

**GPU PROGRAMMING** 

**HPC** courses

UvA, January 2016

Ana Lucia Varbanescu

SHANE COOK

M<

#### UNIVERSITY OF AMSTERDAM

#### Graphics in 1980





#### Graphics in 2000

3

eric ate Ares's rocket Mr Elusive ate Ares's rocket Ares was melted by Willits's plasmagun

×.

#### You fragged Ares 2nd place with 28

res

Rocket Launcher

100 29 28

# Graphics in 2015



#### **GPUs in movies**

#### □ From Ariel in Little Mermaid to Brave





#### So ...

- GPUs are a steady market
  - **Gaming**
  - CAD-like activities
    - Traditional or not ...
  - Visualisation
    - Scientific or not ...
- GPUs are increasingly used for other types of applications
  - Number crunching in science, finance, image processing
  - (fast) Memory operations in big data processing

#### Another GPGPU history



### **TODO** List

- 1. Briefly on performance
- 2. GPGPUs
- 3. CUDA
- 4. If (time\_left && vote) advanced CUDA

else

talk more about performance

# Performance [1]

#### Latency/delay

- The time for one operation (instruction) to finish, L
- To improve: minimize L
  - Lower is better
- Throughput
  - The number of operations (instructions) per time unit, T
  - To improve: maximize T
    - Higher is better
    - Thus, time per instruction decreases, on average
- Example: 1 man builds a house in 10 days.
  - Latency improvement: ...
  - Throughput improvement: ...

# Performance [2]

□ How do we get faster computers?

- Faster processors and memory
  - Increase clock frequency  $\rightarrow$  latency boost
- Better memory techniques
  - Use memory hierarchies  $\rightarrow$  latency boost
  - More memory closer to processor  $\rightarrow$  latency boost
- Better processing techniques
  - Use pipelining  $\rightarrow$  throughput boost
- More processing units (cores, threads, ...)
- Accelerators
  - Use specialized functional units  $\rightarrow$  latency+throughput boost

# -1 Why multi- and many-cores?

Multi-cores = processors with multiple,

homogeneous cores

Many-cores = GPUs & alikes

# Moore's Law

- 12
- Gordon Moore (co-founder of Intel) predicted in 1965 that the transistor density of semiconductor chips would double roughly every 18 months.



"The complexity for minimum component costs has increased at a rate of roughly a factor of two per year ... Certainly over the short term this rate can be expected to continue, if not to increase...." Electronics Magazine 1965

#### **Transistor Counts**

13



### Traditionally ...

- 14
- □ More transistors = more functionality
- □ Improved technology = faster clocks = more speed

Thus, every 18 months, we expect better and faster processors.

They were all sequential: they execute one operation per clock cycle.

#### **Revolution in Processors**





#### New ways to use transistors

- 16
- Parallelism on-chip: multi-core processors
  - Transformed in many-core processors.
- "Multicore revolution"
  - Every machine is a parallel machine.
  - Accelerators start playing an important role.
    - Specialized
    - Energy efficient
    - Used on demand
- Can all applications use this parallelism?
- Can we program all these architectures efficiently?
  Performance? Productivity?

### GPU vs. CPU performance

Theoretical GFLOP/s

 $1 \text{ GFLOPs} = 10^9 \text{ ops} / \text{second}$ 



#### GPU vs. CPU performance

Theoretical GB/s

 $1 \text{ GB/s} = 8 \times 10^{9} \text{ bits / second}$ 



Why do we use many-cores?

- 19
- Performance
  - Large scale parallelism
- Power Efficiency
  - Use transistors more efficiently
- Price (GPUs)
  - Game market is huge, bigger than Hollywood
    - Gaming pays for our HPC needs!
  - Mass production, economy of scale
- Prestige
  - Reach ExaFLOP by 2019/2022 ...



GPUs = the hardware GPGPU = general purpose GPU (typically related to software/programming)

#### **GPGPU** History



1995

2000

2005

2010

- Current generation: NVIDIA Kepler
  - 7.1B transistors
  - More cores, more parallelism, more performance

# GPGPU @ NVIDIA



GPUs @ AMD

#### **AMD Radeon Graphics Roadmap**



GPUs @ ARM



# (NVIDIA) GPUs

- Architecture
  - Many (100s) slim cores
  - Sets of (32 or 192) cores grouped into "multiprocessors" with shared memory
    - SM(X) = stream multiprocessors
  - Work as accelerators
- Memory
  - Shared L2 cache
  - Per-core caches + shared caches
  - Off-chip global memory
- Programming
  - Symmetric multi-threading
  - Hardware scheduler

#### NVIDIA's GPU Architecture

26



### Parallelism

#### Data parallelism (fine-grain)

Restricted forms of task parallelism possible with newest generation of NVIDIA GPUs

#### SIMT (Single Instruction Multiple Thread) execution

- Many threads execute concurrently
  - Same instruction
  - Different data elements
  - HW automatically handles divergence
- Not same as SIMD because of multiple register sets, addresses, and flow paths\*
- Hardware multithreading
  - HW resource allocation & thread scheduling
    - Excess of threads to hide latency
    - Context switching is (basically) free

\*http://yosefk.com/blog/simd-simt-smt-parallelism-in-nvidia-gpus.html

#### Integration into host system

- Typically PCI Express 2.0
- Theoretical speed 8 GB/s
  - **Effective**  $\leq 6 \text{ GB/s}$
  - □ In reality: 4 6 GB/s
- V3.0 recently available
  Double bandwidth
  Less protocol overhead







### CPU vs. GPU







#### Why so different?

- Different goals produce different designs!
  - CPU must be good at everything
  - GPUs focus on massive parallelism
    - Less flexible, more specialized
- CPU: minimize latency experienced by 1 thread
  - big on-chip caches
  - sophisticated control logic
- □ GPU: maximize throughput of all threads
  - # threads in flight limited by resources => lots of resources (registers, etc.)
  - multithreading can hide latency => no big caches
  - share control logic across many threads

# CPU vs. GPU

- 31
- 🗆 Movie
- The Mythbusters
  - Jamie Hyneman & Adam Savage
  - Discovery Channel
- Appearance at NVIDIA's NVISION 2008



# <sup>32</sup> GPU Hardware: NVIDIA



# Fermi

- Consumer: GTX 480, 580
- HPC: Tesla C2050
  - More memory, ECC
  - 1.0 Tlop SP
  - 515 GFlop SP
- 16 streaming multiprocessors (SM)
  - **GTX 580: 16**
  - **GTX 480: 15**
  - **C2050:** 14
- SMs are independent768 KB L2 cache



#### Fermi Streaming Multiprocessor (SM)



### **CUDA Core Architecture**

- Decoupled floating point and integer data paths
- Double precision throughput
  is 50% of single precision
- Integer operations optimized for extended precision
  - 64 bit and wider data element size
- Predication field for all instructions
- Fused-multiply-add

	SM						
	Instruction Cache						
	Warp Scheduler						
	Dispatch Unit				Dispatch Unit		
070	-						
Port	Register File (32,768 x 32-bit)						
NT Unit eue							
	Core	Co					
	Core	Core					
	Core						
	Core						
	Core						
	Core						
	Core						
	Core						
	Interconnect Network						
	64 KB Shared Memory / L1 Cache						
	Uniform Cache						
	Tex						
	Texture Cache						
	PolyMorph Engine						
	Vertex Fetch Tessellator Transform						
	Attribute Setup Stream Output						

CUDA

# Memory architecture (since Fermi)

36

- Configurable L1 cache per SM
  - 16KB L1 cache / 48KB Shared
  - 48KB L1 cache / 16KB Shared


#### Kepler: the new SMX

- Consumer:
   GTX680, GTX780, GTX-Titan
- - Tesla K10..K40
- SMX features
  - 192 CUDA cores
    - 32 in Fermi
  - 32 Special Function Units (SFU)
     4 for Fermi
  - 32 Load/Store units (LD/ST)
     16 for Fermi
- 3x Perf/Watt improvement

SMX															
PolyMorph Engine 2.0															
Attribute Seture															
	Instruction Cache														
Dispat	ch Unit	Dispat	ch Unit	Dispatch Unit Dispatch Unit			Dispatch Unit Dispatch Unit			Dispatch Unit Dispatch Unit					
	+ + + + + + + + +														
L.	Register File (65,536 x 32-bit)														
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
C	0.000			0.000	0.000	I DIST	OF11	0.000					C	I D/PT	0511
Core	Core	Core	Core	Core	Core	LUIST	aru	Core	Core	Core	Core	Core	Core	LU/31	350
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	Core	Core	LD/ST	SFU
	Texture Cache														
	64 KB Shared Memory / L1 Cache														
	Uniform Cache														
Tex Tex Tex Tex Tex Tex Te							эх								
Te	Tex Tex Tex Tex Tex Tex Tex Tex														
•	Interconnect Network														

### A comparison

38

	FERMI GF100	FERMI GF104	KEPLER GK104	KEPLER GK110
Compute Capability	2.0	2.1	3.0	3.5
Threads / Warp	32	32	32	32
Max Warps / Multiprocessor	48	48	64	64
Max Threads / Multiprocessor	1536	1536	2048	2048
Max Thread Blocks / Multiprocessor	8	8	16	16
32-bit Registers / Multiprocessor	32768	32768	65536	65536
Max Registers / Thread	63	63	63	255
Max Threads / Thread Block	1024	1024	1024	1024
Shared Memory Size Configurations (bytes)	16K	16K	16K	16K
	48K	48K	32K	32K
			48K	48K
Max X Grid Dimension	2 <b>^</b> 16-1	2^16-1	2^32-1	2^32-1
Hyper-Q	No	No	No	Yes
Dynamic Parallelism	No	No	No	Yes

#### Maxwell: the newest SMM

- Consumer:
   GTX 970, GTX 980, ...
- - □ Ś
- SMM Features:
  - 4 subblocks of 32 cores
  - Dedicated L1/LM per 64 cores
  - Dispatch/ecode/registers per 32 cores
- □ L2 cache: 2MB (~3x vs. Kepler)
- 40 texture units
- Lower power consumption

SMM												
PolyMorph Engine 2.0												
Vertex Fetch Tessellator Viewport Transform												
Attribute Setup Stream Output												
Instruction Cache												
	-	nstructio	on Buffe	er	Instruction Buffer							
Di	spatch Uni	it st	D	)ispatch Ur	nit	Dispatch Unit Dispatch Unit						
	+ Periot	ar Eila //	6 204 -	<b>1</b>		Register File (16 384 x 32-bit)						
	Regist		10,304 X	52-011)								
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU	
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU	
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU	
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU	
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU	
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU	
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU	
Core	Core	Core	Core	LD/ST	SFU	Core	Core	Core	Core	LD/ST	SFU	
Texture / L1 Cache												
					Texture /	L1 Cache						
	Tex			Tex	Texture /	L1 Cache	Төх			Tex	_	
	Tex	nstructi	on Buffe	Tex er	Texture /	L1 Cache	Төх	nstructi	on Buffe	Tex er		
	Tex I	nstructi Warp Si	on Buffe	Tex er	nit	L1 Cache	Tex I	nstructi Warp Se	on Buffe	Tex er	nit	
D	Tex I ispatch Un	nstructi Warp Si it	on Buffe cheduler	Tex er Dispatch U	nit	D	Tex I ispatch Un	nstructi Warp So it	on Buffe cheduler C	Tex er Dispatch Un	nit	
D	Tex I Ispatch Un I Regist	nstructi Warp So it er File (	on Buffe cheduler 16,384 x	Tex er Dispatch Ui 32-bit)	nit	D	Tex Ispatch Un Regist	nstructio Warp So it er File ('	on Buffe cheduler C 16,384 x	Tex Pr Dispatch Un 32-bit)	nit	
D	Tex Ispatch Un Regist	nstructi Warp Si it er File ( Core	on Buffe cheduler 16,384 x Core	Tex er Dispatch Ur 32-bit) LD/ST	nit SFU	Core	Tex I ispatch Un Regist Core	nstructi Warp So it er File (' Core	on Buffe cheduler 16,384 x Core	Tex er Dispatch Ur 32-bit) LD/ST	nit SFU	
Core	Tex ispatch Un Regist Core Core	nstructi Warp Si it er File ( Core Core	on Buffe cheduler 16,384 x Core Core	Tex er Dispatch U 32-bit) LD/ST LD/ST	nit SFU SFU	Core	Tex I ispatch Un Regist Core Core	nstructi Warp So it er File (' Core Core	on Buffe cheduler 16,384 x Core Core	Tex er Dispatch Ur 32-bit) LD/ST LD/ST	nit SFU SFU	
Core Core Core	Tex Ispatch Un Regist Core Core Core	it Warp So it er File ( Core Core Core	on Buffe cheduler 16,384 x Core Core Core	Tex ar Dispatch Ur 32-bit) LD/ST LD/ST LD/ST	SFU SFU SFU	Core Core	Tex ispatch Un Regist Core Core Core	nstruction Warp So it er File ( Core Core	on Buffe cheduler 16,384 x Core Core Core	Tex Pispatch Ur 32-bit) LD/ST LD/ST	nit SFU SFU SFU	
Core Core Core	Tex Ispatch Un Regist Core Core Core	warp Si it Core Core Core Core	on Buffe cheduler 16,384 x Core Core Core	Tex Dispatch U/ 32-bit) LD/ST LD/ST LD/ST	SFU SFU SFU SFU	Core Core Core	Tex ispatch Un Regist Core Core Core	it Warp So it Core Core Core Core	on Buffe cheduler 16,384 x Core Core Core	Tex Dispatch Un 32-bit) LD/ST LD/ST LD/ST LD/ST	nit SFU SFU SFU SFU	
Core Core Core Core Core	Tex ispatch Un Regist Core Core Core Core	Warp So it Core Core Core Core Core	on Buffe cheduler 16,384 x Core Core Core Core Core	Tex Dispatch U 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST	nit SFU SFU SFU SFU SFU	Core Core Core Core	Tex ispatch Um Regist Core Core Core Core Core	warp Si it er File (' Core Core Core Core	Core Core Core Core Core Core	Tex Pispatch U/ 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST	nit SFU SFU SFU SFU SFU	
Core Core Core Core Core Core	Tex ispatch Um Regist Core Core Core Core Core	warp Si it Core Core Core Core Core Core	on Buffe cheduler 16,384 x Core Core Core Core Core	Tex Dispatch U/ 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST	SFU SFU SFU SFU SFU SFU	Core Core Core Core Core	Tex ispatch Un Regist Core Core Core Core Core Core	Instruction Warp So it Core Core Core Core Core Core Core	on Buffe heduler 16,384 x Core Core Core Core Core	Tex in ispatch Ur 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST	nit SFU SFU SFU SFU SFU SFU	
Core Core Core Core Core Core Core	Tex Ispatch Um Regist Core Core Core Core Core Core	Warp Si Warp Si it Core Core Core Core Core Core	Core Core Core Core Core Core	Tex Sispatch Unit Dispatch Unit Dispatch Unit LDIST LDIST LDIST LDIST LDIST	SFU SFU SFU SFU SFU SFU	Core Core Core Core Core	Tex ispatch Un Regist Core Core Core Core Core Core	Warp Si Warp Si er File ( Core Core Core Core Core	on Buffe cheduler Core Core Core Core Core Core Core	Tex ar Jispatch Un 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST	nit SFU SFU SFU SFU SFU SFU SFU	
D D Core Core Core Core Core Core	Tex I Ispatch Unit Regist Core Core Core Core Core Core Core	Warp Sa Warp Sa at er File ( Core Core Core Core Core Core	Core Core Core Core Core Core Core	Tex ar Jispatch Ui 2-bit) LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST	SFU SFU SFU SFU SFU SFU SFU SFU SFU	Core Core Core Core Core Core Core	Tex ispatch Um Regist Core Core Core Core Core Core Core Core	warp Si Warp Si it Core Core Core Core Core Core Core	Core Core Core Core Core Core Core Core	Tex ar ar ar ar ar ar ar ar ar ar	sFU SFU SFU SFU SFU SFU SFU SFU SFU	
Core Core Core Core Core Core Core	Tex ispatch Unit Regist Core Core Core Core Core Core Core	Warp Sr R Ter File ( Core Core Core Core Core Core Core Core	Core Core Core Core Core Core Core Core	Tex ar Jispatch U. 4 4 1 1 1 1 1 1 1 1 1 1 1 1 1	SFU SFU SFU SFU SFU SFU SFU SFU SFU SFU	L1 Cache D Core Core Core Core Core Core Core	Tex I I Ispatch Un Regist Core Core Core Core Core Core Core	warp St Warp St er File ( Core Core Core Core Core Core Core	Core Core Core Core Core Core Core Core	Tex sr lispatch Ur 32-bit) LD/ST LD/ST LD/ST LD/ST LD/ST LD/ST	SFU SFU SFU SFU SFU SFU SFU SFU	
Core Core Core Core Core Core Core	Tex I spatch Um Regist Core Core Core Core Core Core Core Core	Warp Sr Reverse File ( Core Core Core Core Core Core Core Core	Core Core Core Core Core Core Core	Tex ar 32-bit) LD:ST LD:ST LD:ST LD:ST LD:ST LD:ST LD:ST LD:ST	SFU SFU SFU SFU SFU SFU SFU SFU SFU	L1 Cache	Tex I Inspatch Um Regist Core Core Core Core Core Core Core Core	warp Sri Warp Sri er File (* Core Core Core Core Core Core	Core Core Core Core Core Core Core Core	Tex ar ar a2-bit) LD:ST LD:ST LD:ST LD:ST LD:ST LD:ST LD:ST LD:ST	nit SFU SFU SFU SFU SFU SFU SFU	



#### Parallelism

- 41
- Threads
  - Independent units of computation
  - Expected to execute in parallel
  - Write once, instantiate many times
- Concurrent execution
  - Threads execute in the same time if there are sufficient resources
- Assume a processor P with 10 cores and an application A with:
  - 10 threads: how long does A take?
  - 20 threads: how long does A take?
  - 33 threads: how long does A take?

#### Parallelism

- **42**
- Synchronization = a thread's execution must depend on other threads
  - Barrier = all threads wait to get to barrier before they continue
  - Shared variables = more threads RD/WR them
    - Locks = threads can use locks to protect the WR sections
  - Atomic operation = operation completed by a single thread at a time
- Thread scheduling = the order in which the threads are executed on the machine
  - User-based: programmer decides
  - OS-based: OS decides (e.g., Linux, Windows)
  - Hardware-based: hardware decides (e.g., GPUs)

### Programming many-cores

- **43**
- = parallel programming:
  - Choose/design algorithm
  - Parallelize algorithm
    - Expose enough layers of parallelism
    - Minimize communication, synchronization, dependencies
    - Overlap computation and communication
  - Implement parallel algorithm
    - Choose parallel programming model
    - (?) Choose many-core platform
  - Tune/optimize application
    - Understand performance bottlenecks & expectations
    - Apply platform specific optimizations
    - (?) Apply application & data specific optimizations



#### CUDA

#### CUDA: Scalable parallel programming

- C/C++ extensions
  - Other wrappers exist
- Straightforward mapping onto hardware
  - Hierarchy of threads (to map to cores)
    - Configurable at logical level
  - Various memory spaces (to map to physical spaces)
    - Usable via variable scopes
- Scale to 1000s of cores & 100,000s of threads
  - GPU threads are lightweight
  - GPUs need 1000s of threads for full utilization

#### **CUDA Model of Parallelism**

- CUDA virtualizes the physical hardware
  - A block is a virtualized streaming multiprocessor
    - threads, shared memory
  - A thread is a virtualized scalar processor
    - registers, PC, state
- Threads are scheduled onto physical hardware without pre-emption
  - threads/blocks launch & run to completion
  - blocks must be independent

#### **CUDA Model of Parallelism**



Software GPU Thread Thread Processor Thread Block Multi--processor ...



...

#### Hierarchy of threads



#### Using CUDA

Two parts of the code:

Device code = GPU code = kernel(s)

- Sequential program
- Write for 1 thread, execute for all
- Host code = CPU code
  - Instantiate grid + run the kernel
  - Memory allocation, management, deallocation
  - C/C++/Java/Python/...
- Host-device communication
  - Explicit / implicit via PCI/e
  - Minimum: data input/output

#### Processing flow



Image courtesy of Wikipedia

#### Grids, Thread Blocks and Threads



#### Kernels and grids

- 52
- $\Box$  Launch kernel (12 x 6 = 72 instances)
- myKernel<<<numBlocks,threadsPerBlock>>>(...);
  - dim3 threadsPerBlock(3,4);
    - threadsPerBlock.x = 3
    - threadsPerBlock.y = 4
    - Each thread:
    - (threadIdx.x, threadIdx.y)
  - dim3 numBlocks(2,3);
    - blockDim.x = 2
    - blockDim.y=3
    - Each block :

(blockIdx.x,blockIdx.y)

Grid								
Thread Block 0, 0	Thread Block 0, 1	Thread Block 0, 2						
$\begin{array}{c c} 1,0 \\ \downarrow \end{array} \begin{array}{c} 1,1 \\ \downarrow \end{array} \begin{array}{c} 1,2 \\ \downarrow \end{array} \begin{array}{c} 2,3 \\ \downarrow \end{array}$		$\begin{array}{c c} 1,0 \\ \bullet \end{array} \begin{array}{c} 1,1 \\ \bullet \end{array} \begin{array}{c} 1,2 \\ \bullet \end{array} \begin{array}{c} 2,3 \\ \bullet \end{array}$						
$\begin{array}{c c} 2,0 \\ \downarrow \\ \bullet \end{array} \begin{array}{c} 2,1 \\ \bullet \end{array} \begin{array}{c} 2,2 \\ \bullet \end{array} \begin{array}{c} 2,3 \\ \bullet \end{array}$		$\begin{array}{c c} 2,0 \\ \bullet \end{array} \begin{array}{c} 2,1 \\ \bullet \end{array} \begin{array}{c} 2,2 \\ \bullet \end{array} \begin{array}{c} 2,3 \\ \bullet \end{array} \end{array}$						
Thread Block 1, 0	Thread Block 1, 1	Thread Block 1, 2						
0,0 0,1 0,2 0,3 ♥ ♥ ♥ ♥	0,0 0,1 0,2 0,3	0,0 0,1 0,2 0,3						
$\begin{array}{c c} 1,0 \\ \downarrow \end{array} \begin{array}{c} 1,1 \\ \downarrow \end{array} \begin{array}{c} 1,2 \\ \downarrow \end{array} \begin{array}{c} 2,3 \\ \downarrow \end{array}$		$\begin{array}{c c} 1,0 \\ \bullet \end{array} \begin{array}{c} 1,1 \\ \bullet \end{array} \begin{array}{c} 1,2 \\ \bullet \end{array} \begin{array}{c} 2,3 \\ \bullet \end{array}$						
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{c} 2_{1}0 \\ \clubsuit \end{array} \begin{array}{c} 2_{1}1 \\ \clubsuit \end{array} \begin{array}{c} 2_{2}2 \\ \clubsuit \end{array} \begin{array}{c} 2_{3}2 \\ \clubsuit \end{array}$	2,0 2,1 2,2 2,3 ↓ ↓ ↓						

## Multiple Device Memory Scopes

- 53
- Per-thread private memory
  - Each thread has its own local memory
  - Stacks, other private data, registers
- Per-SM shared memory
  - Small memory close to the processor, low latency
- Device memory
  - GPU frame buffer
  - Can be accessed by any thread in any SM



#### Memory Spaces in CUDA



### **Device Memory**

- CPU and GPU have separate memory spaces
  - Data is moved across PCI-e bus
  - Use functions to allocate/set/copy memory on GPU
  - Very similar to corresponding C functions
- Pointers are just addresses
  - Can't tell from the pointer value whether the address is on CPU or GPU
  - Must exercise care when dereferencing:
    - Dereferencing CPU pointer on GPU will likely crash
    - Same for vice versa

#### **Additional memories**

- - Read-only
  - Data resides in device memory
  - Different read path, includes specialized caches
- Constant memory
  - Data resides in device memory
  - Manually managed
  - Small (e.g., 64KB)
  - Assumes all threads in a block read the same addresses
    - Serializes otherwise

#### C for CUDA

57

Philosophy: provide minimal set of extensions necessary

Function qualifiers:
 \_\_global\_\_\_ void my\_kernel() { }
 \_\_device\_\_ float my\_device\_func() { }

Execution configuration: dim3 gridDim(100, 50); // 5000 thread blocks dim3 blockDim(4, 8, 8); // 256 threads per block (1.3M total) my\_kernel <<< gridDim, blockDim >>> (...); // Launch kernel

Built-in variables and functions valid in device code:

dim3 gridDim; // Grid dimension
dim3 blockDim; // Block dimension
dim3 blockIdx; // Block index
dim3 threadIdx; // Thread index

void syncthreads(); // Thread synchronization



# First CUDA program

- Determine mapping of operations and data to threads
- Write kernel(s)
  - Sequential code
  - Written per-thread
- Determine block geometry
  - Threads per block, blocks per grid
  - Number of grids (>= number of kernels)
- Write host code
  - Memory initialization and copying to device
  - Kernel(s) launch(es)
  - Results copying to host
- Optimize the kernels

#### Vector add: sequential

```
void vector_add(int size, float* a, float* b, float* c) {
   for(int i=0; i<size; i++) {
      c[i] = a[i] + b[i];
   }</pre>
```

#### How do we parallelize this?

- What does each thread compute?
  - One addition per thread
  - Each thread deals with \*different\* elements
  - How do we know which element?
    - Compute a mapping of the grid to the data
      - Any mapping will do!

#### Processing flow



Vector add: Kernel

}

```
// compute vector sum c = a + b
// each thread performs one pair-wise addition
__global___ void vector_add(float* A, float* B, float* C) {
    int i = ?
    C[i] = A[i] + B[i];
```

### Calculating the global thread index



"global" thread index:
 blockDim.x \* blockIdx.x + threadIdx.x;

### Calculating the global thread index



"global" thread index:

blockDim.x \* blockIdx.x + threadIdx.x;



Vector add: Kernel

}

#### Done with the kernel!

#### Processing flow



#### Vector add: Launch kernel

**68** 

```
// compute vector sum c = a + b
// each thread performs one pair-wise addition
__global___ void vector_add(float* A, float* B, float* C) {
    int i = threadIdx.x + blockDim.x * blockIdx.x;
    C[i] = A[i] + B[i];
} GPU code
```

```
int main() {
    // initialization code here ...
    N = 5120;
    // launch N/256 blocks of 256 threads each
    vector_add<<< N/256, 256 >>>(deviceA, deviceB, deviceC);
    // cleanup code here ...
}
```

#### Vector add: Launch kernel



#### What if N = 5000?

```
int main() { Host code
   // initialization code here ...
   N = 5000;
   // launch N/256 blocks of 256 threads each
   vector_add<<< N/256, 256 >>>(deviceA, deviceB, deviceC);
   // cleanup code here ...
}
```

#### Vector add: Launch kernel



#### What if N = 5000?

## Memory Allocation / Release

71

- □ Host (CPU) manages device (GPU) memory:
  - cudaMalloc(void \*\*pointer, size\_t nbytes)
  - cudaMemset(void \*pointer, int val, size\_t count)

cudaFree(void\* pointer)

int n = 1024; int nbytes = n \* sizeof(int); int\* data = 0; cudaMalloc(&data, nbytes); cudaMemset(data, 0, nbytes); cudaFree(data);

#### **Data Copies**

 cudaMemcpy(void \*dst, void \*src, size\_t nbytes, enum cudaMemcpyKind direction);
 returns after the copy is complete
 blocks CPU thread until all bytes have been copied
 doesn't start copying until previous CUDA calls complete

#### enum cudaMemcpyKind

- cudaMemcpyHostToDevice
- cudaMemcpyDeviceToHost
- cudaMemcpyDeviceToDevice

Non-blocking copies are also available

DMA transfers, overlap computation and communication
#### Vector add: Host

```
int main(int argc, char** argv) {
  float *hostA, *deviceA, *hostB, *deviceB, *hostC, *deviceC;
  int size = N * sizeof(float);
```

```
// allocate host memory
hostA = malloc(size);
hostB = malloc(size);
hostC = malloc(size);
```

// initialize A, B arrays here...

```
// allocate device memory
cudaMalloc(&deviceA, size);
cudaMalloc(&deviceB, size);
cudaMalloc(&deviceC, size);
```

Vector add: Host

}

// transfer the data from the host to the device cudaMemcpy(deviceA, hostA, size, cudaMemcpyHostToDevice); cudaMemcpy(deviceB, hostB, size, cudaMemcpyHostToDevice);

// launch N/256 blocks of 256 threads each
vector\_add<<<N/256, 256>>>(deviceA, deviceB, deviceC);

// transfer the result back from the GPU to the host
cudaMemcpy(hostC, deviceC, size, cudaMemcpyDeviceToHost);

Done with the host code!

# Summary

- Determine mapping of operations and data to threads
- Write kernel(s)
  - Sequential code
  - Written per-thread
- Determine block geometry
  - Threads per block, blocks per grid
  - Number of grids (>= number of kernels)
- Write host code
  - Memory initialization and copying to device
  - Kernel(s) launch(es)
  - Results copying to host
- Optimize the kernels

#### Practice ?

Let's try this in practice

- Run on DAS4
- You need an ssh client
  - On Linux/Mac: terminal will do
  - On Windows: download putty
    - http://www.chiark.greenend.org.uk/~sgtatham/putty/ download.html
- Use vim to see the files
- Follow the directions in the manual for Assignment 1 and 2

# Practice [cont]

- Vim commands:
  - i enables editing
  - Esc finishes editing
  - :qw exit and save changes
  - :q! exit without saving changes

□ "make" – to compile



Scheduling and synchronization

## **Thread Scheduling**

- 79
- Order in which thread blocks are scheduled is undefined!
  - any possible interleaving of blocks should be valid
  - presumed to run to completion without preemption
  - can run in any order
  - can run concurrently OR sequentially

Order of threads within a block is also undefined!

Q: How do we do global synchronization with these scheduling semantics?

- 81
- Q: How do we do global synchronization with these scheduling semantics?
- □ A1: Not possible!

- 82
- Q: How do we do global synchronization with these scheduling semantics?
- □ A1: Not possible!
- □ A2: Finish a grid, and start a new one!

- 83
- Q: How do we do global synchronization with these scheduling semantics?
- □ A1: Not possible!
- □ A2: Finish a grid, and start a new one!

```
step1<<<grid1,blk1>>>(...);
```

// CUDA ensures that all writes from step1 are complete.
step2<<<grid2,blk2>>>(...);

We don't have to copy the data back and forth!

#### **Atomics**

- 84
- Guarantee that only a single thread has access to a piece of memory during an operation
  - No loss of data
  - Ordering is still arbitrary
- Different types of atomic instructions
  - Add, Sub, Exch, Min, Max, Inc, Dec, CAS, And, Or, Xor
  - On device memory and/or shared memory
- Much more expensive than load + operation + store

#### **Example: Histogram**

- 85
- // Determine frequency of colors in a picture.
- // Colors have already been converted into integers
- // between 0 and 255.
- // Each thread looks at one pixel,
- // and increments a counter

```
global void histogram(int* colors, int* buckets)
```

```
int i = threadIdx.x + blockDim.x * blockIdx.x;
int c = colors[i];
buckets[c] += 1;
```

#### **Example: Histogram**



// Determine frequency of colors in a picture. **by**ors have already been converted into integers between 0 and 255. Each thread looks at one pixel, // and increments a counter global void historra (int\* colors, int\* buckets) int i = threadIdx.x + blockpim.x \* blockIdx.x; Orrec/ int c = colors[i]; buckets[c] += 1;

#### **Example: Histogram**

87

- // Determine frequency of colors in a picture.
- // Colors have already been converted into integers
- // between 0 and 255.
- // Each thread looks at one pixel,
- // and increments a counter atomically

\_global\_\_ void histogram(int\* colors, int\* buckets)

```
int i = threadIdx.x + blockDim.x * blockIdx.x;
int c = colors[i];
atomicAdd(&buckets[c], 1);
```

## 4 CUDA: optimizing your application

- 1. Occupancy
- 2. Shared Memory
- 3. Coalescing
- 4. Streams
- 5. Shared Memory Bank Conflicts

# **Thread Scheduling**

- 89
  - SMs implement zero-overhead warp scheduling
    - A warp is a group of 32 threads that runs concurrently on an SM
    - At any time, the number of warps concurrently executed by an SM is limited by its number of cores.
    - Warps whose next instruction has its inputs ready for consumption are eligible for execution
    - Eligible Warps are selected for execution on a prioritized scheduling policy
    - All threads in a warp execute the same instruction when selected



# Stalling warps

- 90
  - □ What happens if all warps are stalled?
     No instruction issued → performance lost

- Most common reason for stalling?
  - Waiting on global memory
- If your code reads global memory every couple of instructions
  - You should try to maximize occupancy

# Occupancy

- What determines occupancy?
- □ Limited resources!
  - Register usage per thread
  - Shared memory per thread block

### Resource Limits (1)



- Pool of registers and shared memory per SM
  - Each thread block grabs registers & shared memory
  - $\square$  If one or the other is fully utilized  $\Longrightarrow$  no more thread blocks

## Resource Limits (2)

- 93
  - Can only have P thread blocks per SM
    - If they're too small, can't fill up the SM
    - Need 128 threads / block on gt200 (4 cycles/instruction)
    - Need 192 threads / block on Fermi (6 cycles/instruction)

Higher occupancy has diminishing returns for hiding latency

#### Hiding Latency with more threads



#### How do you know what you're using?

Use compiler flags to get register and shared memory usage

"nvcc -Xptxas -v"

- Use the NVIDIA Profiler
- Plug those numbers into CUDA Occupancy Calculator

Maximize occupancy for improved performance
 Empirical rule! Don't overuse!



#### **Thread divergence**

```
"I heard GPU branching is expensive. Is this true?"
 global void Divergence(float* dst,float* src )
{
    float value = 0.0f;
    if ( threadIdx.x \% 2 == 0 )
// active threads : 50%
        value = src[0] + 5.0f;
    else
// active threads : 50%
        value = src[0] - 5.0f;
    dst[index] = value;
```

#### Execution



Worst case performance loss:

50% compared with the non divergent case.

#### Another example





(assume logic below is to be executed for each element in input array 'A', producing output into



### Performance penalty?

- Depends on the amount of divergence
  - Worst case: 1/32 performance
    - When each thread does something different
- Depends on whether branching is data- or IDdependent
  - If ID consider grouping threads differently
  - If data consider sorting
- Non-diverging warps => NO performance penalty
   In this case, branches are not expensive ...

## 4 CUDA: optimizing your application

- 1. Occupancy
- 2. Shared Memory
- 3. Coalescing
- 4. Streams
- 5. Shared Memory Bank Conflicts

## Matrix multiplication example

 $\Box C = A * B$ 

102

- 🗆 Each element C,i,j
  - = dot(row(A,i),col(B,j))
- Parallelization strategy
  - Each thread computes element in C
  - 2D kernel



#### Matrix multiplication implementation

103

global void mat\_mul(float \*a, float \*b,
 float \*c, int width)

// calc row & column index of output element
int row = blockIdx.y\*blockDim.y + threadIdx.y;
int col = blockIdx.x\*blockDim.x + threadIdx.x;

```
float result = 0;
```

```
// do dot product between row of a and column of b
for(int k = 0; k < width; k++) {
    result += a[row*width+k] * b[k*width+col];
}
c[row*width+col] = result;</pre>
```

B

### Matrix multiplication performance

104

Loads per dot product term	2 (a and b) = 8 bytes
FLOPS	2 (multiply and add)
AI	2 / 8 = 0.25
Performance GTX 580	1581 GFLOPs
Memory bandwidth GTX 580	192 GB/s
Attainable performance	192 * 0.25 = 48 GFLOPS
Maximum efficiency	3.0 % of theoretical peak

#### Data reuse

- Each input element
   in A and B is read
   WIDTH times
- Load elements into shared memory
- Have several threads use local version to reduce the memory bandwidth



## Using shared memory

- Partition kernel loop into phases
- In each thread block, load a tile of both matrices into shared memory each phase
- Each phase, each thread computes a partial result



#### Matrix multiply with shared memory

107

\_\_\_\_\_\_\_shared\_\_\_\_float s\_b[TILE\_WIDTH][TILE\_WIDTH];

// calculate the row & column index int row = by\*blockDim.y + ty; int col = bx\*blockDim.x + tx;

float result = 0;

#### Matrix multiply with shared memory

// loop over input tiles in phases
for(int p = 0; p < width/TILE\_WIDTH; p++) {
 // collaboratively load tiles into shared memory
 s\_a[ty][tx] = a[row\*width + (p\*TILE\_WIDTH + tx)];
 s\_b[ty][tx] = b[(p\*TILE\_WIDTH + ty)\*width + col];
 \_\_syncthreads();</pre>

```
c[row*width+col] = result;
```



108
## Use of Barriers in mat\_mul

- Two barriers per phase:
  - syncthreads after all data is loaded into shared memory
  - syncthreads after all data is read from shared memory
  - Second <u>syncthreads</u> in phase p guards the load in phase p+1
- Use barriers to guard data
   Guard against using uninitialized data
   Guard against corrupting live data

## Matrix multiplication performance

	Original	shared memory
Global loads	$2N^3 * 4$ bytes	(2N <sup>3</sup> / TILE_WIDTH) * 4 bytes
Total ops	2N <sup>3</sup>	2N <sup>3</sup>
AI	0.25	0.25 * TILE_WIDTH

Performance GTX 580	1581 GFLOPs
Memory bandwidth GTX 580	192 GB/s
Al needed for peak	1581 / 192 = <b>8.23</b>
TILE_WIDTH required to achieve peak	0.25 * TILE_WIDTH = 8.23, TILE_WIDTH = 32.9

## 4 CUDA: optimizing your application

- 1. Occupancy
- 2. Shared Memory
- 3. Coalescing
- 4. Streams
- 5. Shared Memory Bank Conflicts

## Coalescing

#### 112

traditional multi-core optimal memory access pattern



many-core GPU optimal memory access pattern



#### Consider the stride of your accesses

113

}

\_global\_\_\_void foo(int\* input, float3\* input2) { int i = blockDim.x \* blockIdx.x + threadIdx.x;

// Stride 1, OK!
int a = input[i];

// Stride 2, half the bandwidth is wasted
int b = input[2\*i];

// Stride 3, 2/3 of the bandwidth wasted
float c = input2[i].x;

### Example: Array of Structures (AoS)

```
114
```

```
struct record {
    int key;
    int value;
    int flag;
};
```

```
record *d_records;
cudaMalloc((void**)&d_records, ...);
```

#### Example: Structure of Arrays (SoA)

```
115
```

```
Struct SoA {
    int* keys;
    int* values;
    int* flags;
};
```

```
SoA d_SoA_data;
cudaMalloc((void**)&d_SoA_data.keys, ...);
cudaMalloc((void**)&d_SoA_data.values, ...);
cudaMalloc((void**)&d_SoA_data.flags, ...);
```

#### Example: SoA vs AoS

}

global void kernel(record\* AoS data, SoA SoA data) { int i = blockDim.x \* blockIdx.x + threadIdx.x; // AoS wastes bandwidth int key1 = AoS data[i].key; // SoA efficient use of bandwidth int key2 = SoA data.keys[i];

## Memory Coalescing

- 117
- Structure of arrays is often better than array of structures
- □ Very clear win on regular, stride 1 access patterns
- Unpredictable or irregular access patterns are case-by-case
- □ Can lose a factor of 10x 30x!

## 118 CUDA: Streams

- 1. Occupancy
- 2. Shared Memory
- 3. Coalescing
- 4. Streams
- 5. Shared Memory Bank Conflicts

## What are streams?

- Stream = a sequence of operations that execute on the device in the order in which they are issued by the host code.
- □ Same stream: In-Order execution
- Different streams: Out-of-Order execution
- Default stream = Synchronizing stream
  - No operation in the default stream can begin until all previously issued operations in any stream on the device have completed.
  - An operation in the default stream must complete before any other operation in any stream on the device can begin.

## Default stream: example

120

cudaMemcpy(d\_a, a, numBytes,cudaMemcpyHostToDevice); increment<<<1,N>>>(d\_a);

CpuFunction(b);

cudaMemcpy(a, d\_a, numBytes,cudaMemcpyDeviceToHost);

All operations happen in the same stream

- Device (GPU)
  - Synchronous execution
    - all operations execute (in order), one after the previous has finished
  - Unaware of CpuFunction()
- Host (CPU)
  - Launches increment and regains control
  - \*May\* execute CpuFunction \*before\* increment has finished
  - Final copy starts \*after\* both increment and CpuFunction() have finished

## Non-default streams

- Enable asynchronous execution and overlaps
  - Require special creation/deletion of streams
    - cudaStreamCreate(&stream1)
    - cudaStreamDestroy(stream1)
  - Special memory operations
    - cudaMemcpyAsync(deviceMem, hostMem, size, cudaMemcpyHostToDevice, stream1)
  - Special kernel parameter (the 4<sup>th</sup> one)
    - increment<<<1, N, 0, stream1>>>(d\_a)
- Synchronization
  - All streams
    - cudaDeviceSynchronize()
  - Specific stream:
    - cudaStreamSyncrhonize(stream1)

## Computation vs. communication

122

//Single stream, numBytes = 16M, numElements = 4M
cudaMemcpy(d\_a, a, numBytes, cudaMemcpyHostToDevice);
kernel<<blocks,threads>>(d\_a, firstElement);
cudaMemcpy(a, d\_a, numBytes, cudaMemcpyDeviceToHost);

#### C1060 (pre-Fermi): 12.9ms

Copy Engine	H2D - Stream 0		D2H - 0
Kernel Engine		0	
C2050 (F	<sup>-</sup> ermi): 9.9ms		
H2D Engine	Stream 0		
Kernel Engine		0	
D2H Engine			0

## Computation-communication overlap[1]\*

```
for (int i = 0; i < nStreams; ++i) {
    int offset = i * streamSize;
    cudaMemcpyAsync(&d_a[offset], &a[offset], streamBytes,
    stream[i]);
    kernel<<blocks,threads,0,stream[i]>>(d_a, offset);
    cudaMemcpyAsync(&a[offset], &d_a[offset], streamBytes,
    stream[i]);
}
```

#### C1060 (pre-Fermi): 13.63 ms (worse than sequential)

Copy Engine	H2D - I		D2H - 1	H2D - 2		D2H - 2 H2D - 3		D2H - 3 H2D - 4		D2H - 4
Kernel Engine		1			2		3		4	

#### C2050 (Fermi): 5.73 ms (better than sequential)

123

H2D Engine	1.1	2	3	4	
Kernel Engine		1	2	3	4
D2H Engine			1	2	3

https://github.com/parallel-forall/code-samples/blob/master/series/cuda-cpp/overlap-data-transfers/async.cu

## Computation-communication overlap[2]\*

124

```
for (int i = 0; i < nStreams; ++i) offset[i]=i * streamSize;</pre>
for (int i = 0; i < nStreams; ++i)</pre>
    cudaMemcpyAsync(&d a[offset[i]], &a[offset[i]], streamBytes,
        cudaMemcpyHostToDevice, stream[i]);
for (int i = 0; i < nStreams; ++i)</pre>
        kernel<<blocks,threads,0,stream[i]>>(d a, offset);
for (int i = 0; i < nStreams; ++i)</pre>
    cudaMemcpyAsync(&a[offset], &d a[offset], streamBytes,
        cudaMemcpyDeviceToHost, stream[i]);
C1060 (pre-Fermi): 8.84 ms (better than sequential)
         H2D - I H2D - 2 H2D - 3 H2D - 4 D2H - I D2H - 2 D2H - 3 D2H - 4
Copy Engine
Kernel Engine
                      2
                           3
C2050 (Fermi): 7.59 ms (better than sequential, worse than v1)
H2D Engine
                 2
                      3
                           4
                      2
Kernel Engine
                           3
                                 4
D2H Engine
                                           2
                                                3
```

http://devblogs.nvidia.com/parallelforall/how-overlap-data-transfers-cuda-cc/

## 125 CUDA: optimizing your application

- 1. Occupancy
- 2. Shared Memory
- 3. Coalescing
- 4. Streams
- 5. Shared Memory Bank Conflicts

## Shared Memory Banks

- Shared memory is banked
  - Only matters for threads within a warp
  - Full performance with some restrictions
    - Threads can each access different banks
    - Or can all access the same value
- Consecutive words are in different banks
- If two or more threads access the same bank but different value, we get bank conflicts

#### **Bank Addressing Examples: OK**





#### Bank Addressing Examples: BAD





## Trick to Assess Performance Impact

- Change all shared memory reads to the same value
- □ All broadcasts = no conflicts
- Will show how much performance could be improved by eliminating bank conflicts
- The same doesn't work for shared memory writes
   So, replace shared memory array indices with threadIdx.x
  - (Could also be done for the reads)



## Portability

- Inter-family vs inter-vendor
  - NVIDIA Cuda runs on all NVIDIA GPU families
  - OpenCL runs on all GPUs, Cell, CPUs
- Parallelism portability
  - Different architecture requires different granularity
  - Task vs data parallel
- Performance portability
  - Can we express platform-specific optimizations?

## The Khronos group



# **OpenCL: Open Compute Language**

- Architecture independent
- Explicit support for many-cores
- Low-level host API
  - Uses C library, no language extensions
- Separate high-level kernel language
  - Explicit support for vectorization
- Run-time compilation
- Architecture-dependent optimizations
  - Still needed
  - Possible

#### Cuda vs OpenCL Terminology

CUDA	OpenCL
Thread	Work item
Thread block	Work group
Device memory	Global memory
Constant memory	Constant memory
Shared memory	Local memory
Local memory	Private memory

### Cuda vs OpenCL Qualifiers

135



Variables		
CUDA	OpenCL	
constant	constant	
device	global	
shared	_local	

#### Cuda vs OpenCL Indexing

136

CUDA	OpenCL
gridDim	get_num_groups()
blockDim	get_local_size()
blockldx	get_group_id()
threadIdx	get_local_id()
Calculate manually	get_global_id()
Calculate manually	get_global_size()

 $_syncthreads() \rightarrow barrier()$ 

#### Vector add: Cuda vs OpenCL kernel

137





### OpenCL VectorAdd host code (1)

```
138
```

```
const size_t workGroupSize = 256;
const size_t nrWorkGroups = 3;
const size_t totalSize = nrWorkGroups * workGroupSize;
```

```
cl_platform_id platform;
clGetPlatformIDs(1, &platform, NULL);
```

```
// create properties list of key/values, 0-terminated.
cl_context_properties props[] = {
    CL_CONTEXT_PLATFORM, (cl_context_properties)platform,
    0
};
```

cl\_context context = clCreateContextFromType(props, CL\_DEVICE\_TYPE\_GPU, 0, 0, 0);

### OpenCL VectorAdd host code (2)

139

// create command queue on 1st device the context reported
cl\_command\_queue commandQueue =

clCreateCommandQueue(context, device, 0, 0);

// create & compile program
cl\_program program = clCreateProgramWithSource(context, 1,
 &programSource, 0, 0);
clBuildProgram(program, 0, 0, 0, 0, 0);

// create kernel
cl\_kernel kernel = clCreateKernel(program, "vectorAdd",0);

#### OpenCL VectorAdd host code (3)

140

float\* A, B, C = new float[totalSize]; // alloc host vecs
// initialize host memory here...

// allocate device memory
cl\_mem deviceA = clCreateBuffer(context,
 CL\_MEM\_READ\_ONLY | CL\_MEM\_COPY\_HOST\_PTR,
 totalSize \* sizeof(cl float), A, 0);

cl\_mem deviceB = clCreateBuffer(context, CL\_MEM\_READ\_ONLY | CL\_MEM\_COPY\_HOST\_PTR, totalSize \* sizeof(cl\_float), B, 0);

cl\_mem deviceC = clCreateBuffer(context, CL\_MEM\_WRITE\_ONLY, totalSize \* sizeof(cl\_float), 0, 0);

### OpenCL VectorAdd host code (4)

141

- // setup parameter values
- clSetKernelArg(kernel, 0, sizeof(cl\_mem), &deviceA); clSetKernelArg(kernel, 1, sizeof(cl\_mem), &deviceB); clSetKernelArg(kernel, 2, sizeof(cl\_mem), &deviceC);
- clEnqueueNDRangeKernel(commandQueue, kernel, 1, 0,
   &totalSize, &workGroupSize, 0,0,0); // execute kernel
- // copy results from device back to host, blocking
  clEnqueueReadBuffer(commandQueue, deviceC, CL\_TRUE, 0,
   totalSize \* sizeof(cl\_float), C, 0, 0, 0);
- delete[] A, B, C; // cleanup clReleaseMemObject(deviceA); clReleaseMemObject(deviceB); clReleaseMemObject(deviceC);



## Summary and conclusions

- Higher performance cannot be reached by increasing clock frequencies anymore
- Solution: introduction of large-scale parallelism
- Multiple cores on a chip
  - Today:
    - Up to 48 CPU cores in a node
    - Up to 3200 compute elements on a single GPU
  - Host system can contain multiple GPUs: 10,000+ cores
  - We can build clusters of these nodes!
- □ Future: 100,000s millions of cores?

## Summary and conclusions

- Many different types of many-core hardware
- Very different properties
  - Performance
  - Programmability
  - Portability
- It's all about the memory
- Choose the right platform for your application
  - Arithmetic intensity / Operational intensity
  - Roofline model
# Open questions

- New application domains e.g., signal processing, graph processing.
  - Performance analysis
  - Peformance prediction
  - Modeling
- Memory patterns understanding, description, detection, automatic improvement
  - Local memory usage
- Heterogeneous computing
  - Using both the host and the device
- Application-device fitting

### Questions?

- Slides are/will be available
- If you are interested in working with us on using GPUs for new applications, let us know!

A.L.Varbanescu@uva.nl

# Backup slides

### Example: Work queue

148

{

// For algorithms where the amount of work per item
// is highly non-uniform, it often makes sense to

// continuously grab work from a queue.

> int i = threadIdx.x + blockDim.x \* blockIdx.x; int q\_index = atomicInc(q\_counter, queue\_max); int result = do\_work(work\_q[q\_index]); output[q\_index] = result;

- // Adjacent Difference application:
- // compute result[i] = input[i] input[i-1]

```
__global__ void adj_diff_naive(int *result, int *input) {
    // compute this thread's global index
    unsigned int i = blockDim.x * blockIdx.x + threadIdx.x;
```

```
if(i > 0) {
   // each thread loads two elements from device memory
   int x_i = input[i];
   int x i minus one = input[i-1];
```

```
result[i] = x_i - x_i_minus_one;
```





```
152
```

}

\_global\_\_\_ void adj\_diff(int \*result, int \*input) { unsigned int i = blockDim.x \* blockIdx.x + threadIdx.x;

\_\_shared\_\_ int s\_data[BLOCK\_SIZE]; // shared, 1 elt / thread // each thread reads 1 device memory elt, stores it in s\_data s\_data[threadIdx.x] = input[i];

```
if(threadIdx.x > 0) {
   result[i] = s_data[threadIdx.x] - s_data[threadIdx.x-1];
} else if(i > 0) {
   // I am thread 0 in this block: handle thread block boundary
   result[i] = s_data[threadIdx.x] - input[i-1];
```

### Using shared memory: coalescing

153

}

\_global\_\_\_ void adj\_diff(int \*result, int \*input) { unsigned int i = blockDim.x \* blockIdx.x + threadIdx.x;

\_\_shared\_\_ int s\_data[BLOCK\_SIZE]; // shared, 1 elt / thread // each thread reads 1 device memory elt, stores it in s\_data s\_data[threadIdx.x] = input[i]; // COALESCED ACCESS!

```
if(threadIdx.x > 0) {
   result[i] = s_data[threadIdx.x] - s_data[threadIdx.x-1];
} else if(i > 0) {
   // I am thread 0 in this block: handle thread block boundary
   result[i] = s_data[threadIdx.x] - input[i-1];
```

# Backup slides

# CPU vs GPU Chip

155

#### AMD Magny-Cours (6 cores)

#### ATI 4870 (800 cores)



2 billion transistors 346 mm<sup>2</sup> 1 billion transistors 256 mm<sup>2</sup>









# Latest generation ATI

- 157
- Southern Islands
- □ 1 chip: HD 7970
  - 2048 cores
  - 264 GB/sec memory bandwidth
  - 3.8 tflops single, 947 gflops double precision
  - Maximum power: 250 Watts
  - **399** euros!
- □ 2 chips: HD 7990
  - 4096 cores, 7.6 tflops
- Comparison: entire 72-node DAS-4 VU cluster has
   4.4 tflops

#### ATI 5870 architecture overview



# ATI 5870 SIMD engine

- □ Each of the 20 SIMD engines has:
  - **1** 16 thread processors x 5 stream cores = 80 scalar stream processing units
  - 20 \* 16 \* 5 = 1600 cores total
  - 32KB Local Data Share
  - its own control logic and runs from a shared set of threads
  - a dedicated fetch unit with 8KB L1 cache
  - a 64KB global data share to communicate with other SIMD engines



# ATI 5870 thread processor



- Each thread processor includes:
  - 4 stream cores + 1 special function
    - stream core
  - general purpose registers
- □ FMA in a single clock



# ATI 5870 Memory Hierarchy

- EDC (Error Detection Code)
  - CRC Checks on Data Transfers for Improved Reliability at High Clock Speeds
- Bandwidths
  - Up to 1 TB/sec L1 texture fetch bandwidth
  - Up to 435 GB/sec between L1 & L2
  - 153.6 GB/s to device memory
  - PCI-e 2.0, 16x: 8GB/s to main memory

# Unified Load/Store Addressing



#### **Unified Address Space**







# Partition data into subsets that fit into shared memory



#### Handle each data subset with one thread block



Load the subset from device memory to shared memory, using multiple threads to exploit memorylevel parallelism





Perform the computation on the subset from shared memory



Copy the result from shared memory back to device memory