



CUDA in NA62

Luca Pontisso (INFN) on behalf of GAP and NaNet collaborations

Workshop di CCR: La Biodola 2016





- Quick glance to GPU
- CUDA: a parallel computing platform/programming model
- Physics case: GPU based L0 trigger for NA62 RICH detector
- Algorithms implemented
 - Histogram
 - Almagest
- Concluding remarks



Graphics Processing Unit (GPU)

INFN

- GPU's advanced capabilities were originally used primarily for 3D game rendering
- Since 2007 high-Level programming languages (CUDA, OpenCL) have been introduced
- Now this devices are largely deployed in General Purpose applications (GPGPU)
- Based on a massively parallel architecture
- Thousands of cores to process parallel workloads efficiently
- Faster evolution with respect to traditional CPU
- Easy to have a desktop PC with teraflops of computing power, with thousands of cores.





CUDA



- A general purpose parallel computing platform and programming model
- It allows developers to use C as a high-level programming language
- Supporting various languages and application programming interfaces



• CUDA extends C defining functions (kernels) that, when called, are executed N times in parallel by N different CUDA threads,



Memory hierarchy





Physics case: Low Level Trigger system in NA62



- Measurement of ultra-rare decay $K^+ \rightarrow \pi^+ v \overline{v}$ (BR ~8x10⁻¹¹)
- Kaon decays in flight
- High intensity unseparated hadron beam (6% kaons)
- L0 trigger: synchronous level must reduce rate from 10 MHz to 1 MHz
 - 1 ms max. latency



NA62 DAQ and Trigger





L. Pontisso – INFN

WORKSHOP DI CCR 2016

20/05/2016 7

GPU-based L0 trigger general scheme



NA62 Rich detector

INFN

- Distinguish between pions and muons from 15 to 35 GeV (inefficiency < 1%)
- 2 spots of 1000 PMs each
- 2 read-out boards for each







Latency communication problem





Total latency dominated by double copy in Host RAM:

- Data are copied from kernel buffer to destination buffer in user space
- Data are copied from CPU memory to GPU memory

How to reduce data transfer time:

- DMA (Direct Memory Access) from the network channel directly in GPU memory
- Custom management of NIC buffers



NaNet: a PCIe NIC family for HEP



OBJECTIVES:

- Bridging the front-end electronics and the software trigger computing nodes.
- Supporting multiple link, multiple network protocols.
- Low and stable communication latency.
- Having a high bandwidth.
- Processing data streams from detectors on the fly.
- <u>Optimizing data transfers with</u> <u>GPU accelerators.</u>







NaNet-10 at CERN

NaNet: GPUDirect and Software





Host

- Linux Kernel Driver
- User space Library (open/close, buf reg, wait recv evts, ...)

Nios II Microcontroller

 Single process program performing System Configuration & Initialization tasks

- GPUDirect allows direct data exchange on the PCIe bus with no CPU involvement
- No bounce buffers on host memory
- Zero copy I/O
- Buffers on GPU arranged in a circular list of persistent receiving buffers (CLOP)
- nVIDIA Fermi/Kepler/Maxwell



GAP RT

Don't do it in GPU





Merging the events coming from the RICH on GPU... NO WAY

- it requires synchronization and serialization
- computing kernel launched after merging







- Merging stage will be performed in HW
- Events are arranged in CLOPs with a new format more suitable for GPU's threads memory access Multi Merged Event GPU Packet (M²EGP).
 - Problem: searching for events position inside a CLOP using 1 thread on GPU <u>takes > 100us</u> for hundreds of events
 - Solution: it must be parallelized. We can use all the threads looking for a known bytes pattern at the begin of every event: it takes
 ~ 35us for 1000 events in a buffer

	* * * * *		11	55	* * *		+ f		* * * *							
STR 3 MGP	STR 2 MGP	STR 1 MGP	STR 0 MGP	STR 3 HIT STR 2 HIT		STR 1 HIT	STR 0 HIT	T PATTERN		TOTAL HIT		TIMES		ТАМР		
STREAM	1; HIT 1	L STREAM 1; HIT 0		STREAM 0; HIT 5 STREA		STREAM	0; HIT 4 STREAM 0; HIT 3		0; HIT 3	STREAM 0; HIT 2		STREAM 0; HIT 1		STREAM 0; HIT 0		
STREAM 2; HIT 0 STREAM 1; HIT 8		STREAM 1; HIT 7		STREAM 1; HIT 6		STREAM 1; HIT 5		STREAM 1; HIT 4		STREAM 1; HIT 3		STREAM 1; HIT 2				
STREAM	2; HIT 8	STREAM 2; HIT 7		STREAM 2; HIT 6		STREAM 2; HIT 5		STREAM 2; HIT 4		STREAM 2; HIT 3		STREAM 2; HIT 2		STREAM 2; HIT 1		
STREAM 3; HIT 4 ST		STREAM	STREAM 3; HIT 3		STREAM 3; HIT 2		STREAM 3; HIT 1		STREAM 3; HIT 0		STREAM 2; HIT 11		STREAM 2; HIT 10		STREAM 2; HIT 9	
PADDING										STREAM 3; HIT 7		STREAM 3; HIT 6		STREAM 3; HIT 5		
127120	119112	111104	10396	9588	8780	7972	7164	6356	5548	4740	3932	3124	2316	158	70	

2222 2222 2222 2222





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
- Events rate at ~10 MHz
- Low latency
 - **Online (synchronous) trigger**
- Accurate
 - **Offline resolution required**





Histogram: a pattern recognition algorithm



- XY plane divided into a grid
- An histogram is created with distances from these points and hits of the physics event
- Rings are identified looking at distance bins whose contents exceed a threshold value





Pros: naturally mapped on the GPU threads grid Element of the grid <--> thread Event <--> block

Cons: memory limited, performances depending on number of hits



Histogram: a pattern recognition algorithm



First implementation 16x16 grid -> 256 threads x event

_	<u> </u>	-	_	-	_	00	100	00	_	_	_		_	_	_
				000	00	00	000	000	00						
_	-	-		00	0000	0000	00	99	000	100	-	<u> </u>	-	-	-
		°	000		00	00			00		poo	1			
	-		00									00	p O O	h -	_
		1 ~	000	000			mm			b m m	000	mm	000	P I	
	_		000		mm	mm	mm	-	Lan an	mm		hana	00	6	_
	00	000	00	00	000	000	00	00		DODO	00	00	000	5	
	00	00	000	000	00	00	000	000	00	00	000	000	00	00	b
	00		00	00	000		00	00		000	00	00		00	
0	00	00			00	00	000		00	00		000	00	00	Þ
	00	page	i can can	(m) (m)	0000	o an a	(m) (m)	a a	000	b an a	66	(m) (m)	D CO C	00	
	00	00		000	00	00			00	00			00	00	P O
	0.0	mm	mm	D CD C	mm	mm	mm	D CD C	mm	mm	0.000	b m m	mm	mm	0
0	0.00		000	തത	000	000	mm	00	000	ama	00	mm		000	00
00	00	000	000	0000	00	00	000	000	000	00	000	0000	000	000	00
0	0 00		00	00	000	000	00	00	000	000	00	00	000	000	00
-0-(00	00	000	0.000				0.000	00		0.00	0.000	00	100	00
00	000	000	00	00	000	000	00	00	000	000	00	00	000	000	0
00	00	00		000	00	00	000		00	00	000	000	00	00	00
00	000			00			00				00			000	0
00		00				00			00				00	00	Ρ
00	000					mm							000	600	0
	0.0	hma	mm	mm	mma		mm	mm	mm	hana	mm	mm		000	r
	00	00		000	00	00	000		00	00	000	000	00	00	b.
	00	boo	00	00	000	0000	00	00	0.000	000	00	00	000	00	
	00	00	000	0000	00	00	000	000	00	00	000	000	00	00	_
	- 2	poo	00	00	000	000	00	00	000	000	00	00	000	00	
_	-	0.0		0.000				0.000	00			000	0		-
		Poo	00	00		000	00	00		pood	00	00	P.O		
_	- 9	00		0000				0000				000	0	-	-
		1	100	www.	mm	mm		a a a	mm			For o	ľo.		
	-	-		100	- ww	00		1000	- www	La co	000	000	<u>۲</u>		_
					~ ~	00	000	000	0 0	0.0	000	5			
						<u> </u>	0.0	00	0.0						

2-step implementation 8x8 grid -> 64 threads x event 4x4 grid only around maximum



- More blocks run concurrently
- Increased precision



Histogram: a pattern recognition algorithm





Almagest: a new multi-ring algorithm



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"







Almagest: a new multi-ring algorithm



This algorithm exposes two levels of parallelism...



Almagest: a new multi-ring algorithm



L. Pontisso – INFN

WORKSHOP DI CCR 2016

INFN



Almagest: selecting triplets on GPU

INFN



2 warps Shared memory is needed

Sorting time 64 hits coordinates in both x and y



1 thread <-> 2 hit 1 warp Sorting using «shuffle» btw threads No shared memory No __syncthreads()

No multiple kernels





nvcc -m64 --ptxas-options=-v -O2 -arch=sm_35

ptxas info : Compiling entry function '_Z5sortfv' for 'sm_35'
ptxas info : Function properties for _Z5sortfv
0 bytes stack frame, 0 bytes spill stores, 0 bytes spill loads
ptxas info : Used 17 registers 24576 bytes smem, 320 bytes cmem[0]

nvcc -m64 --ptxas-options=-v -O2 -arch=sm_35

ptxas info : Compiling entry function '_Z5sortfv' for 'sm_35'
ptxas info : Function properties for _Z5sortfv
 0 bytes stack frame, 0 bytes spill stores, 0 bytes spill loads
ptxas info : Used 21 registers, 320 bytes cmem[0]

More resources available for the computing stage...



Almagest



1 thread <-> 1 hit 2 warps Every thread checks if Tolomeo's condition is satisfied starting with one of the triplets previously found



Vote function __ballot() Results propagated through the warp



Shared memory and barrier synchronization because of working with 2 warps



Almagest





- Tesla K20c
- **Only computing time**
- <0.5 us per event (multi-rings) for large



64

96 128





- Tested a working solution with the NaNet-1 board during NA62 2015 Run
- Merging events on GPU not feasible
- Multi-ring algorithms such as Histogram and Almagest are implemented on GPU
 - **There is still room for improvement** -> **Optimize computing kernels**
- The GPU-based L0 trigger with the new board NaNet-10 is installed and being currently tested during the 2016 run



R. Ammendola ^(b), A. Biagioni^(a), S. Chiozzi^(d), A. Cotta Ramusino^(d), S. Di Lorenzo^(f,g), R. Fantechi^(f), M. Fiorini^(d,e), O. Frezza^(a), G. Lamanna^(c), F. Lo Cicero^(a), A. Lonardo ^(a), M. Martinelli^(a), I. Neri^(d), P.S. Paolucci^(a), E. Pastorelli^(a), R. Piandani^(f), L. Pontisso^(f), F. Simula^(a), M. Sozzi^(f,g), P. Vicini^(a).

^(a) INFN Sezione di Roma
^(b) INFN Sezione di Roma Tor Vergata
^(c) INFN - Laboratori Nazionali di Frascati
^(d) INFN Sezione di Ferrara
^(e) Università di Ferrara
^(f) INFN Sezione di Pisa
^(g) Università di Pisa



Thank You