

# A new paradigm using GPUs for fast triggering and pattern matching at the NA62 CERN experiment

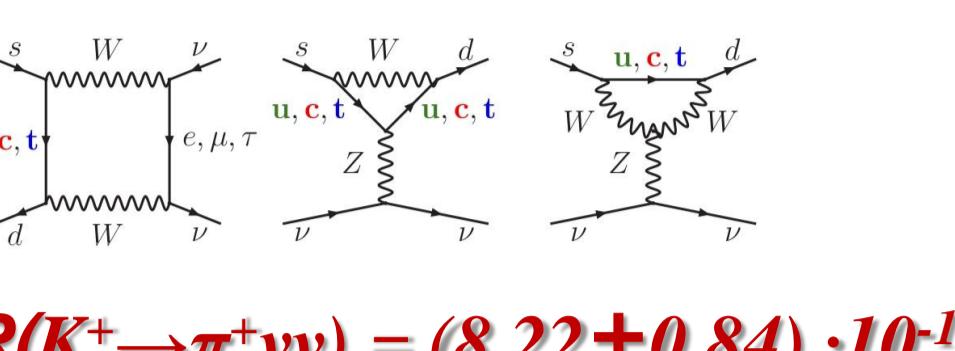


G.Collazuol<sup>1,2</sup>, G.Lamanna<sup>1,2</sup>, M.Sozzi<sup>2,3</sup>

(<sup>1</sup>Scuola Normale Superiore of Pisa, <sup>2</sup>INFN sez.Pisa, <sup>3</sup>University of Pisa)

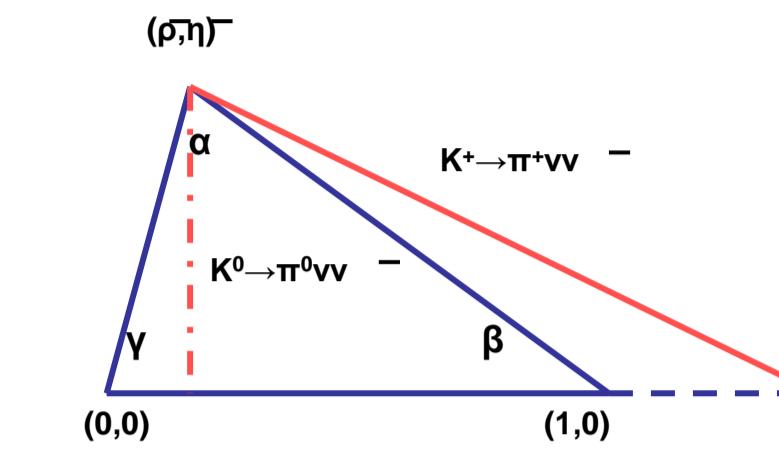


- $K^+ \rightarrow \pi^+ \bar{v} \bar{v}$  process is unique as extremely accurate and clean probe for the non-trivial flavour structure of physics beyond the Standard Model



$$BR(K^+ \rightarrow \pi^+ \bar{v} \bar{v}) = (8.22 \pm 0.84) \cdot 10^{-11}$$

- $K^+ \rightarrow \pi^+ \bar{v} \bar{v}$  is computed in a clean theoretical environment due to the small contribution by hadronic matrix elements and long distance terms.



- FCNC process forbidden at tree level (GIM mechanism)
- top contribution is dominant in loops: cleanest way to extract  $V_{td}$  and to give independent determination of the unitarity triangle

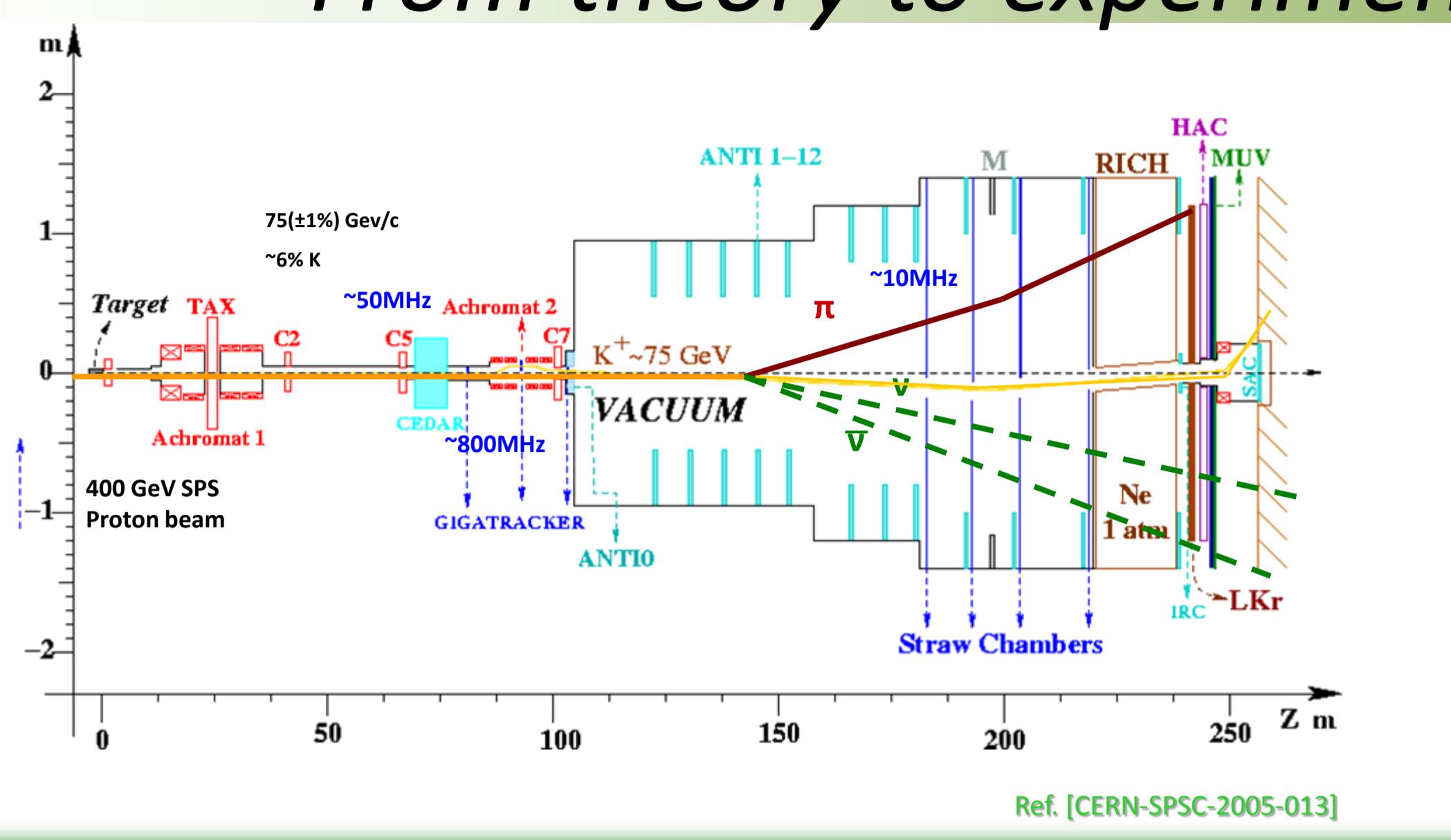
- NA62 aims at measuring  $K^+ \rightarrow \pi^+ \bar{v} \bar{v}$  with  $O(100)$  SM events, in two years of data taking

- Very challenging experiment: veto capability on not-interesting decays and high beam intensity
- tracking of single beam and decay particles, redundant time, momentum, angle measurements, high vetoing power for charged and neutral particles, highly effective particle identification, Ultra high rate capability...

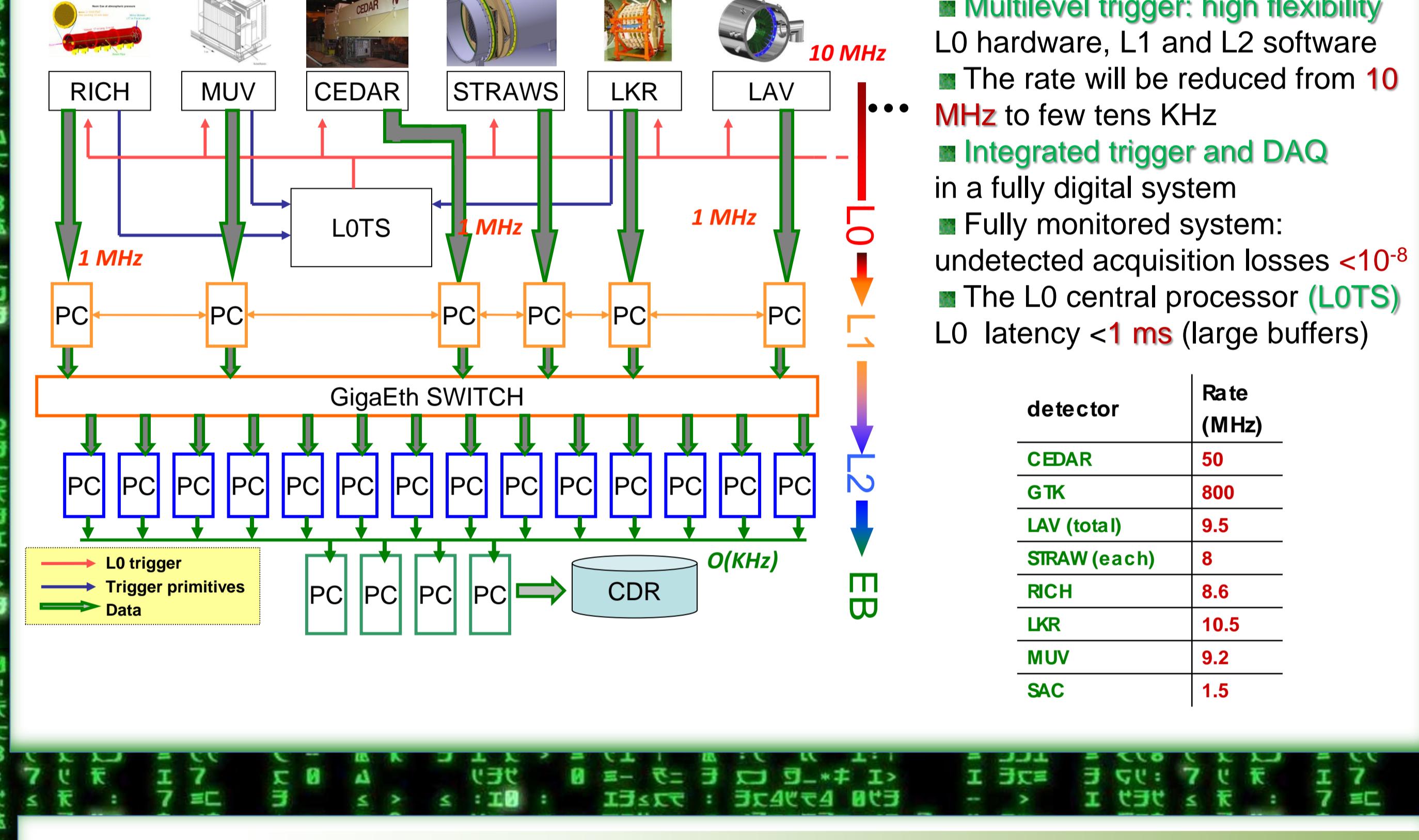
## Main backgrounds

$BR(K^+ \rightarrow \pi^+ \pi^0)$	$20.7 \cdot 10^{-2}$
$BR(K^+ \rightarrow \mu^+ \nu)$	$63.5 \cdot 10^{-2}$
$BR(K^+ \rightarrow e^+ \pi^0 \nu)$	$5.1 \cdot 10^{-2}$
$BR(K^+ \rightarrow \mu^+ \pi^0 \nu)$	$3.4 \cdot 10^{-2}$

## From theory to experiment

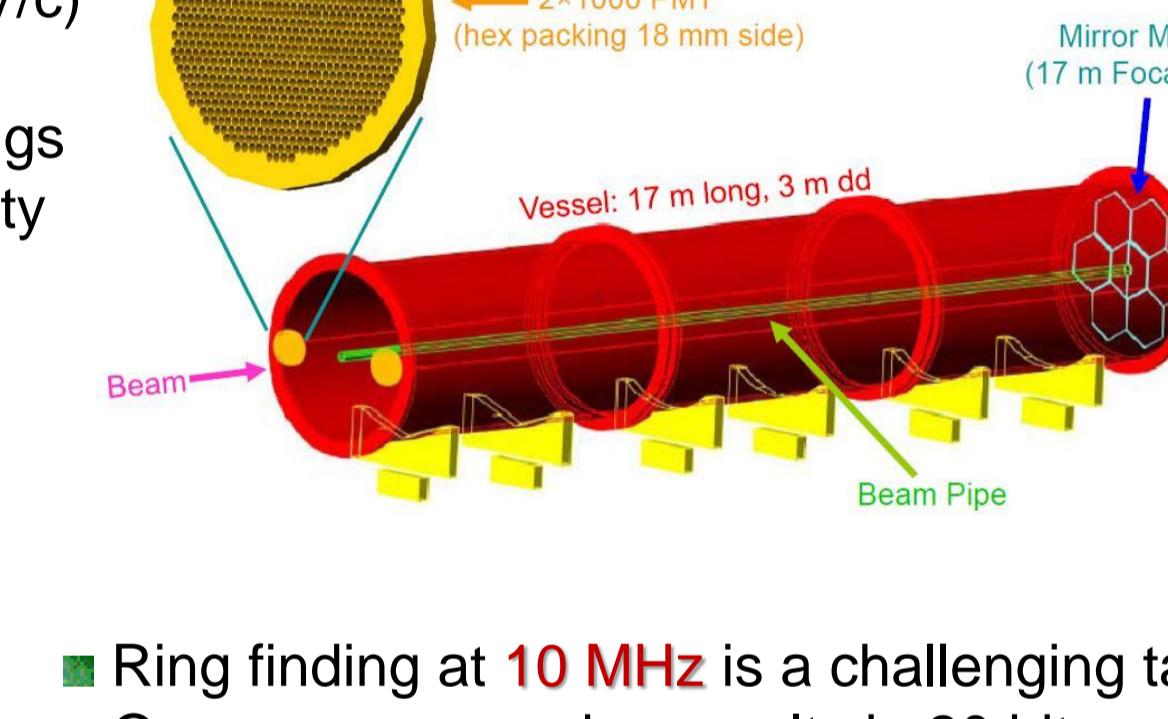


## The NA62 trigger



## RICH & Rings

- RICH counter ( $Ne$ , 1 atm) for particle identification downstream the decay region.
- separation  $\pi^+ - \mu^+$  at  $5 \times 10^{-3}$  level (range 15-35 GeV/c)
- Good time resolution: below 100 ps
- $\sim 2000$  photomultipliers to identify Cherenkov rings
- Participate to L0 trigger definition with multiplicity (by default)



- Rings at trigger level  $\rightarrow$  selective conditions
- Ring's radius and center  $\rightarrow$  velocity and angle
- $\rightarrow$  coarse on-line measurement of particle momentum (with assumption on particle type)

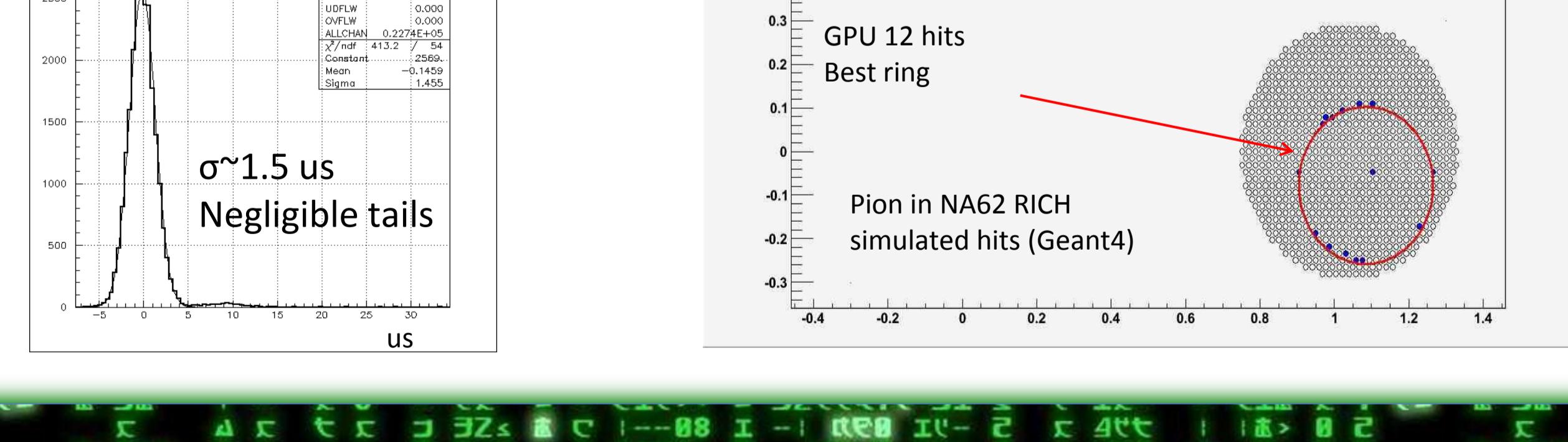
- Ring finding at **10 MHz** is a challenging task
- On average one pion results in 20 hits  $\rightarrow 1.6 \text{ GB/s}$  (160 B/event)

## Fast Ring finder

- Simplest parallelized algorithm: one process for each possible center (each photomultiplier)
- Without optimizations one ring found in  $< 150 \mu\text{s}$  (stable timing results with negligible tails)
- On track for use in  $< 1\text{ms}$  latency on-line system

- Goal:  $100 \mu\text{s}$  per ring
- best algorithm (under test): **Generalized Hough Transform**

- System used for tests: PC with
  - CPU Intel i7-950 + 12GB RAM,
  - GPU: 2 x NVIDIA TESLA C1060 i.e. 2 x {240 cores + 4 GB VRAM}

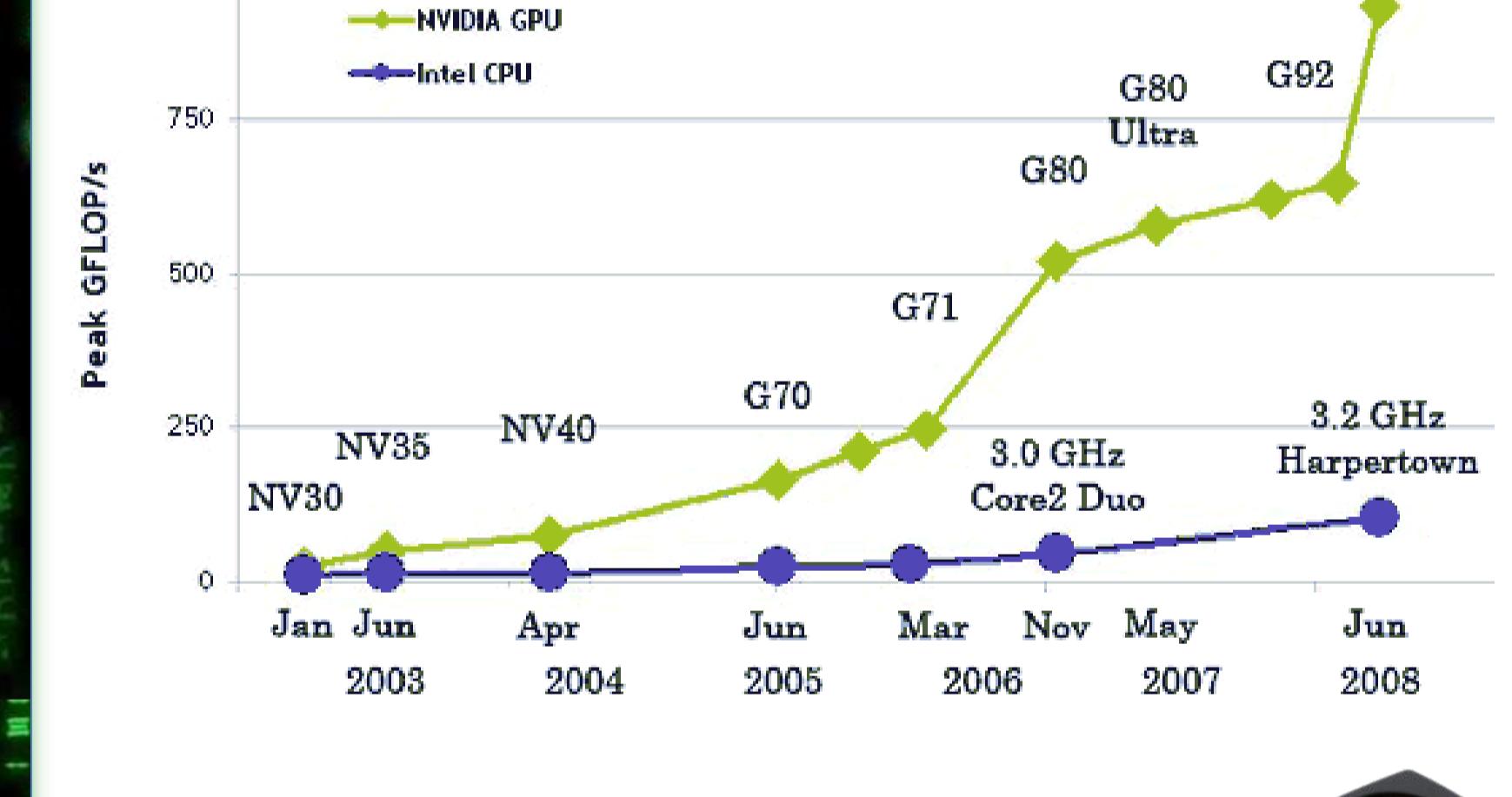


- Simpler scheme for GPU **cooperation** with L0TS: links through PCI-e lanes or better SLI fast links

- Other possibilities in case of effective GPU performance not respecting the L0 requirements:
  - L0.5:** preselection by L0TS before sending to GPU
  - L1:** for decreasing the dimension of the L1 farm

- Many of the detectors can participate to lowest trigger level thanks to the GPU approach (spectrometer, calorimeter, etc...)

## GPU (Graphics Processor Unit)



- Fast expanding field of scientific computation on GPUs (GPGPU)

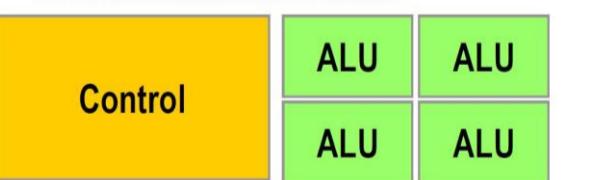
- This project: attempt to exploit GPUs in moderately large REAL-TIME environment

- So far large latencies restricted this to high-performance computing; now  $O(100\mu\text{s})$  latencies easily obtained

- Easy and cheap on a workstation
  - a few TFlops  $\rightarrow$  fast answers
  - hundreds GB/s  $\rightarrow$  high data rates

- Fast links and real-time Linux for appropriate connection with the external world

- The Graphics Cards sector is one of the most supported by IT industry (Digital imaging, video processing, Video Games,...): can profit for free!



- GPU = graphics processor in commercial Graphics Cards

- GPUs power increasing faster than the CPUs due to:

- Different architecture:** more transistors for processing rather than caching, resource usage prediction and flow control
- Different scalability:** the GPU power is related to the dimension and numbers of chips (more and more cores) and not to the chip speed (as in the CPU)

- GPU operate simultaneous processing on thousands points: **large parallel computing power, large memory, huge bandwidths**

- Examples:

- ATI-RADEON 4870:** 10cores x16processors x5sub-units  $\sim 1\text{TFlops}$ ,  $\sim 100 \text{ GB/s}$  bandwidth

- NVIDIA-TESLA C1060:** 30 cores x 8 processors + 4GB RAM  $\sim 1\text{TFlops}$ ,  $\sim 100 \text{ GB/s}$  bandwidth

GPUs	1 Tesla GPU
Single Precision Performance	933 Gigaflops
Double Precision Performance	78 Gigaflops
Memory	4 GB DDR3
Memory speed	800 MHz
Bandwidth	102 GB/s

## GPU integration in NA62 trigger

- Measurements with high resolution at L0 = great opportunity for efficient trigger in ultra high rate environment
- Total trigger bandwidth for main physics channels (no zero-suppression) can be reduced by selective conditions early in trigger chain
- TDAQ will not be a bottleneck in case of higher beam intensity for collecting more statistics
- The scheme will give the unique possibility to collect additional interesting decays channels otherwise suppressed by the trigger for the main channels
- Computing power pressure on the high software trigger levels highly reduced

Gianmario.collazuol@cern.ch, gianluca.lamanna@cern.ch, marco.sozzi@cern.ch