Programming Paradigms for GPU devices

> 1-3 March 2017 CINECA (Roma)





Rights & Credits

These slides are CINECA 2014 and are released under the Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) Creative Commons license, version 3.0.

Uses not allowed by the above license need explicit, written permission from the copyright owner. For more information see:

http://creativecommons.org/licenses/by-nc-nd/3.0/

Slides and examples were authored by:

Isabella Baccarelli, Luca Ferraro, Sergio Orlandini



What is a GPU

Graphics Processing Unit

a device equipped with an highly parallel microprocessor (*many-core*) and a private memory with very high bandwidth

 born in response to the growing demand for high definition 3D rendering graphic applications





CPU vs GPU Architectures

- GPU hardware is specialized for problems which can be classified as *intense data-parallel computations*
 - the same set of operations are executed on different data
 - designed such that more transistors are devoted to data processing rather than data caching and flow control





GPU Architecture Scheme

- A tipical GPU architecture consists of
- main global memory
 - high bandwidth
- Streaming Processor
 - grouping independent cores and control units
- each SM unit has
 - many ALU cores
 - instruction scheduler dispatchers
 - a shared memory with very fast access to data





The concurrency revolution

- CPU vendors tended to increase the computational power of single processing unit by increasing the working frequency and adding more higher level control logic and pipelines
- GPU increased the number of processing units, less logic, lowering frequency and dropping down power consumption





GPGPU (General Purpose GPU) and GPU computing

- many applications that process large data sets can use a data-parallel programming model to speed up the computations
- many algorithms outside the field of image rendering are accelerated by data-parallel processing
- In so why not using GPU power for applications out of the graphic domain?
- many attemps where made by brave programmers and researchers in order to force GPU APIs to treat their scientific data (atoms, signals, events, etc) as pixel or vertex to be crunched by GPUs
- not many survived, yet the era of GPGPU computing was just begun ...



GPGPU Programming Tools

- nVIDIA CUDA (Compute Unified Device Architecture)
 - a set of extensions to higher level programming language to use GPU as a coprocessor for heavy parallel task
 - a developer toolkit to compile, debug, profile programs and run them easily in a heterogeneous systems
- OpenCL (Open Computing Language):
 - a standard open-source programming model developed by major brands of hardware manufacters (Apple, Intel, AMD/ATI, nVIDIA).
 - like CUDA, provides extentions to C/C++ and a developer toolkit
 - extensions for specific hardware (GPUs, FPGAs, MICs, etc)
 - it's very low level (verbose) programming
- Accelerator Directives Approach
 - OpenACC
 - OpenMP v4.x accelerator directives
 - you hope your compiler understand what you want, and do a good job
- Library Based:
 - MAGMA, CUDA Libraries, StarPu, ArrayFire, etc



GPGPU Programming Model

- GPU is seen as an auxilirary coprocessor equiped with
 - thousands of cores
 - global memory with high bandwidth
- computational-intensive data-parallel regions of a program can be exectued on the GPU device
 - thousands of threads will be executed on the GPU
 - each thread will insist on a different GPU core
 - each thread can acts on a different data element independently
 - the GPU parallelism is very close to the SPMD paradigm
- the more the working thread, the better are the performances
 - GPU threads are very *light*
 - no penalty is paid in case of *context-switch* (each thread has its own registers)
 - the more the threads, the more the chance to hide memory or computational latencies



CUDA Execution Model

- serial parts of a program, or those with low level of parallelism, keep running on the CPU (host)
- data parallel and computational intensive parts are executed on the GPU (device)
- required data is moved on GPU memory and back to HOST memory





GPU Thread Hierarchy

- In order to compute N elements on the GPU in parallel, at least N concurrent threads must be created on the device
- GPU threads are grouped togheter in *teams* or *blocks* of threads
- Threads belonging to the same block or team can cooperate togheter exchanging data through a shared memory cache area





more on the GPU Execution Model



when a GPU kernel is invoked:

- each thread block is assigned to a SM in a roundrobin mode
 - a maximum number of blocks can be assigned to each SM, depending on hardware generation and on how many resorces each block needs to be executed (registers, shared memory, etc)
 - the runtime system maintains a list of blocks that need to execute and assigns new blocks to SMs as they complete the execution of blocks previously assigned to them
 - once a block is assigned to a SM, it remains on that SM until the work for all threads in the block is completed
 - each block execution is independent from the other (no synchronization is possible among them)
 - thread of each block are partitioned into warps of 32 *threads* each, so to map each thread with a unique consecutive thread index in the block, starting from index 0.
- the scheduler select for execution a warp from one of the residing blocks in each SM.
- A warp execute one common instruction at a time
 - each GPU core take care of one thread in the warp
 - fully efficiency when all threads agree on their execution path

11

Trasparent Scalability

- the GPU runtime system can execute blocks in any order relative to each other
- This flexibility enables to execute the same application code on hardware with different numbers of SM



Data movement

- data must be moved from HOST to DEVICE memory in order to be processed by a CUDA kernel
- when data is processed, and no more needed on the GPU, it is transferred back to HOST



outine Applications and function

Connection Scheme of *host/device*



CINEC/

SuperComputing Applications and Impedior

D2H and H2D Data Transfers

- GPU devices are connected to the host with a PCIe bus
 - PCIe bus is characterized by very low latency, but also by a low bandwidth with respect to other bus

Technology	Peak Bandwidth
PCIex GEN2 (16x, full duplex)	8 GB/s (peak)
PCIex GEN3 (16x, full duplex)	16 GB/s (peak)
DDR3 (full duplex)	26 GB/s (single channel)

- data transfer can easily become a bottleneck in heterogeneous environment equipped with accelerators
 - strive to minimize transfers or execute them in overlap with computations (advanced technique, more on this later)



- nVIDIA GPU HPC Architectures
 - FERMI, KEPLER and PASCAL generation
 - computing capability
 - more on GPU Execution Model





NVIDIA Architectures naming

- Mainstream & laptops: GeForce
 - Target: videogames and multi-media

- Workstation: Quadro
 - Target: professional graphic applications such as CAD, modeling 3D, animation and visual effects

- GPGPU: Tesla
 - Target: High Performance Computing



NVIDIA Fermi Architecture (2009)

- 16 Streaming Multiprocessors (SM)
- 4-6 GB global memory with ECC
- first model with a cache hierarchy:
 - L1 (16-48KB) per SM
 - L2 (768KB) shared among all SM
- 2 independent controllers for data transfer from/to host through PCI-Express
- Global thread scheduler
 (GigaThread global scheduler) which manage and distribute thread blocks to be processed on SM resources





Fermi Streaming Multiprocessor (SM)

Streaming Multiprocessor sports:

- 32/48 CUDA cores with an arithmetic logic unit (ALU) and a floating point unit (FPU) fully pipelined
- floating point operations are fully IEEE 754-2008 a 32-bit e a 64-bit
 - **fused multiply-add** (FMA) for both single and double precision
- 32768 registers (32-bit)
- 64KB configurable L1 shared-memory/cache
 - 48-16KB or 16-48KB shared/L1 cache
- 16 load/store units
- 4 Special Function Unit (SFU) to handle trascendental mathematical functions (sin, sqrt, recp-sqrt,..)







NVIDIA Kepler Architecture (2012)

- x3 performance/watt with respect to FERMI
 - 28nm litography
- 192 CUDA cores
- 4 warp scheduler (2 dispatcher)
 - 2 independent instruction/warp
- standard IEEE 754-2008
- 65536 registers per SM (32-bit)
- 32 load/store units
- 32 Special Function Unit
- 1534KB L2 cache (x2 vs Fermi)
- 64KB shared-memory/cache + 48KB read-only L1 cache
- 16 texture units (x4 vs Fermi)

MX	X								SM													
	Verlex Patch Tessellator Viewport Tessellator									-1			. In	structio	in Cach							
	Alteria Series Steam Output											Warp Scheduler Warp Scheduler					ler -					
	Individual Cache									Dispatch Unit Dispatch Unit					HL .							
	Wary Schender									-			-									
	AND											Registe	in File (B	2,708 x	32-brtj							
				-		egitte	r ne (00,000	x 32-0	1									-	-	LOIST	
	Core	Com	Cree	Core	C	LOUT	aru.	G ***	Caste	10++++	Core	-	Cum	LINET		Co	n.	Core	Com	Core	LD/ST	
-	Cere:	Gen	-	C -++	Cere	10.01			Come	Gene	C		Core	1.0007	111						LOIST	aru
	C	cim.	Care	Curry	C	LOW	-		· · · ·	-	-	-	C.me	-	-	Co	110	Core	Gore	Core	LOIST	
		Cian	Gam			LINKE			Code		Cores.		Corre	LINET	111			and the second	and a	and the second	LDIST	
																GO		Con	GOM	Contra .	LD/ST	
							-											Com	Core	Com	LD/ST	
	•	-	-	Core		L.Cower			-	-	Gam	-	Cint	-			al est	and the			LD/ST	
1000	Cere	(in the second s	Q.m	Gere	C. I	LOW	800	Cart	Cirre	8. m	5000	-	C	LIVET	i ru	Ce		Core	Core	Core	LD/ST	
	Gen	Con	Gore	Con	C in	LINET	-	Geek	Cim	Con	Core	C.	Cme	Linky	eru:						L0/61	SFU
-	C	Com	Gùne	Cire	Cert	LOW	airu	Con	Cint	10 mm	Core		Cme	10.01	NU.	Co	1741	Core	Core	Core	LD/BT	
1000	C-re	C.m.	0.00	-	(interest	LOW	-	C.e.s.	6 m	-	Core	-	Cure	LINIT	111	-	_				LO/ST	_
linia	tere	due		-		LINE		i.e.			Lon	-		LOW	iru.	Co	ire:	Core	Core	Core	LD/ST	
anna a		Citere 1	1	Cores.		HERE			-		Corre		-	LINE	100	-	-			_	LUIST	SFU
														Sec. 1		Co	100	Core	Gore	Core	LOST	
	Cont.	C.C.C.	C. C. C.	etere.		LISTI	anu.		C IIII	-	C.C.C.			vinar,	-			-			-	
	Care	film.	Q.m.	Gove	6	CONT			Citt	Gine	- Const		Corre	UNIT					enconnec			
-	Core	Geo.	Gom	Crew	C.e.e	10au			See.	C	Core		Sine	LINET	H.H.	64 KB Shared Memory / L1 Cache						
-	Core	Core	s in the	Com	Com	LOW	490	C and	Cure	6e	C	-	Cure	LO T				_	Uniform	Cache	_	
-	Tettas Certe											Tex		Tex	Тех	1	lex					
	64 RB Rhavell Henney / J. Clarks											Textore	Cache									
74								74	-									P	olyMorp	h Engine	Viewo	100
1		in.		-		1		T		1		1		-			(erte)	Petch	Tesse	lator	Transfe	ərm
	Harrison Bernard							-			Attribut	te Setup	Stream (Dutput								



NVIDIA Pascal Architecture (2016)

- SM composed of two independent blocks
- each block sports:
 - 1 warps x 2 dispatchers
 - 32 ALU SIMD units
 - 16FP64 units
 - 8 Load/Store units
 - 8 SFU units
 - 32768 32bits registers
- each block accesses:
 - 64KB shared memory
 - L1 64KB cache
 - 4 texture units

SM	SM															
							Instructi	on Cache								
		1	nstructi	on Buffe	r.		Instruction Buffer									
	Warp Scheduler								Warp Scheduler							
-	Dispatch Unit Dispatch Unit								Dispatch Unit Dispatch Unit							
	Register File (32,768 x 32-bit)								Regist	er File (3	32,768 ×	32-bit)				
Care	Core	0F Unit	Core	Core	DP Linit	LOIST	SFU	Core	Core	UP. Linit	Core	Core	DP Unit	LD/ST	SFU	
Core	Core	DP Unit	Core	Core	DP Upit	LDIST	SFU	Core	Core	Unit	Core	Core	DP Unit	LDIST	SFU	
Core	Core	0P Unit	Core	Core	DP Unit	LOIST	SFU	Core	Core	00 Unit	Core	Core	DP Unit	LO/ST	SFU	
Core	Core	DP Unit	Core	Core	Unit	LDIST	SFU	Core	Core	DP Unit	Core	Core	Unit	LDIST	SFU	
Core	Core	Unit	Core	Coro	DP Unit	LDIST	SFU	Core	Core	Unit.	Gora	Core	Unit	LDIST	SFU	
Core	Core	0P Unit	Core	Core	Unit	LDIST	SFU	Core	Core	unit:	Cara	Core	Unit	LD/ST	SFU	
Core	Corr	Unit	Core	Core	Unit	LDIST	SFU	Core	Core	unit .	Core	Core	OP Unit	LD/ST	SFU	
Core	Core	Unit	Core	Core	Unit	LDIST	SFU	Core	Core	Unit	Core	Core	Unit	LD/ST	SFU	
	Texture / L1 Cache															
	Tex Tex Tox Tox															
						0	4KB Sha	red Memo	W.							



NVIDIA Pascal Architecture (2016)

- 6 Compute Graphic Clusters (CGC) with 10 SM each
- 16nm litography
 - 2X Watt/Flop respect Kepler architecture
- 4MB L2 cache
- High Bandwidth Memory
 - 16GB RAM
 - 760 GB/s bandwidth
- NVLink tecnology
 - 80GB/s bandwidth to host data transfers
 - 5X respect PCIe Gen3 16x



Peak Performance: 5,7 TFlops



Compute Capability

- the compute capability of a device describes its architecture
 - registers, memory sizes, features and capabilities
- the compute capability is identified by a code like "compute Xy"
 - major number (X): identifies base line chipset architecture
 - minor number (y): indentifies variants and releases of the base line chipset
- a compute capability select the set of usable PTX instructions

compute capability	feature support
compute_20	FERMI architecture
compute_30	KEPLER K10 architecture (only single precision)
compute_35	KEPLER K20, K20X, K40 architectures
compute_37	KEPLER K80 architecture (two K40 on a single board)
compute_53	MAXWELL GM200 architecture (only single precision)
compute_60	PASCAL GP100 architecture



Capability: resources constraints

Technical Specifications		Compute Capability							
rechnical specifications	1.0	1.1	1.2	1.3	2.x	3.0	3.5		
Maximum dimensionality of grid of thread blocks			3						
Maximum x-dimension of a grid of thread blocks			65535			2 ³¹ -1			
Maximum y- or z-dimension of a grid of thread blocks				65535					
Maximum dimensionality of thread block				3					
Maximum x- or y-dimension of a block		5	12			1024			
Maximum z-dimension of a block				64					
Maximum number of threads per block		5	12			1024			
Warp size				32					
Maximum number of resident blocks per multiprocessor			16						
Maximum number of resident warps per multiprocessor	2	24 32 48				64			
Maximum number of resident threads per multiprocessor	7	768 1024 1536			1536	2048			
Number of 32-bit registers per multiprocessor	8	8 K 16 K 32 K				64 K			
Maximum number of 32-bit registers per thread	128					63 255			
Maximum amount of shared memory per multiprocessor	16 KB					48 KB			
Number of shared memory banks	16					32			
Amount of local memory per thread	16 KB					512 KB			
Constant memory size	64 KB								
Cache working set per multiprocessor for constant memory	8 KB								
Cache working set per multiprocessor for texture memory	Device dependent, between 6 KB and 8 KB								
Maximum width for a 1D texture reference bound to a CUDA array	8192					65536			



Warps

- The GPU multiprocessor creates, manages, schedules, and executes threads in groups of 32 parallel threads called *warps*.
- Individual threads composing a warp start together at the same program address, but they have their own instruction address counter and register state and are therefore free to branch and execute independently
- each warp can execute instructions on
 - SM cores
 - load/store units
 - SFUs units





Hiding Latencies

- What is latency?
 - the number of clock cycles needed to complete an istruction
 - ... that is, the number of cycles I need to wait for before another dependent operation can start
 - arithmetic latency (~ 18-24 cycles)
 - memory access latency (~ 400-800 cycles)
- We cannot discard latencies (it's an hardware design effect), but we can lesser their effect and hide them.
 - saturating computational pipelines in computational bound problems
 - saturating bandwidth in memory bound problems
- We can organize our code so to provide the scheduler a sufficient number of independent operations, so that the more the warp are available, the more context-switch can hide latencies and proceed with other useful operations
- There are two possible ways and paradigms to use (can be combined too!)
 - Thread-Level Parallelism (TLP)
 - Instruction-Level Parallelism (ILP)



Thread-Level Parallelism (TLP)

- Strive for high SM occupancy: that is try to provide as much threads per SM as possible, so to easy the scheduler find a warp ready to execute, while the others are still busy
- This kind of approach is effective when there is a low level of independet operations per CUDA kernels





Instruction-Level Parallelism (ILP)

- Strive for multiple independent operations inside you CUDA kernel: that is, let your kernel act on more than one data
- this will grant the scheduler to stay on the same warp and fully load each hardware pipeline

 note: the scheduler will not select a new warp untill there are eligible instructions ready to execute on the current warp

thread

$$w = w + b$$

$$z = z + b$$

$$y = y + b$$

$$x = x + b$$

$$w = w + a$$

$$z = z + a$$

$$y = y + a$$

$$x = x + a$$
4 independent operations



Coalesced Access to GPU Memory

- All load/store requests in global memory are issued per warp (as all other instructions)
 - 1. each *thread* in a *warp* compute the address to access
 - 2. load/store units select segments where data resides
 - 3. load/store start transfer of needed segments
- It is very important to align data in memory so to have aligned accesses (*coalesced*) during load/store operation in global memory, reducing the number of segments moved across the bus

Stid	ed based copy	Offset based copy					
Stride	Bandwidth GB/s	Offset	Bandwidth GB/s				
1	106.6	0	106.6				
2	34.8	1	72.2				
8	7.9	8	78.2				
16	4.9	16	83.4				
32	2.7	32	105.7				



