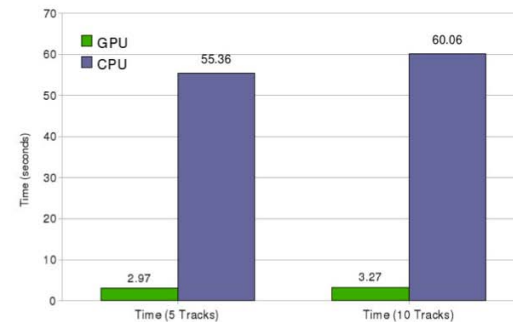


Glance at Performance

TUM  
Technische Universität München

### Hardware setup:

- CPU: Intel Core™2 Quad Q8400 2.66 GHz (single thread)
- GPU: NVIDIA GTX 280 (1GB)



18.4x

Execution time comparison

## ■ “Partial wave analysis at BES III – Harnessing the power of GPUs”

### Parallel PWA on GPU

Events are independent - calculate terms in the sum in parallel

- Use many PCs
- Use parallel hardware and make use of Single Instruction - Multiple Data (SIMD) capabilities
- Very strong here: Graphics processors (GPUs): Cheap and powerful hardware

PWA is embarrassingly parallel:

- Exactly the same (relatively simple) calculation for each event
- Every event has its own data, or fit parameters are shared
- Ideal for GPU implementation
- True for many HEP applications



13

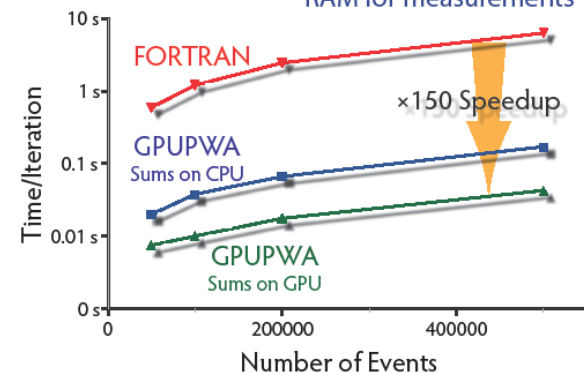
Analysis on GPUs — Niklaus Berger

CHEP 2010, Taipei

### Performance

We use a toy model  $J/\psi \rightarrow \gamma K^+ K^-$  analysis for all performance studies

Using an Intel Core 2 Quad 2.4 GHz workstation with 2 GB of RAM and an ATI Radeon 4870 GPU with 512 MB of RAM for measurements



18

Partial Wave Analysis on GPUs — Niklaus Berger

CHEP 2010, Taipei

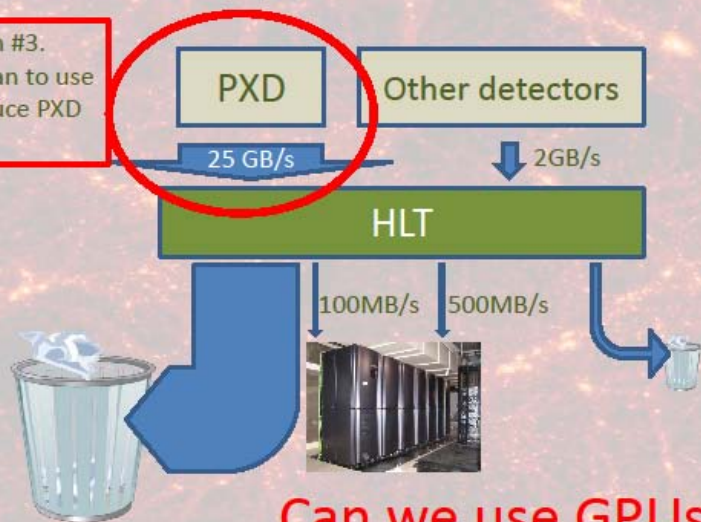




- “HLT (and other things) with GPUs”

## Belle II High Level Trigger

This is option #3.  
There is a plan to use  
FPGA to reduce PXD  
hits



Can we use GPUs?

Sorry. We have just started to work on GPU here and I don't have much to report yet  
You've heard about Belle HLT in the morning

## CPU/GPU Results

	6144×6144 (GFLOPS)	12288×12288 (GFLOPS)
merical recipes)	0.61	0.62
3GHz) MKL(6core)	68	72
(code)	108	115
code)	87	91
	159	190

>×2 faster

ce of C2050 is supposed to be >500 GFLOPS  
d is suppressed to ¼ of C2050 (but is faster)  
ULA 2.1 are used

- Floats are faster as expected
- Copy from/to CPU to/from GPU not included (but is significant)
- Cholesky decomposition cannot saturate GPU
- Cannot do 24576×24576(not enough memory on one C2050)



## “HLT (and other things) with GPUs”

### Inflation

- Inflation started  $10^{-36}$  sec. after the birth of the universe and lasted for  $10^{-34}$  sec.
- During that period, the universe expanded by an order of  $e^{60}$ , from Planck scale to a meter (our observable universe)
- Inflation was caused by a particle (field) of energy scale of  $10^{16}$  GeV
- Cosmic Microwave Background Radiation is the probe to measure its energy scale
  - Let me use the parameter “r” to represent

2010/10/21

Nobu Katayama

### CPU/GPU Results

	6144×6144 (GFLOPS)	12288×12288 (GFLOPS)
CPU (i-7 920) (numerical recipes)	0.61	0.62
CPU(i-7 <a href="#">X980@3.33GHz</a> ) MKL(6core)	68	72
GTX480 (our CUDA code)	108	115
C2050 (our CUDA code)	87	91
C2050 (CULA)	159	190



- Peak performance of C2050 is supposed to be >500 GFLOPS
- GTX480 DP speed is suppressed to ¼ of C2050 (but is faster)
- CUDA 3.1 and CULA 2.1 are used
- Floats are faster as expected
- Copy from/to CPU to/from GPU not included (but is significant)
- Cholesky decomposition cannot saturate GPU
- Cannot do 24576×24576(not enough memory on one C2050)

16



## ■ “Algorithm acceleration from GPGPUs for the ATLAS upgrade”

GPU Computing  
Z Finder  
Kalman Filter  
Summary

GPGPUs  
GPU Projects at Edinburgh  
Project Resources  
ATLAS Trigger

### GPU Projects at Edinburgh

- Number of GPU related projects at Edinburgh over the summer:
- **Chris Jones** - "Porting the Z finder algorithm to GPU" (MSc in High Performance Computing)
- **Maria Rovatsou** - "SIMT design of the High Level Trigger Kalman Fitter" (MSc School of Informatics)
- **James Henderson** - "An Investigation Into Particles Tracking and Simulation Algorithms using GPUs"
- Project reports and source code available at: [ATLAS Edinburgh GPU Computing](#)

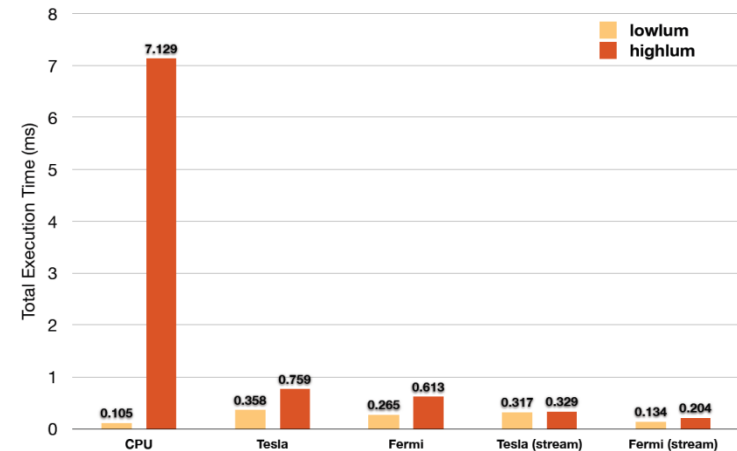
Navigation icons

3 / 22

GPU Computing  
Z Finder  
Kalman Filter  
Summary

GPU Motivation  
Algorithm and Test case  
Z Finder Kernel  
Timing Results

### Timing Results



- Results for spacepoint pairs show up to 35x speed-up (Fermi).
- Initial results for spacepoint triplets also show speed-up.

Navigation icons

12 / 22





# Alfio Lazzaro/CERN openlab

CERN  
openlab

- “Maximum likelihood fits using GPUs”



## Test environment

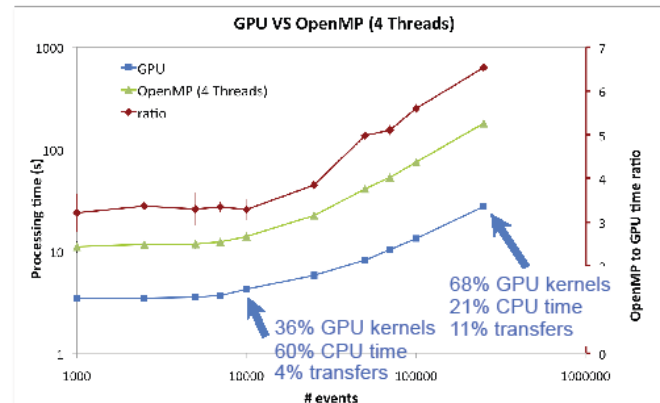
- PCs
  - CPU: Nehalem @ 3.2GHz: 4 cores – 8 hw-threads
  - OS: SLC5 64bit - GCC 4.3.4
  - ROOT trunk (October 11<sup>th</sup>, 2010)
- GPU: ASUS nVidia GTX470 PCI-e 2.0
  - Commodity card (for gamers)
  - Architecture: GF100 (Fermi)
  - Memory: 1280MB DDR5
  - Core/Memory Clock: 607MHz/837MHz
  - Maximum # of Threads per Block: 1024
  - Number of SMs: 14
  - CUDA Toolkit 3.1 06/2010
  - Developer Driver 256.40
  - Power Consumption 200W
  - Price ~\$340

Alfio Lazzaro (alfio.lazzaro@cern.ch)



## PDF-event-base: GPU VS OpenMP

- Fair comparison
  - Same algorithm
  - Algorithm on CPU optimized and parallelized (4 threads)
  - CPU does the final sum of the *NLL* and normalization integral calculations
- Check that the results are compatible: asymmetry less than  $10^{-12}$



- Speed-up increases with the dimension of the sample, taking benefit from the data streaming on GPU and the integral calculation only on the CPU
- ~3x for small samples, up to ~7x for large samples

Alfio Lazzaro (alfio.lazzaro@cern.ch)

15



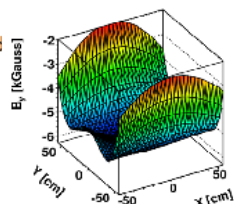


## “Fast track reconstruction of the muon system and transition radiation detector”

### Optimization of the algorithm

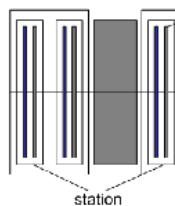
- Minimize access to global memory

- o Approximation of the 70 MB large magnetic field map
  - 7 degree polynomial in the detector planes was proven the best



- Simplification of the detector geometry

- o Problem
  - Monte-Carlo geometry consists of 800000 nodes
  - Geometry navigation based on ROOT TGeo
  - Take into account absorbers and staggered Z positions of the stations
- o Solution
  - Create simplified geometry by converting Monte-Carlo geometry
  - Implement fast geometry navigation for the simplified geometry



- Computational optimization of the Kalman Filter

- o From **double** to **float**
- o Implicit calculation on non-trivial matrix elements
- o Loop unrolling
- o Branches (if then else ..) have been eliminated

**All these steps are necessary to implement SIMD tracking**

A. Lebedev, "Track reconstruction in the muon system and transition radiation detector of the CBM experiment"

### Performance of the track fit in MUCH

#### Track fit quality

Residuals					Pulls				
X [cm]	Y [cm]	Tx *10 <sup>-3</sup>	Ty *10 <sup>-3</sup>	q/p *10 <sup>-3</sup> [GeV <sup>-1</sup> ]	X	Y	Tx	Ty	q/p
0.38	0.39	9.1	8.7	3.4	1.02	0.99	1.08	1.08	0.92

#### Speedup of the track fitter

	Time [μs/track]	Speedup
Initial	1200	-
Optimization	13	92
SIMDization	4.4	3
Multithreading	0.5	8.8
Final	0.5	2400

**Throughput: 2\*10<sup>6</sup> tracks/s**

Computer with 2xCPU Intel Core i7 (8 cores in total) at 2.67 GHz

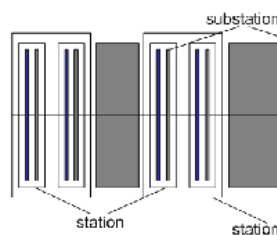
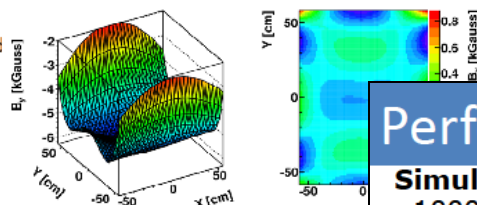
A. Lebedev, "Track reconstruction in the muon system and transition radiation detector of the CBM experiment"



## “Fast track reconstruction of the muon system and transition radiation detector”

### Optimization of the algorithm

- Minimize access to global memory
  - Approximation of the 70 MB large magnetic field map
    - 7 degree polynomial in the detector planes was proven the best
- Simplification of the detector geometry
  - Problem
    - Monte-Carlo geometry consists of 800000 nodes
    - Geometry navigation based on ROOT TGeo
    - Take into account absorbers and staggered Z positions of the stations
  - Solution
    - Create simplified geometry by converting Monte-Carlo geometry
    - Implement fast geometry navigation for the simplified geometry
- Computational optimization of the Kalman Filter
  - From **double** to **float**
  - Implicit calculation on non-trivial matrix elements
  - Loop unrolling
  - Branches (if then else ..) have been eliminated



All these steps are necessary to implement SIMD tracking

A. Lebedev, "Track reconstruction in the muon system and transition radiation detector of the CBM experiment"

### Performance of the tracking in MUCH

#### Simulation:

- 1000 UrQMD events at 25 AGeV Au-Au collisions + 5  $\mu^+$  and 5  $\mu^-$  embedded in each event

	Initial version	Parallel version
Efficiency [%]	94.7	94.0

#### Speedup of the track finder

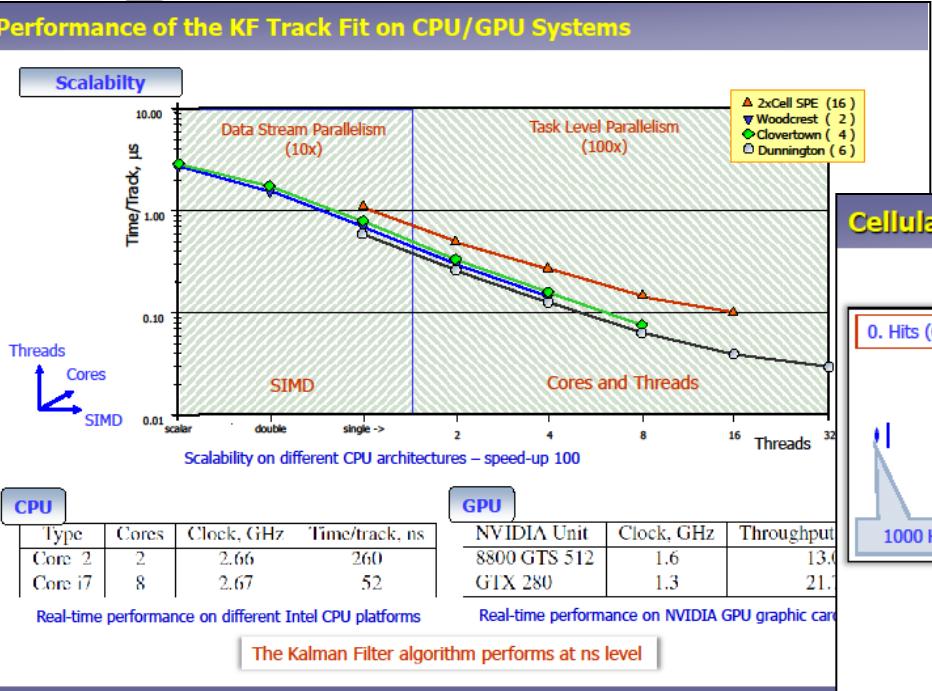
	Time [ms/event]	Speedup
Initial	730	-
Optimization	7.2	101
SIMDization	4.8	1.5
Multithreading	1.5	3.3
Final	1.5	487

Computer with 2xCPUs Intel Core i7 (8 cores in total) at 2.67 GHz

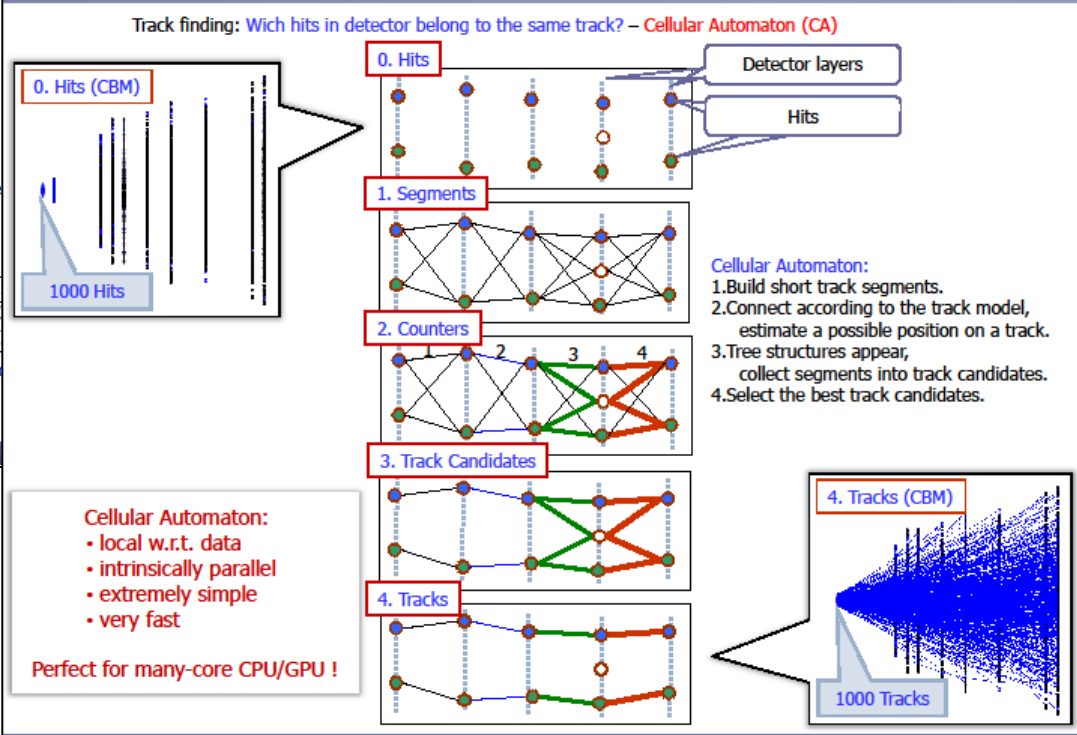
A. Lebedev, "Track reconstruction in the muon system and transition radiation detector of the CBM experiment"



- “Many-core scalability of the online reconstruction in CBM”

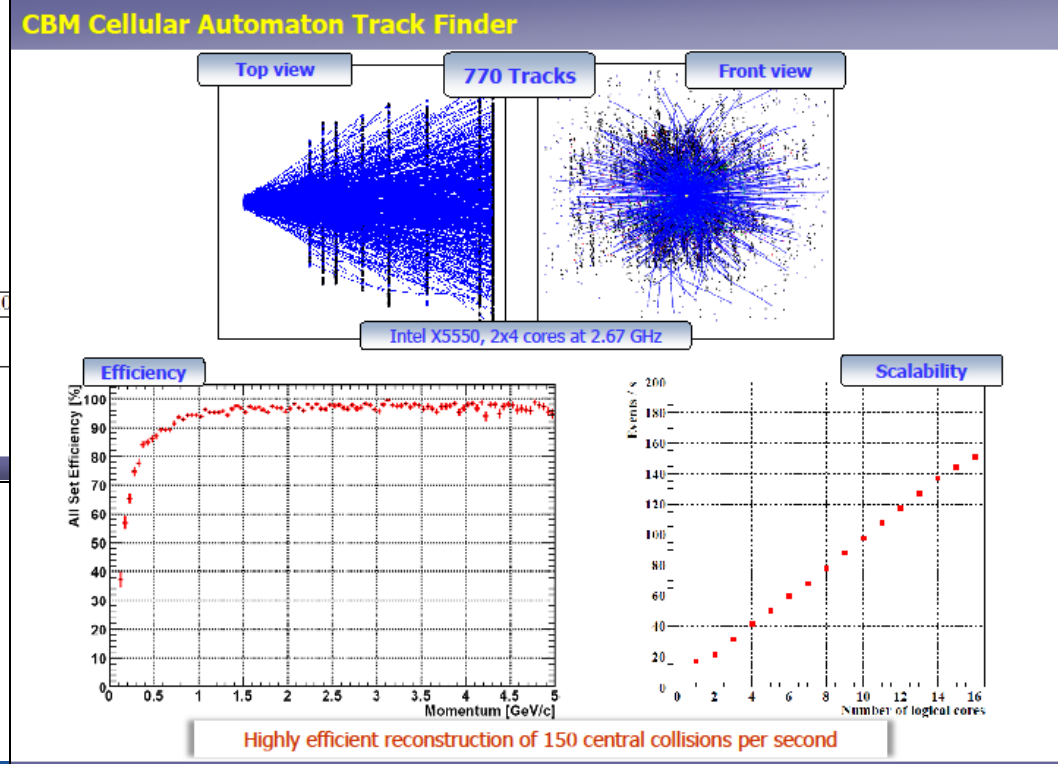
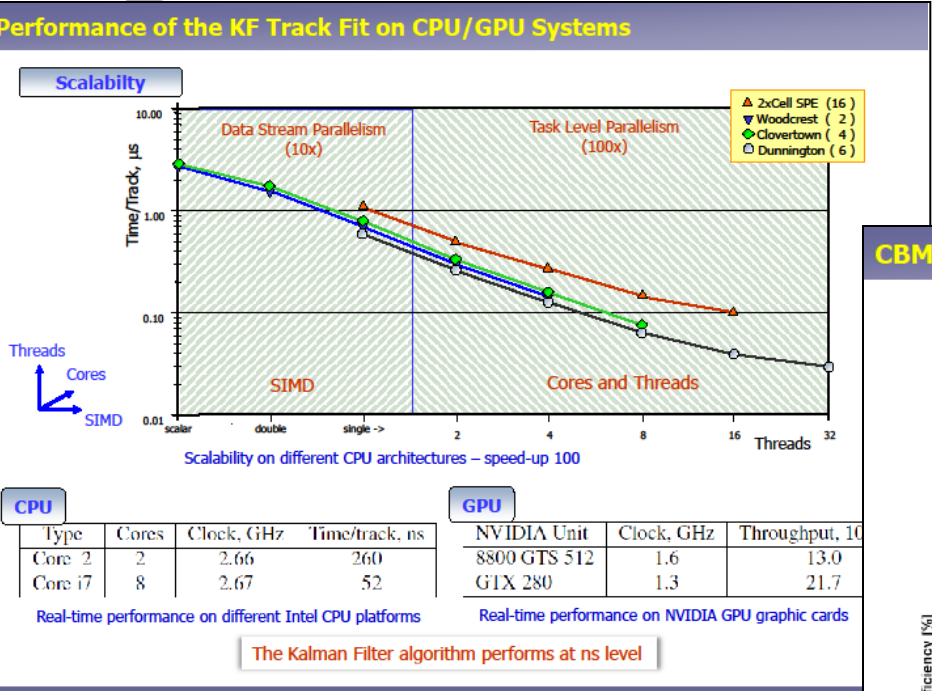


### Cellular Automaton (CA) as Track Finder





- “Many-core scalability of the online reconstruction in CBM”



## ■ GPUs

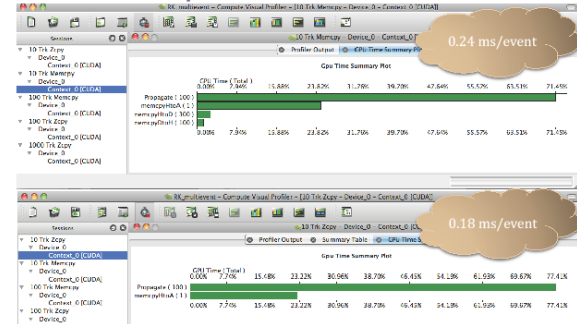
### Cuda Toolkit

- NVCC C compiler
- CUDA FFT and BLAS libraries for the GPU
- CUDA-gdb hardware debugger
- CUDA Visual Profiler
- CUDA runtime driver (also available in the standard NVIDIA GPU driver)
- CUDA programming manual

10/21/10 GPUs: High Performance Co-Processors

14

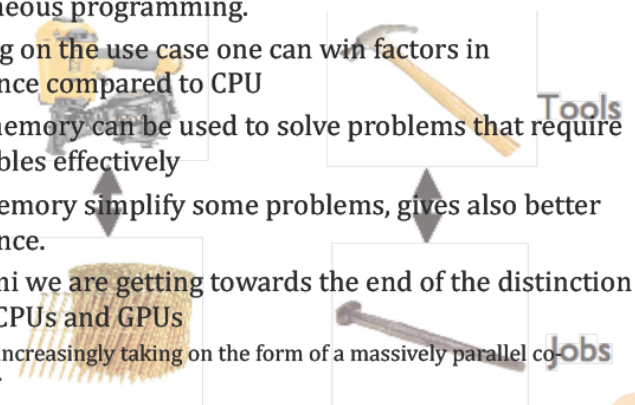
### NVIDIA Compute Visual Profiler (OpenCL/CUDA) (Runge-Kutta propagation Zcpy vs. Memcpy using FERMI GTX 480)



10/21/10 GPUs: High Performance Co-Processors

### Summary

- Cuda is an easy to learn and to use tool that allows heterogeneous programming.
- Depending on the use case one can win factors in performance compared to CPU
- Texture memory can be used to solve problems that require lookup tables effectively
- Pinned Memory simplify some problems, gives also better performance.
- With Fermi we are getting towards the end of the distinction between CPUs and GPUs
  - The GPU increasingly taking on the form of a massively parallel co-processor



10/21/10 GPUs: High Performance Co-Processors

27

- More and more people are working on multi-core, many-core, and accelerators
  - Especially, in the on-line domain
  - And, in new or upgraded experiments
  - Positive outcome: Code is revisited and made more “C-like”
- But, many people forget to do a “fair” comparison:
  - GPU code should expose “rich” loops for threading
    - Transfer times must be included
  - CPU code should exploit vectors + threads
- As usual it is important to perform a comprehensive calculation of
  - Throughput/W/CHF
- **Expect more activity in this domain in the coming months/years**
- In openlab we will continue to participate actively in realistic and relevant evaluations