# GPUs for Online Processing in Low-Level Trigger Systems

Massimiliano Fiorini

(Università di Ferrara and INFN Ferrara)

On behalf of the GAP Collaboration

# GAP Project



- **GAP** (**G**PU **A**pplication **P**roject) for real-time in HEP and medical imaging is a 3 years project funded by the Italian Ministry of research, started in 2013

- Collaboration between INFN Sezione di Pisa, University of Ferrara and University of Roma "La Sapienza"

- Demonstrate the feasibility of using off-the-shelf computer commodities to accelerate real-time scientific computations

- Application in different fields:
  - High Energy Physics (low and high level triggers)
  - Medical Imaging (NMR, CT and PET)

# GAP Project

- **GAP** (**G**PU **A**pplication **P**roject) for real-time in HEP and medical imaging is a 3 years project funded by the Italian Ministry of research, started in 2013
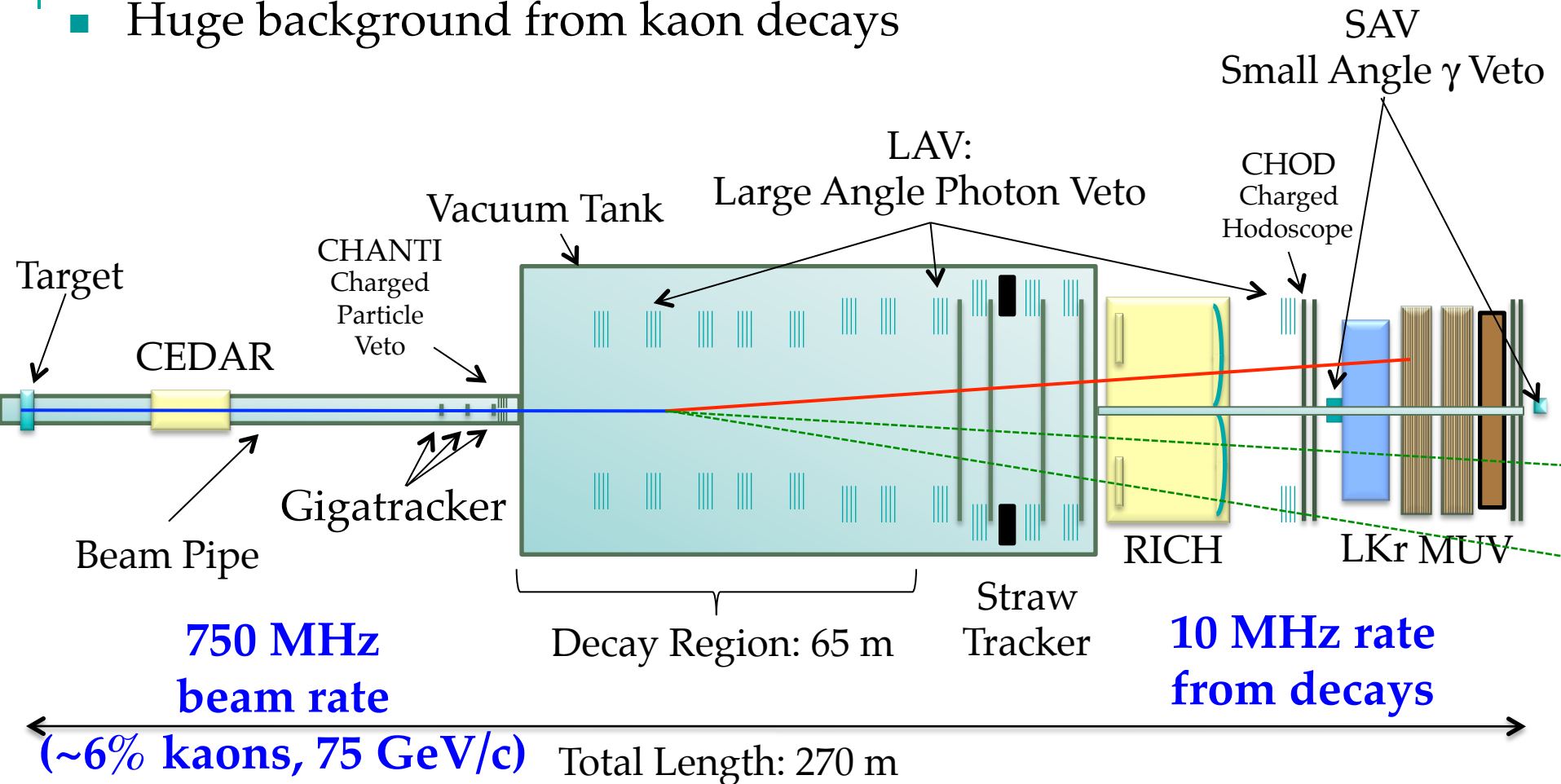
- Collaboration between INFN Sezione di Pisa, University of Ferrara and University of Roma "La Sapienza"

- Demonstrate the feasibility of usin~~g~~ commodities to acce~~l~~

- App~~l~~

  - ~~Energy~~ ~~P~~hysics (low and high level triggers)
  - ~~M~~edical Imaging (NMR, CT and PET)

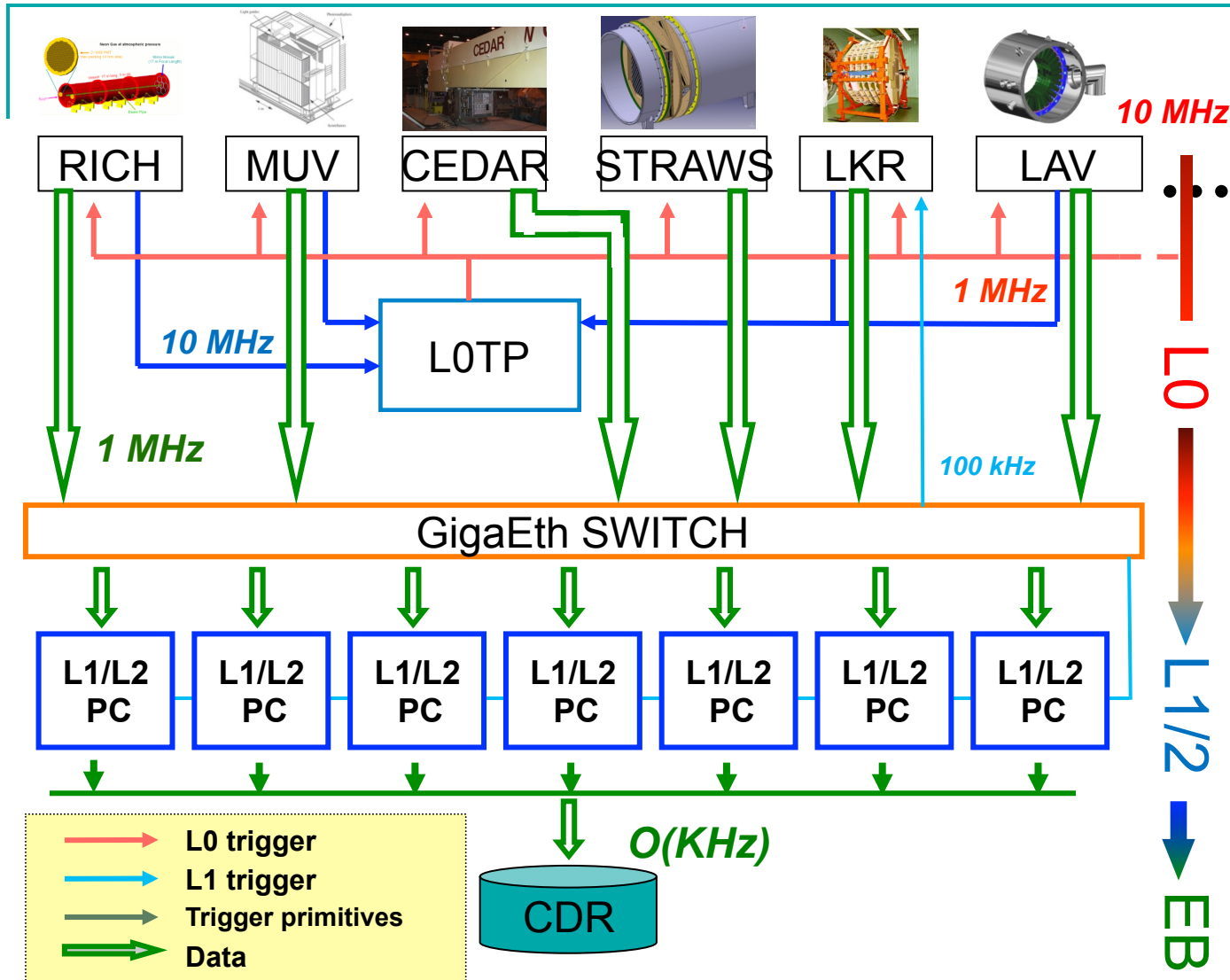**Poster: M. Fiorini "The GAP Project - GPU for Realtime Applications in High Level Trigger and Medical Imaging"**

# Physics case: NA62

- $K^+ \to \pi^+ \nu\nu$ decay (BR~$8\times10^{-11}$)
- Huge background from kaon decays



**750 MHz beam rate**

**(~6% kaons, 75 GeV/c)**

Decay Region: 65 m

Total Length: 270 m

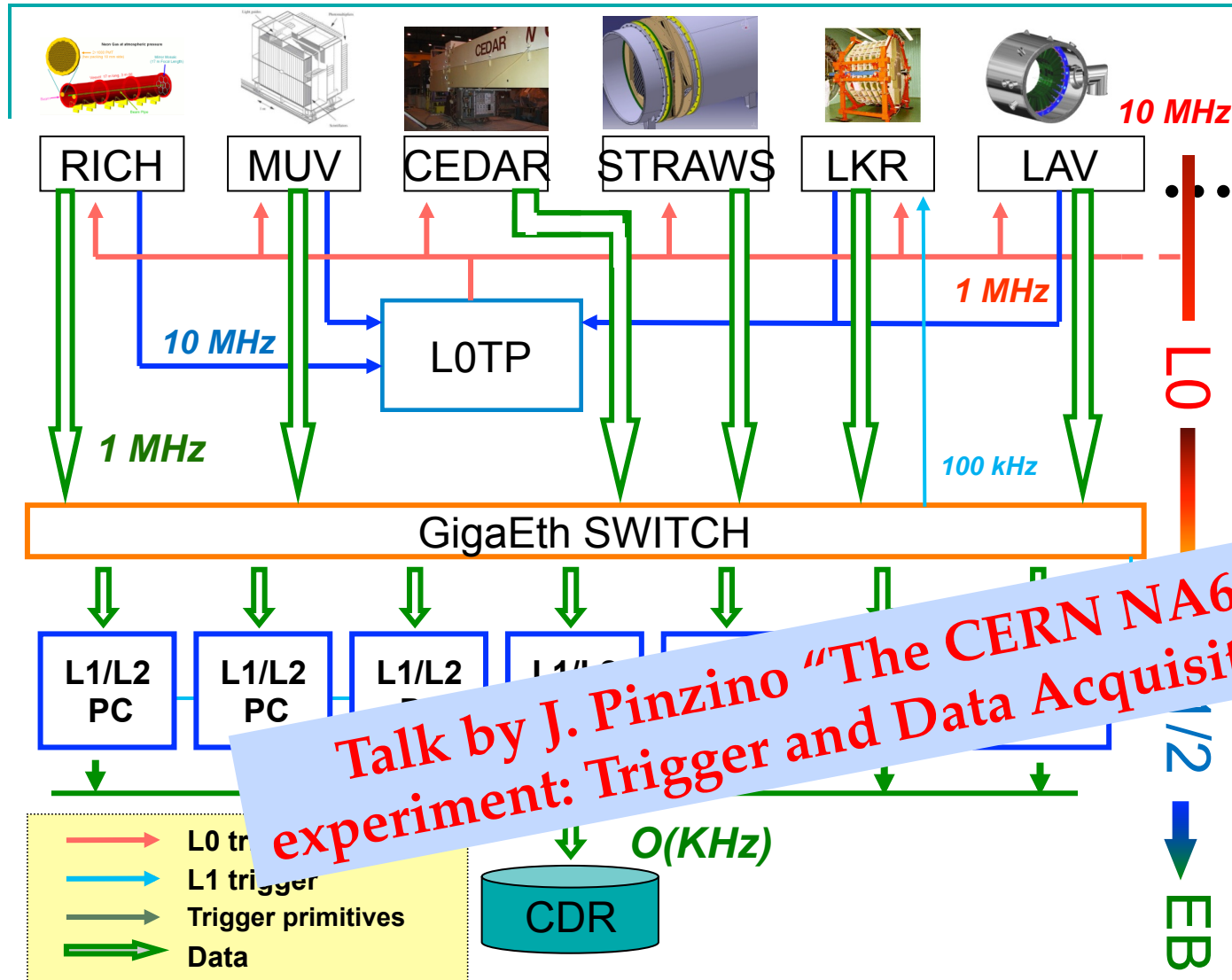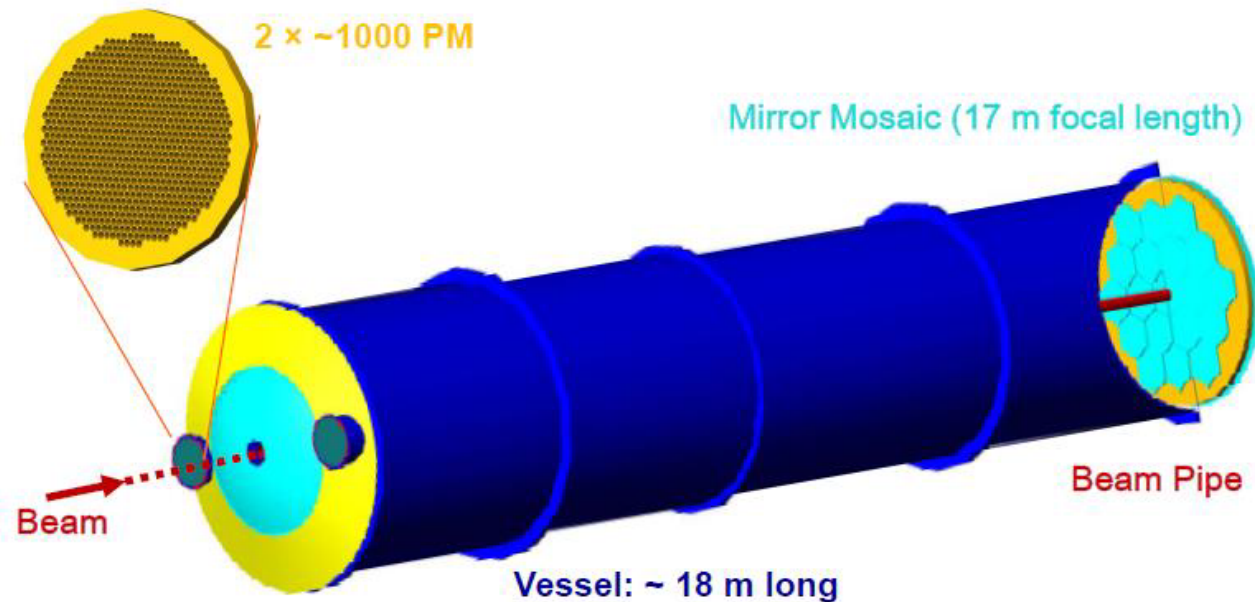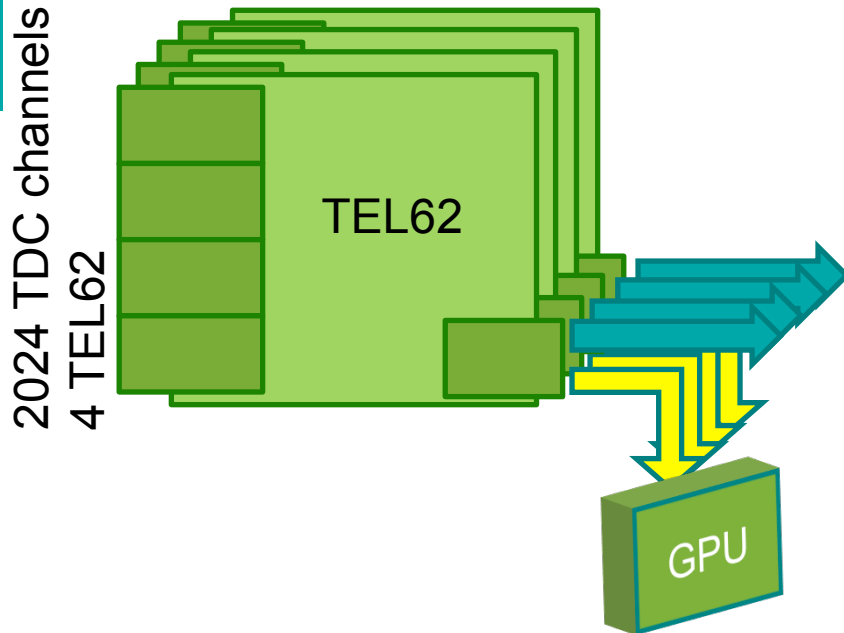**10 MHz rate from decays**

# Trigger and DAQ



- **L0**: Hardware **synchronous** level
  - 10 MHz to 1 MHz, **1 ms max. latency**
  - Primitives (MUV, RICH, LAV, LKR)
- **L1**: Software level
  - "Single detector", 1 MHz to 100 kHz
- **L2**: Software level
  - "Complete information", 100 kHz to 10 kHz

Legend:
- L0 trigger
- L1 trigger
- Trigger primitives
- Data

Diagram labels: RICH, MUV, CEDAR, STRAWS, LKR, LAV, L0TP, GigaEth SWITCH, L1/L2 PC, CDR

10 MHz, 1 MHz, 10 MHz, 100 kHz, O(KHz)

L0, L1/2, EB

# Trigger and DAQ



- **L0**: Hardware synchronous level
  - 10 MHz to 1 MHz, **1 ms max. latency**
  - Primitives (MUV, RICH, LAV, LKR)

- **L1**: Software level "Single detector", 1 MHz to 100 kHz

- **L2**: Software level
  - "Complete information", 100 kHz to 10 kHz

Talk by J. Pinzino "The CERN NA62 experiment: Trigger and Data Acquisition"

# The RICH detector

- 17 m focal length, 3 m in diameter, filled with Ne at 1 atm

- Pion/muon separation in the range 15-35 GeV/c

- Time resolution ~70 ps

- Mis-identification $5 \times 10^{-3}$

- 2 spots of 1000 PMTs each



2 × ~1000 PM

Mirror Mosaic (17 m focal length)

Beam Pipe

Beam

Vessel: ~ 18 m long

- 10 MHz events rate in the RICH (~20 hits/track)
  - Main contribution from kaon decays (~1 MHz from halo muons and pion decays)

# The goal: **GPU in L0 RICH**

**2024 TDC channels**
**4 TEL62**

TEL62

GPU

- **4 TEL62 for RICH detector**
  - ☐ 8×1Gb/s links for data r/o
  - ☐ 4×1Gb/s trigger primitives
  - ☐ 4×1Gb/s GPU trigger
- **Events rate: 10 MHz**
- **L0 trigger rate: 1 MHz**
- **Max Latency: 1 ms**

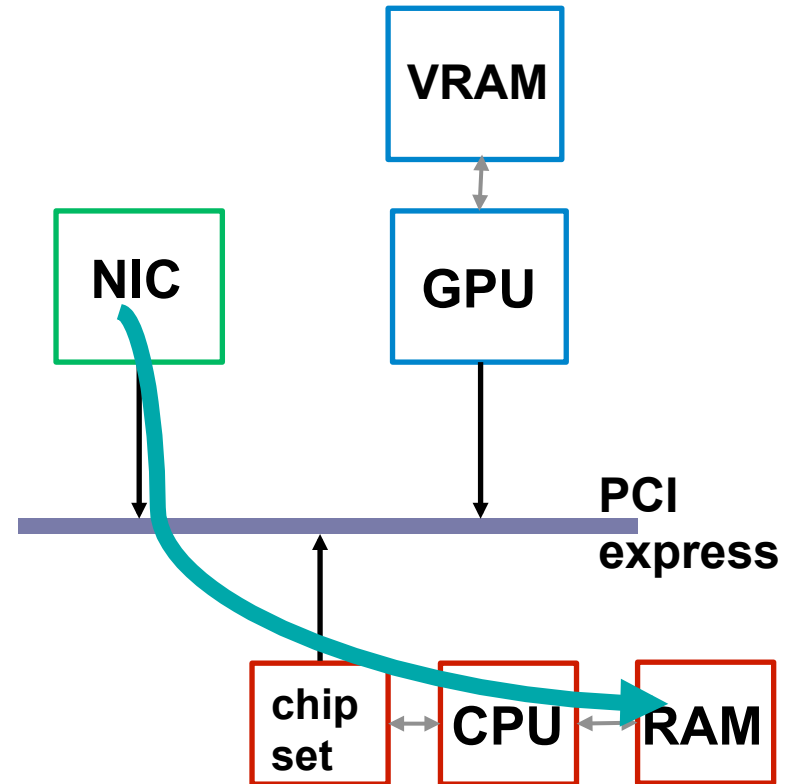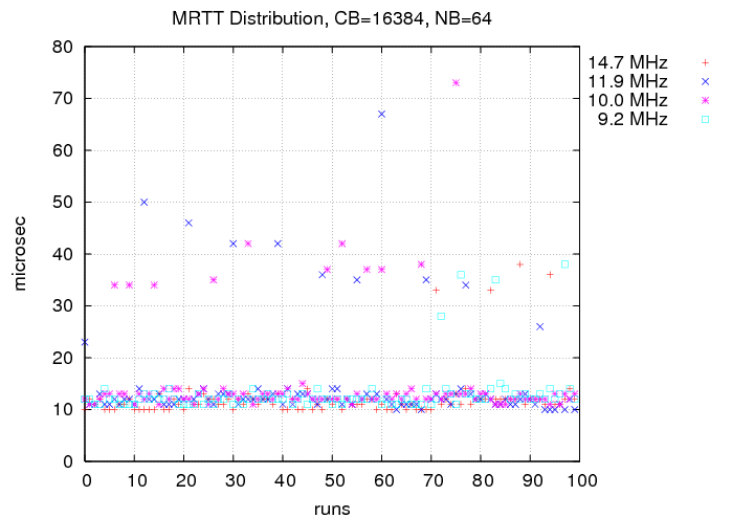- **This is not the L0 trigger baseline version for the NA62 RICH detector**

**TEL62 RO buffer** → **Reduced rate** → **L1** → **L2**

**L0TP**

# GPUs in Low Level Triggers

- Two main issues to be solved:

- **Latency**
  - Is the GPU latency per event small enough to cope with the tiny latency of low level triggers?
  - Is the latency stable enough for usage in synchronous trigger systems?

- **Computing power**
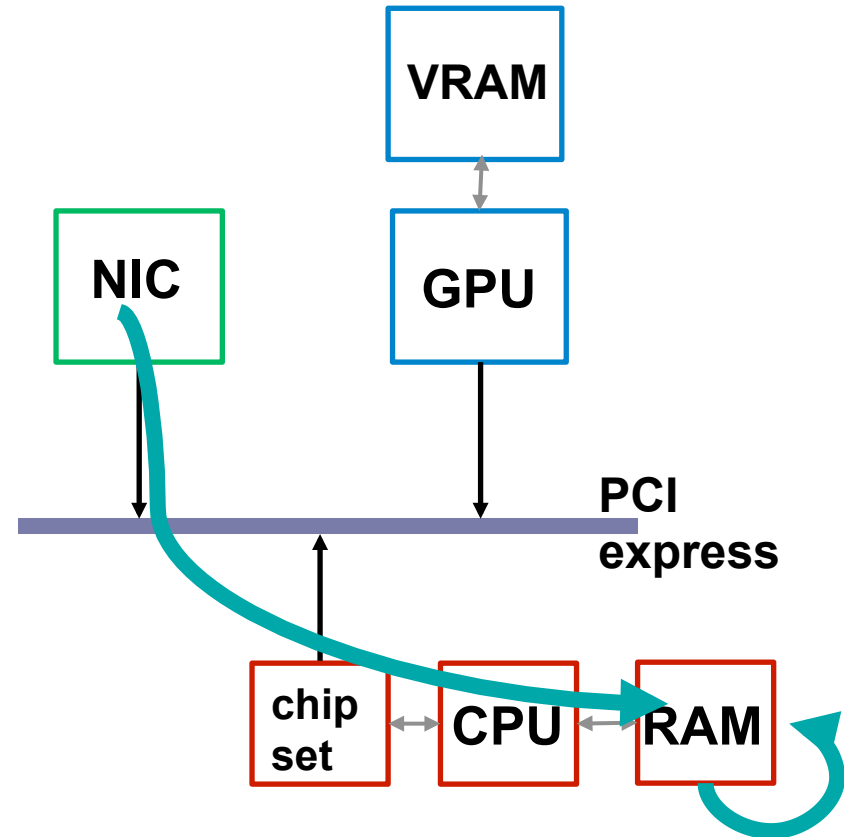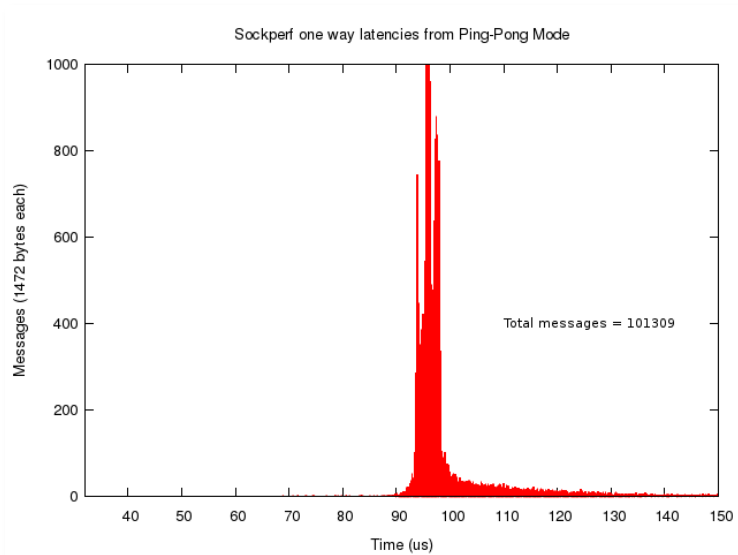  - Is the GPU fast enough to take a trigger decision at tens of MHz events rate?

# GPU Processing

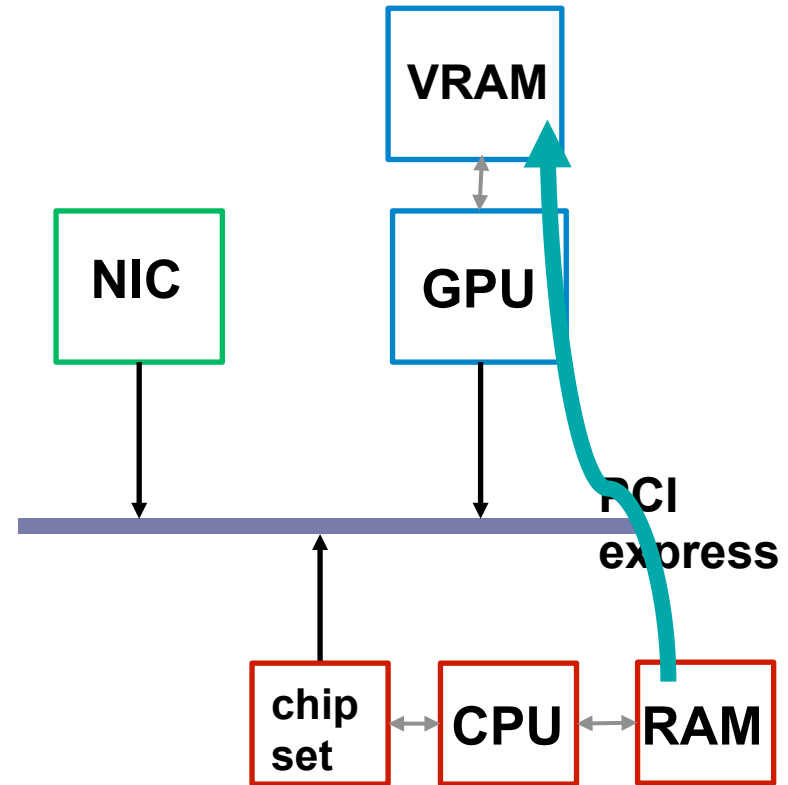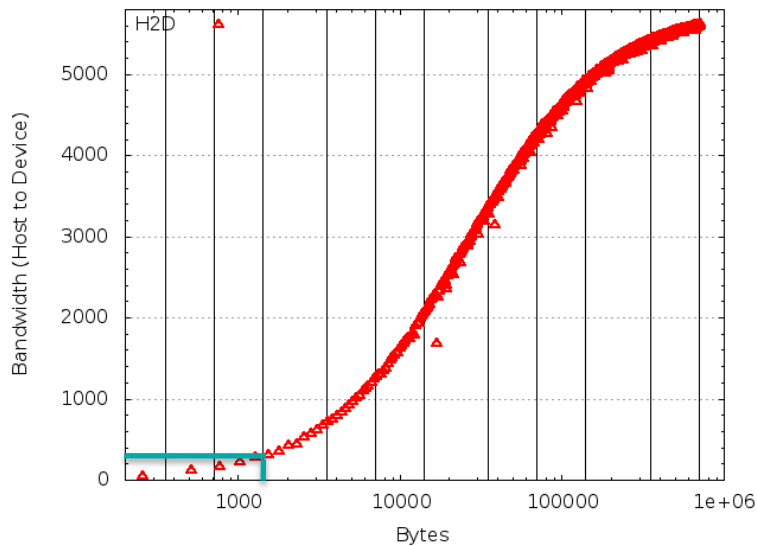- Example: packet with 1404 B (few tens of events in NA62 RICH application)
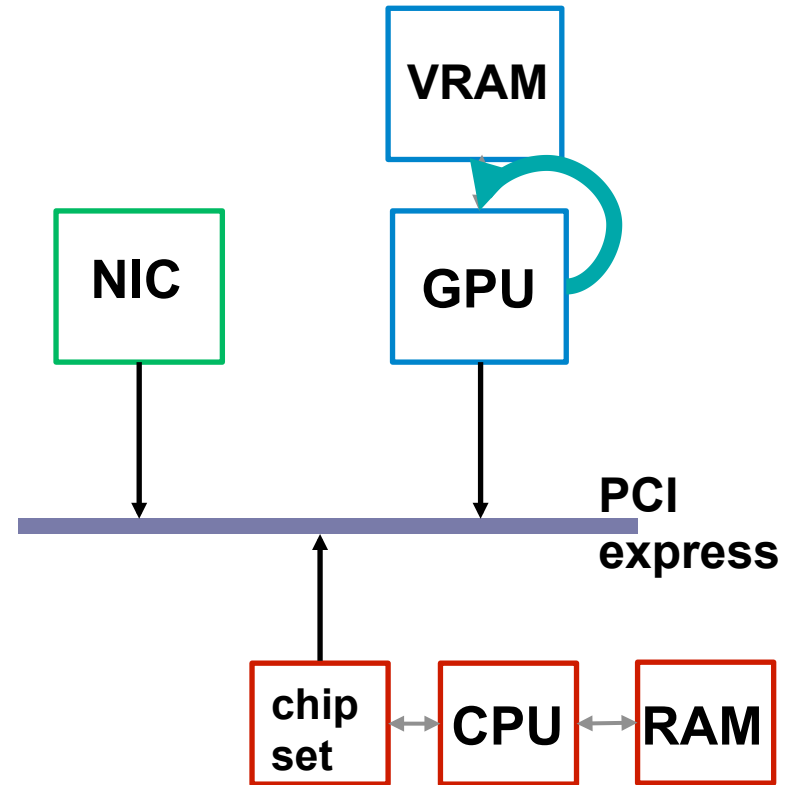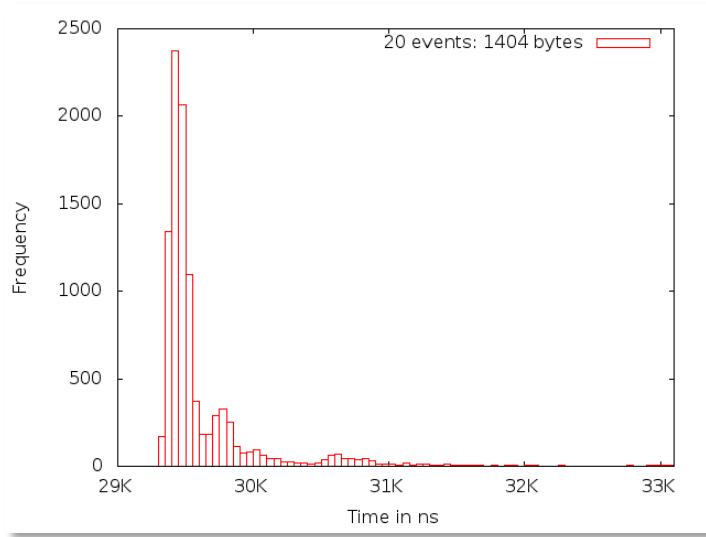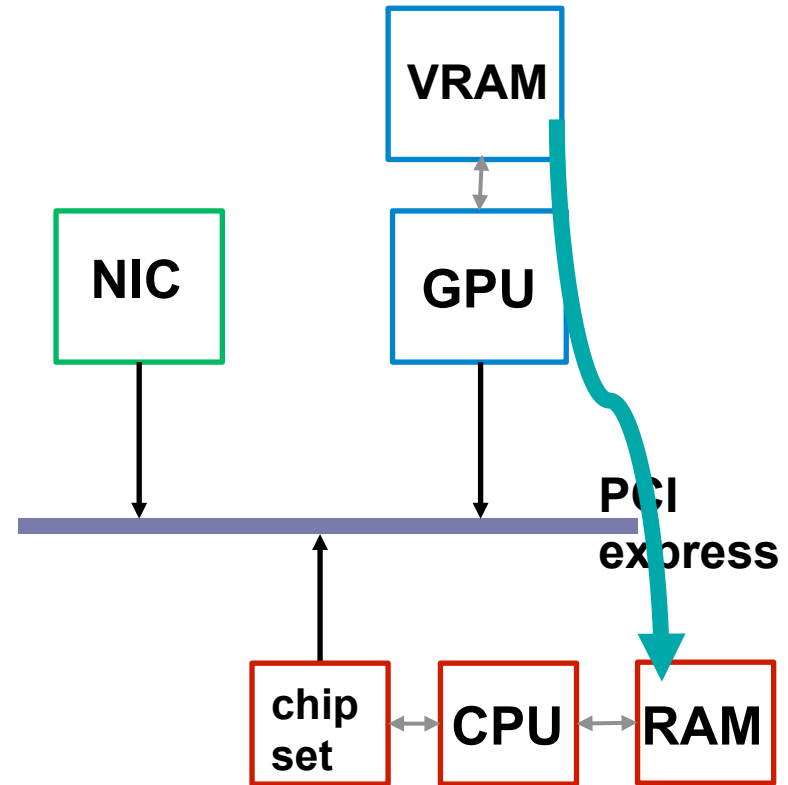- T=0

# GPU Processing
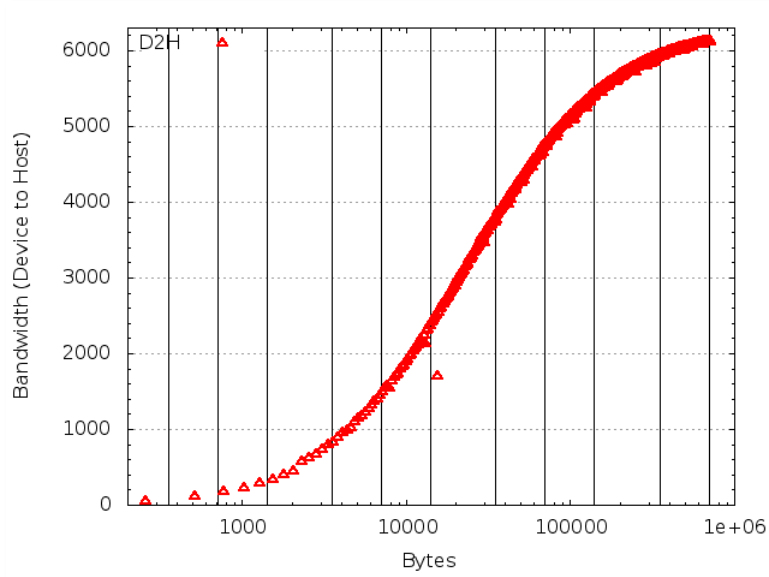
# GPU Processing

# GPU Processing
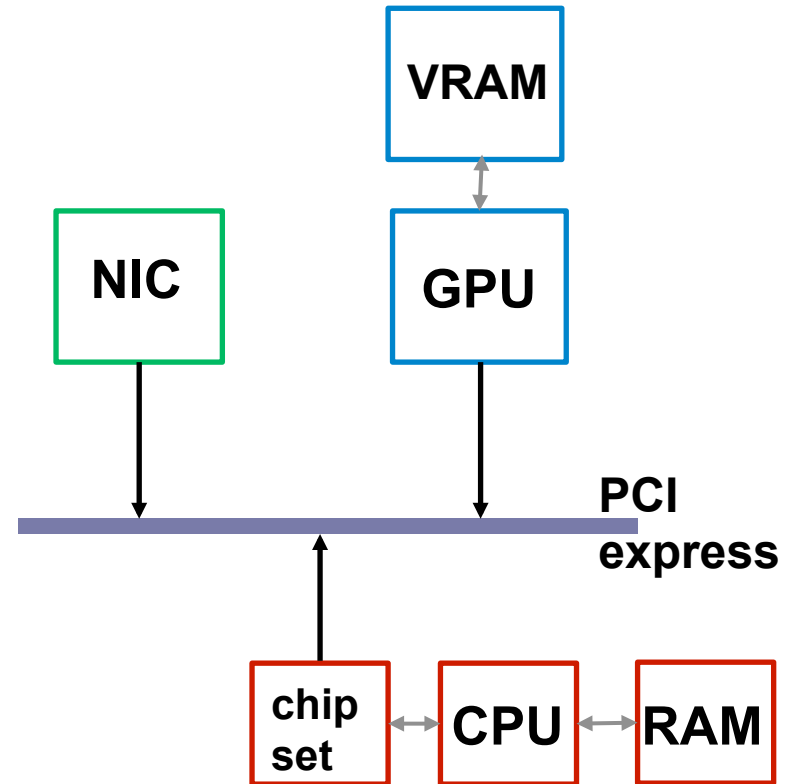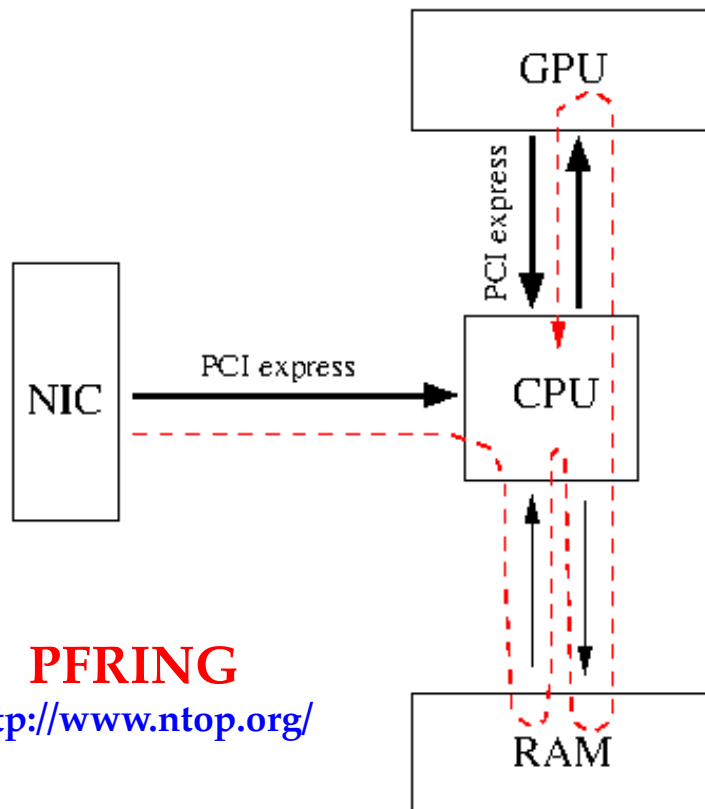
# GPU Processing

# GPU Processing

# GPU Processing

- Latency due to data transfer from the detector to the system is bigger than the latency due to GPU computing

- It scales almost linearly (apart from the overheads) with the data size while the latency due to computing can be hidden exploiting the huge resources

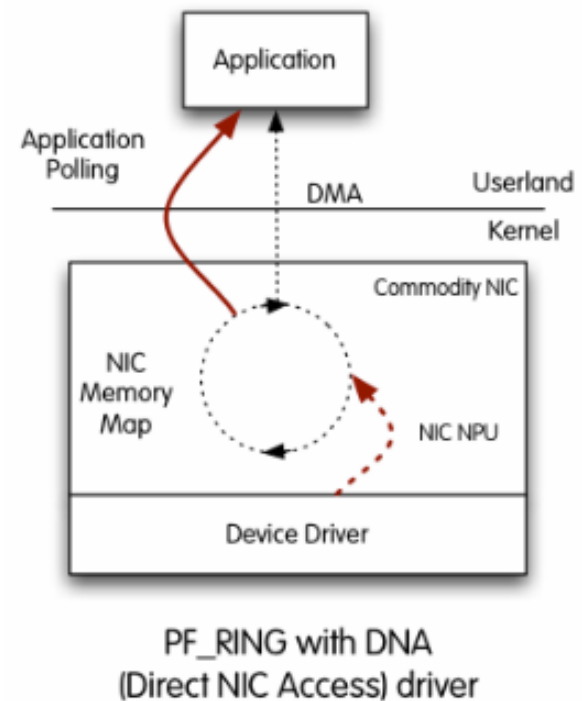- Communication latency fluctuations quite big (~50%)

# First solution: **PFRING**

■ PFRING DNA (Direct NIC Access) is a way to map NIC memory to userland so that there is no additional packet copy besides the DMA transfer done by the NIC
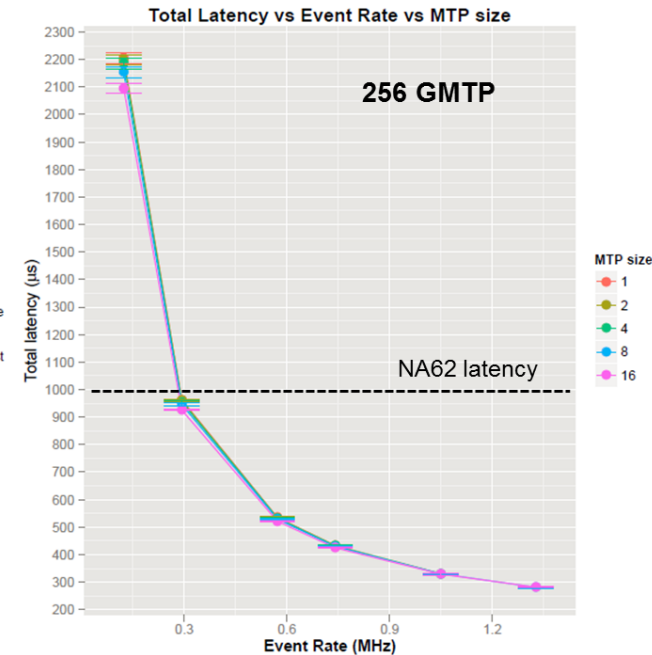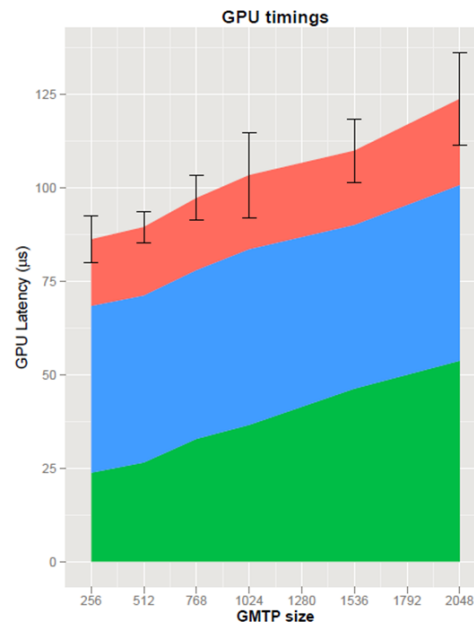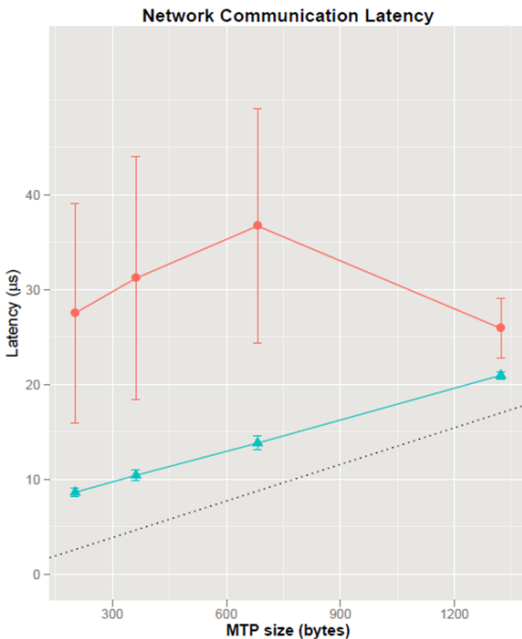


**PFRING**
http://www.ntop.org/
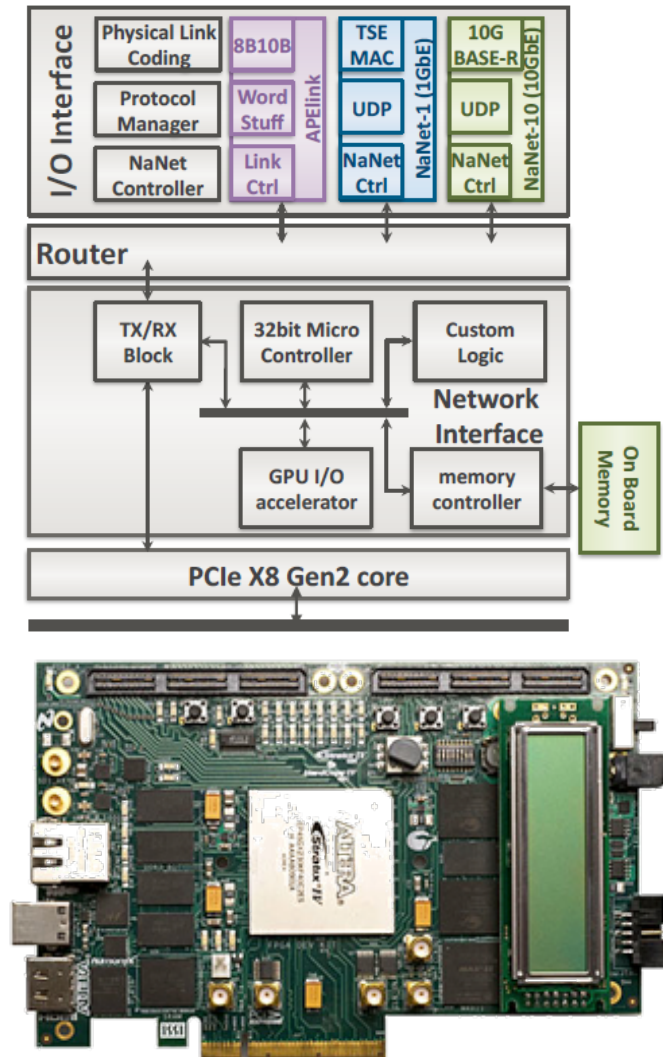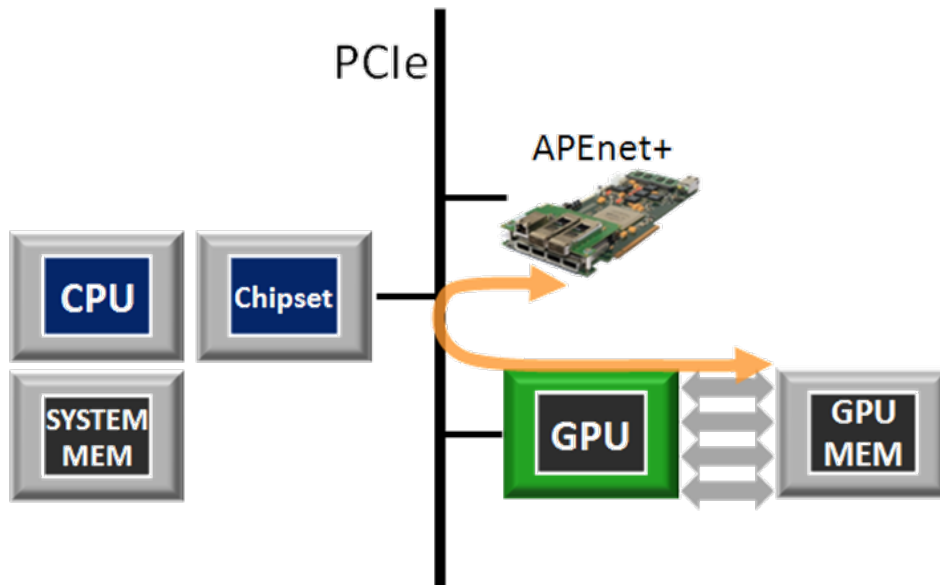
PF_RING with DNA
(Direct NIC Access) driver

# Results: PFRING

- Latency reduced by a factor 3 and negligible fluctuations
- The total latency is given as a function of the number of events to buffer before the start of GPU computation
- For real application the "working point" depends on the events rate and event dimension

# Second solution: NANET

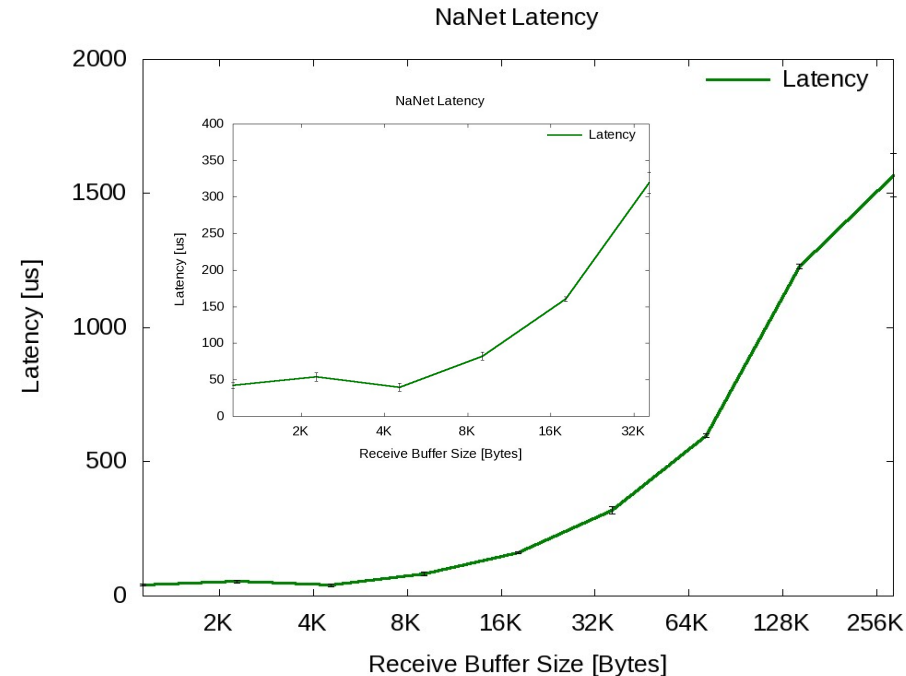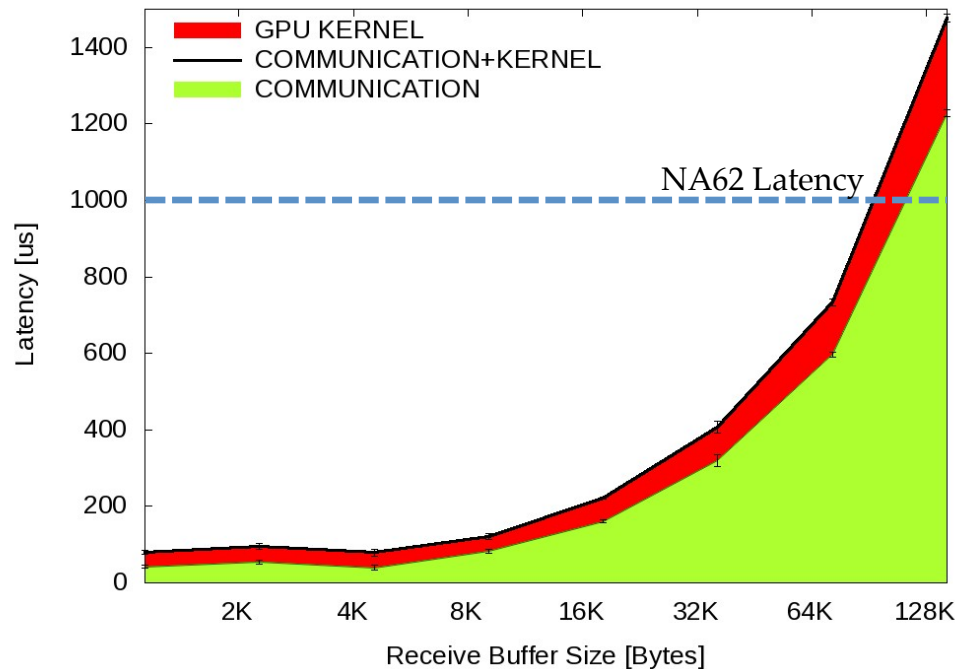- NANET is an FPGA-based NIC that has GPUDirect RDMA capabilities



**APEnet Rome Group**
**R. Ammendola *et al., JINST* 9 C02023, 2014**

# Results: NANET

- Variable size buffers sent to the GPU memory

  - Data sent from GbE NIC from the same PC hosting NANET


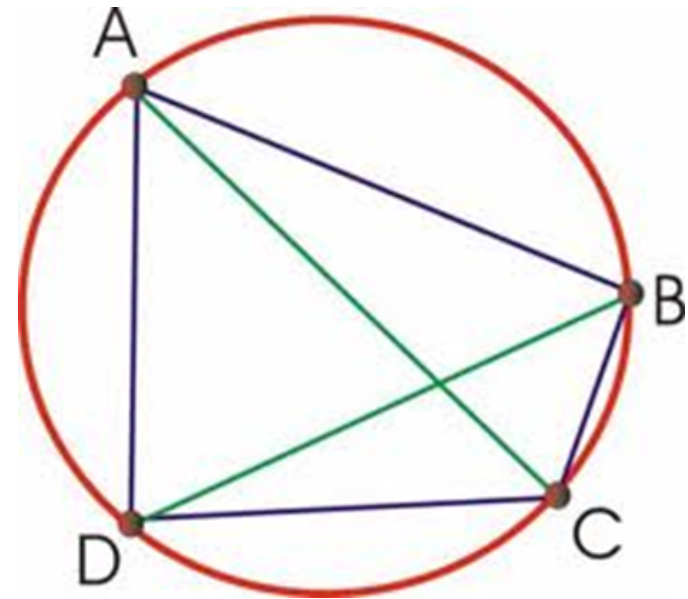Communication and Kernel Calculation Latencies


NaNet Latency

- GPU kernel for pattern matching is then executed

  - Maximum buffer size limited to <90 KBytes
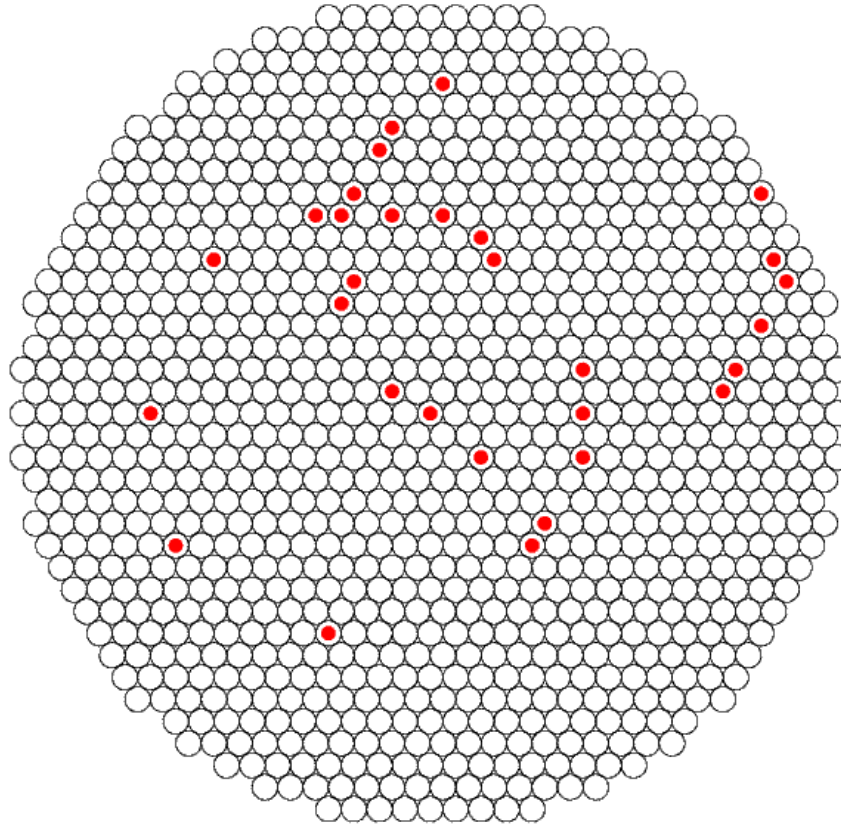
# L0 RICH trigger algorithm

- Requirements for an on-line RICH reconstruction algorithm:
- Trackless
  - No information from the tracker
  - Difficult to merge information from many detectors at L0
- Multi-rings
  - Many-body decay in the RICH acceptance
- Fast
  - Non-iterative procedure
  - Events rate at a level of ~10 MHz
- Low latency
  - Online (synchronous) trigger
- Accurate

# Almagest





- New algorithm (Almagest) based on Ptolemy's theorem: *"A quadrilater is cyclic (the vertex lie on a circle) if and only if is valid the relation: AD\*BC+AB\*DC=AC\*BD "*

- Design a procedure for parallel implementation

# Almagest: example

Massimiliano Fiorini (Ferrara)

i) Select a *triplet* (3 starting points)

i) Select a *triplet* (3 starting points)

ii) Loop on the remaining points: if the next point does not satisfy the Ptolemy's condition then reject it

# Almagest: example

i) Select a *triplet* (3 starting points)

ii) Loop on the remaining points: if the next point does not satisfy the Ptolemy's condition then reject it

iii) If the point satisfy the Ptolemy's condition then consider it for the fit

i) Select a *triplet* (3 starting points)

ii) Loop on the remaining points: if the next point does not satisfy the Ptolemy's condition then reject it

iii) If the point satisfy the Ptolemy's condition then consider it for the fit
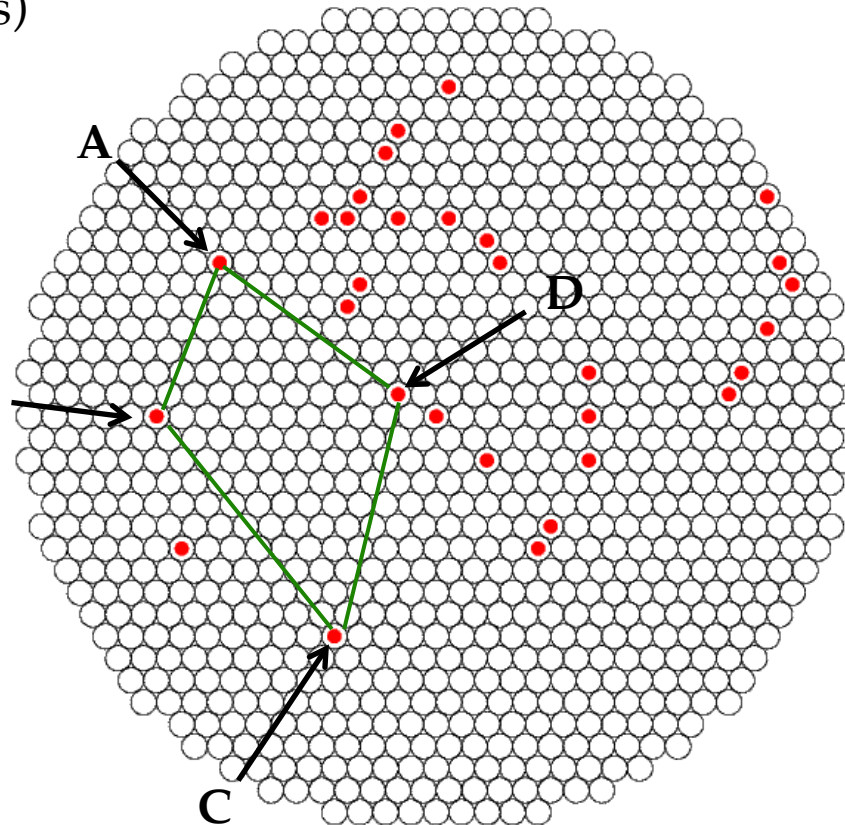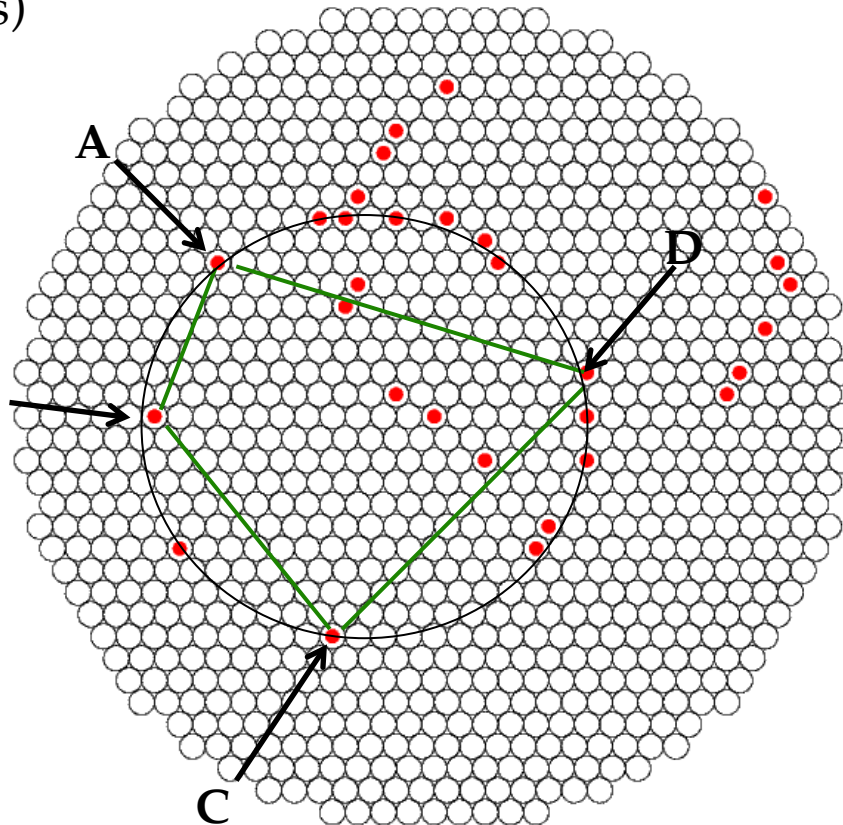
iv) …again…

A

B

C

D

# Almagest: example

i) Select a *triplet* (3 starting points)

ii) Loop on the remaining points: if the next point does not satisfy the Ptolemy's condition then reject it

iii) If the point satisfy the Ptolemy's condition then consider it for the fit

iv) …again…

v) Perform a single ring fit

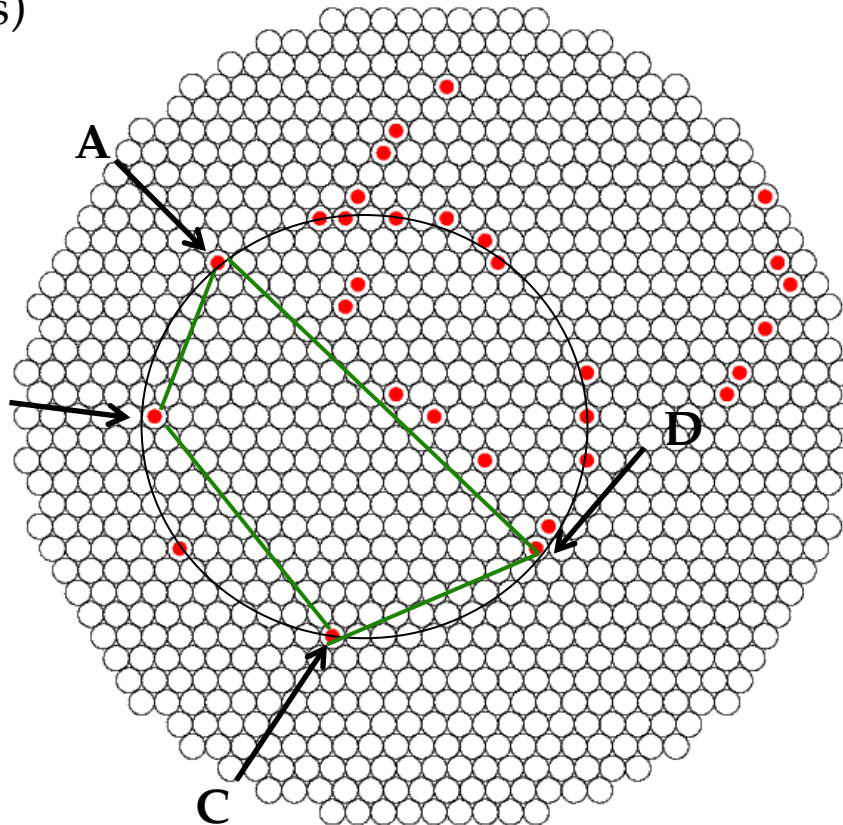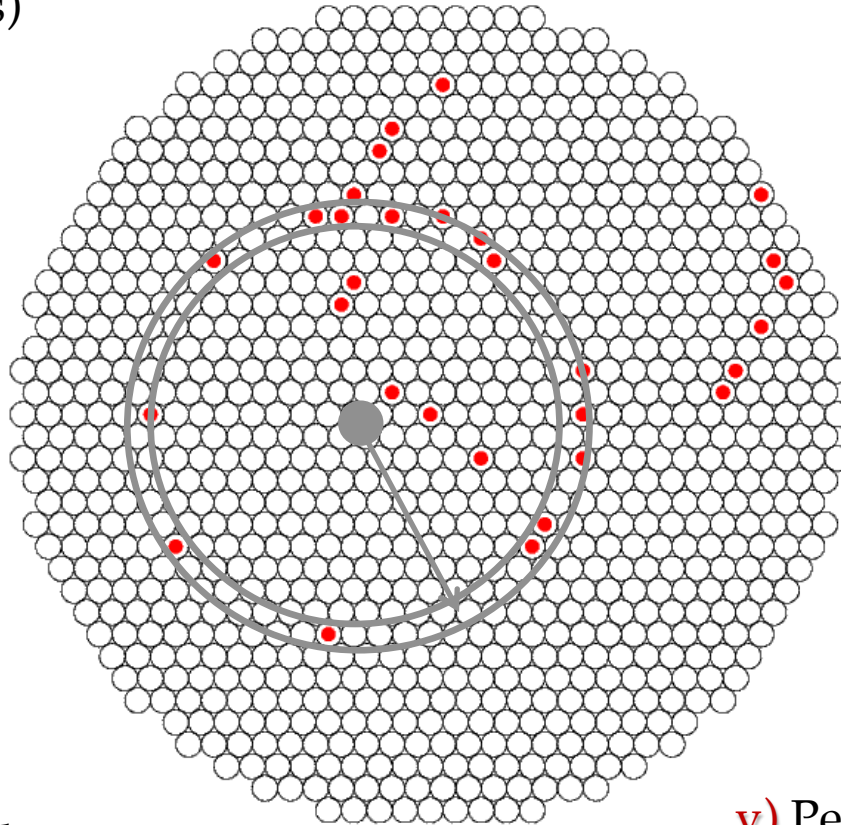# Almagest: example

i) Select a *triplet* (3 starting points)

ii) Loop on the remaining points: if the next point does not satisfy the Ptolemy's condition then reject it

iii) If the point satisfy the Ptolemy's condition then consider it for the fit

iv) …again…

v) Perform a single ring fit

vi) Repeat by excluding the already used points

# Almagest on GPU

- **Very high parallelism**
  - Huge number of computing cores (>2000)
  - Huge memory bandwidth



- **Two levels of parallelism**
  - Several triplets run in parallel
  - Several events at the same time

# Almagest: implementation

- Tests on NVIDIA Tesla K20 GPU

- Total computing time order **a few µs** per event (on single GPU)

- Good efficiency (using 8 triplets)

  - Room for improvement

- Further tests ongoing to study noise immunity, bias, efficiency a function of the number of hits, etc…





Total GPU trigger execution time

# Next steps

- Receive the TTC stream (timing and trigger) from the experiment
  - TTC interface board recently produced
- Integration in the NA62 Trigger and DAQ system
  - Dry runs in summer
  - Parasitic test during NA62 experimental run in October

# Conclusions

- The use of GPUs in HEP trigger systems could give several advantages, but processing performances and latencies should be carefully studied
  - Data transfer is the dominant contribution
- Construction of a demonstrator L0 processor for the NA62 RICH is under way
  - Cherenkov rings pattern recognition within the total L0 latency of 1 ms seems possible
- Integration with the NA62 Trigger and DAQ system
  - Dry runs in summer 2014
  - Parasitic data taking during NA62 experimental run starting October 2014

# SPARES

# Latency measurement

- Not easy to measure (with high precision) the latency between two systems running with uncorrelated clocks
- Use an oscilloscope to have a common time reference
  - Use the start of the packet in the TEL62 as "start" and the "computation done" in the PC as "stop"
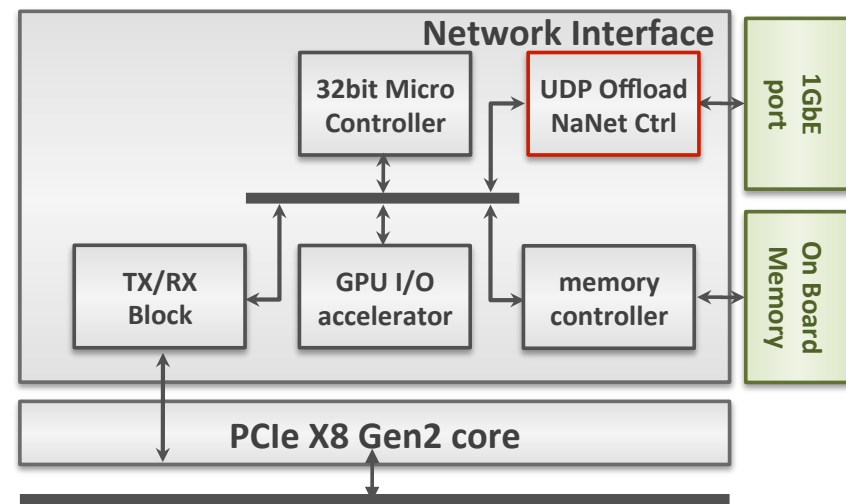
# NANET (1)

- An **FPGA** based PCIe 8x gen 2 board derived from the **APEnet+ 3D NIC** design, able to directly inject an UDP data stream from an external **GbE** source into the memory of a Fermi- or Kepler-class NVIDIA GPU exploiting GPUDirect RDMA capabilities

- PCIe P2P protocol between Nvidia Fermi/Kepler devices and APEnet+

- First non-NVIDIA device (2012) supporting **GPUDirect RDMA**

  - No bounce buffers on host. APEnet+ can target GPU memory with no CPU involvement

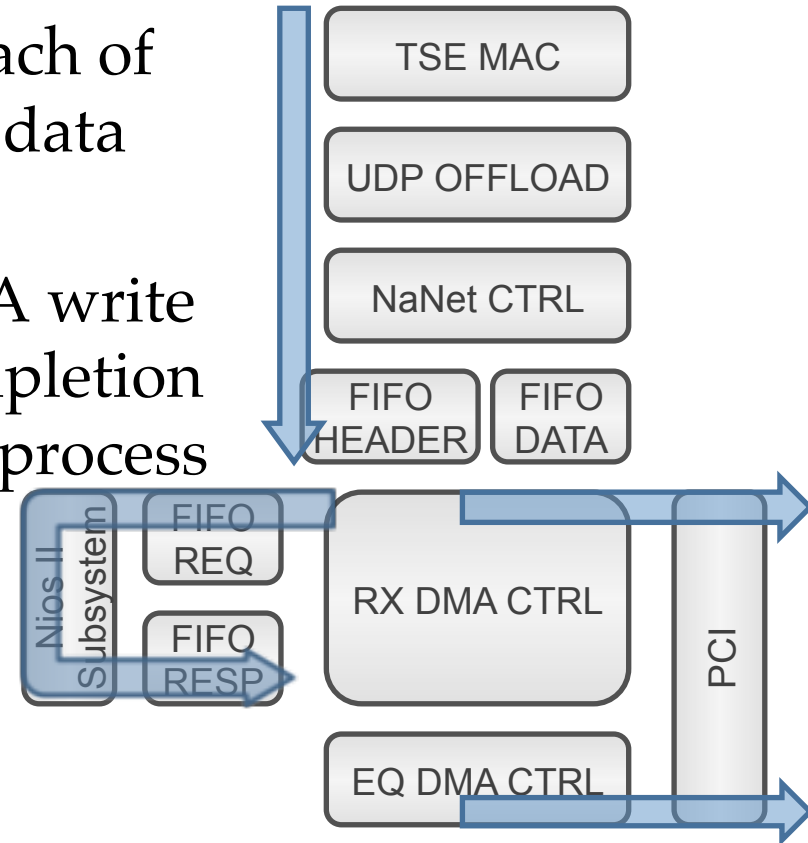  - **Latency reduction for small messages**

# NANET (2)

- **<u>NaNet</u>**: reduce comm. latency and its fluctuations to meet real-time constraints

- UDP offload collects data coming from the Altera Triple-Speed Ethernet Megacore (TSE MAC) and redirects UDP packets into a hardware processing data path

- NaNet Controller encapsulates the UDP payload in a newly forged APEnet+ packet and sends it to the RX Network Interface logic

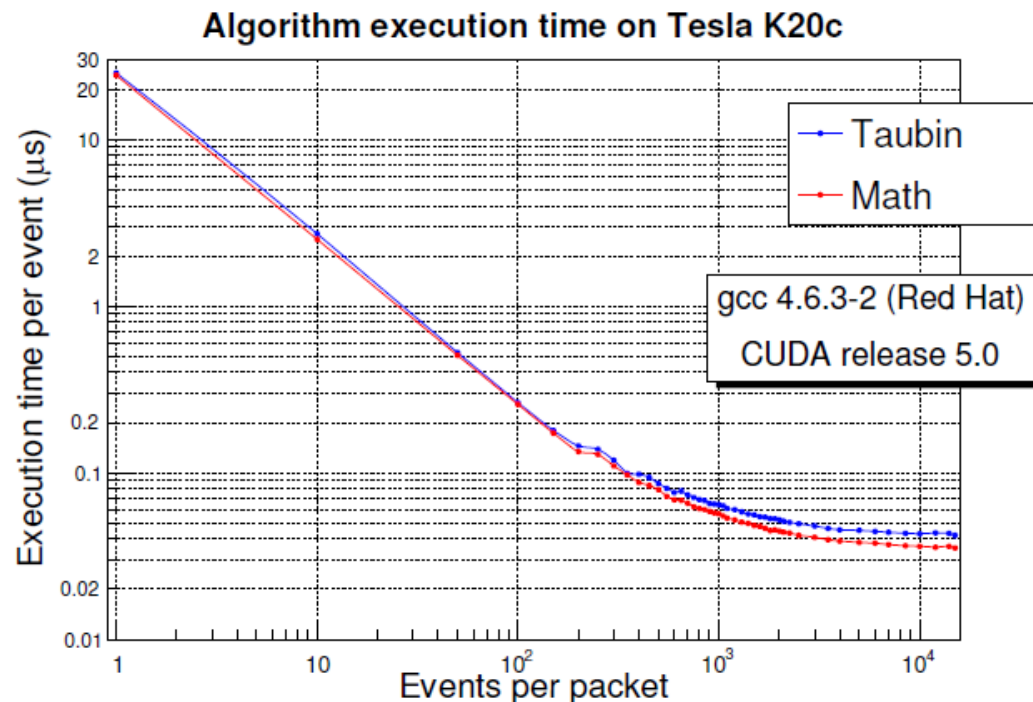- RX DMA CTRL manages CPU/GPU memory write process

# NANET (3)

- Nios II handles all the details pertaining to buffers registered by the application to implement a zero-copy approach of the RDMA protocol (OUT of the data stream)

- EQ DMA CTRL generates a DMA write transfer to communicate the completion of the CPU/GPU memory write process

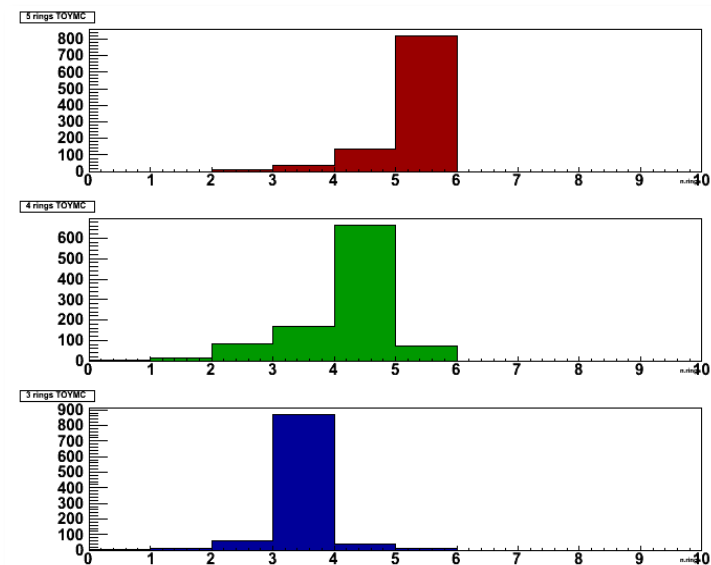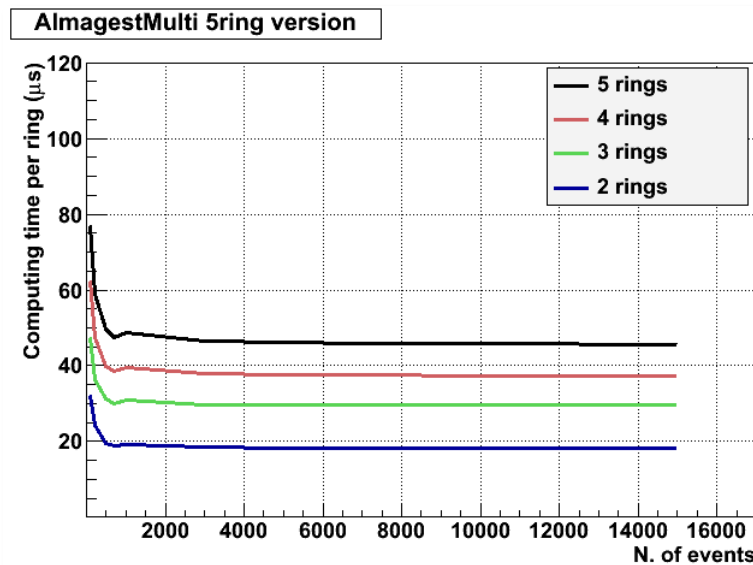- A Performance Counter is used to analyze the latency of the GbE data flow inside the NIC



TSE MAC

UDP OFFLOAD

NaNet CTRL

FIFO HEADER    FIFO DATA

Nios II Subsystem

FIFO REQ

FIFO RESP

RX DMA CTRL

PCI

EQ DMA CTRL

# Single ring

- Crowford method ("math"):
  - Translate in the center of mass
  - Least square minimization → linear
- Taubin method:
  - More efficient: minimize the bias introduced by the Kasa related methods (minimization of simple algebraic distance)
  - Resolution slightly better (on identified rings)
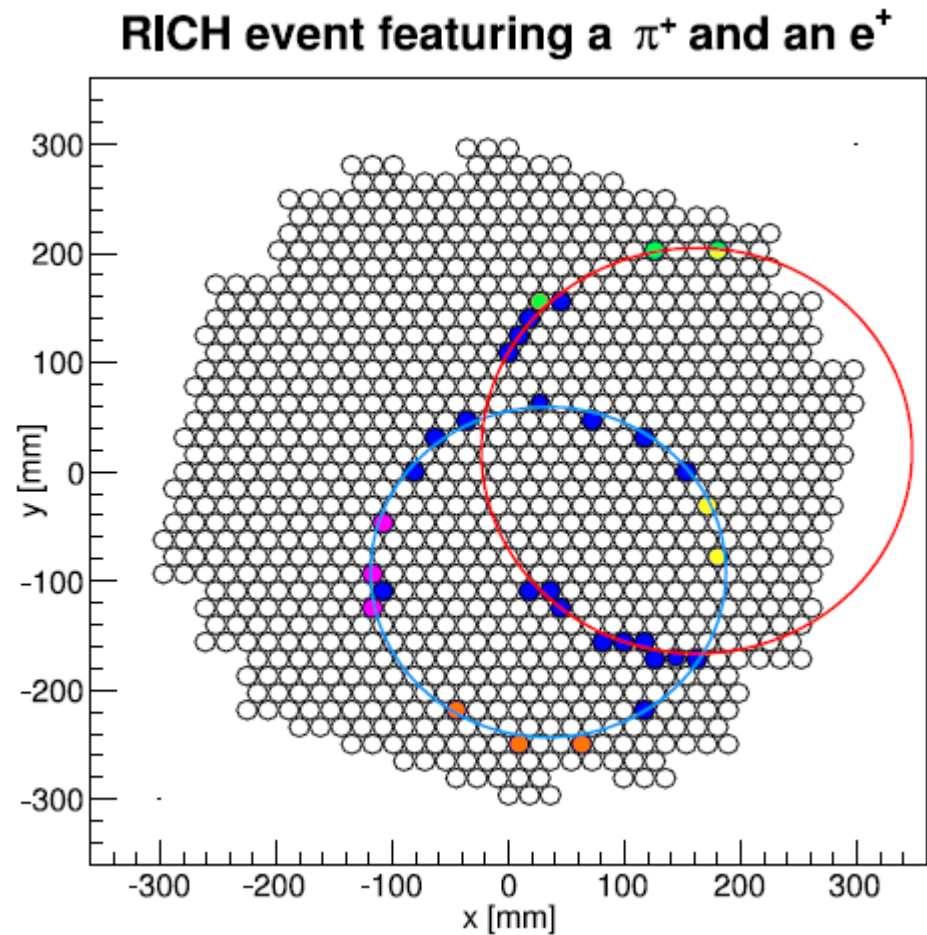- The difference of computing time on the GPU is at level of 10 ns per event



Algorithm execution time on Tesla K20c

gcc 4.6.3-2 (Red Hat)
CUDA release 5.0

# N(=hits) triplets

- Number of triplets equal to the number of hits.
- Relatively high efficiency.
- Computing time depends on number of rings (different number of GPU cores per events)
- Results on TESLA C1060 (240 cores, less than 1 Tflops)
- Room for optimization

# 4 selected triplets

- Only 4 triplets per event are used: left, right, up and down
- Further cuts to avoid too close hits

**RICH event featuring a $\pi^+$ and an $e^+$**

# 4 selected triplets

- Stability with small noise (studies are ongoing)
- Inefficiency due to the order in choosing the rings.
- Dependence on the cuts to define the triplets.