# GPU Programming with CUDA

#### Ph.D. Pedro Velho and Eng. Mônica Liliana Hernandez Ariza

# Disclaimer

I am not an expert on the subject

Working with GPU for less than one year

I wish I have had this course at the time!

#### Download the slides and examples

# http://190.168.117.201/~velho/cuda.html

Meeting the audience!

How many of you used concurrent programming before?

How many threads?

How many already used CUDA?



- A few general purpose cores
- Bíg cache memory
- Eg.: Nehalem i7 quad-core
  - 4 cores (8 threads)
  - Cache is about 50% of die area



# Lltografia

- A few general purpose cores
- Big cache memory
- Eg.: Nehalem í7 quad-core
  - 4 cores (8 threads)
  - Cache is about 50% of die area



C. C.Laure		THE PARTY OF			Service Parties	N Ma	- Kaset
	1000		Margari,		STORAGY.		ALC: UNK
Mano//	Herrory	Medically	Weiner.	Minney	Messory	thermsp	Marcusy



- Design goal massívely parallel graphics
- A lot of replicated functional units
- Small cache síze
- Eg .: NVIDIA GTX280
  - 240 SP (streaming processors)
  - support for 30720 simultaneous threads

# Computer Graphics is a Computational intensive application





# Computer Graphics is a Computational intensive application

A lot of \$\$ from game industry

Video game software sales: 1998-2007 \$ Billions



# Computer Graphics is a Computational intensive application

A lot of \$\$ from game industry

Expressive gain in performance for parallel graphics rendering

Caught attention from the scientific community

Video game software sales: 1998-2007 \$ Billions



# GPU is also adapted to several scientific applications







Molecular Biology

Fluid Simulation

Weather Forecast



Wednesday, July 11, 12



#### Potential Gain in Performance



Several guys from Intel



Victor W Lee *et. al.*, Debunking the 100X GPU vs. CPU Myth: An Evaluation of Throughput Computing on CPU and GPU



#### Potential Gain in Performance

> 100 times faster!

Several guys from Intel



Victor W Lee *et. al.*, Debunking the 100X GPU vs. CPU Myth: An Evaluation of Throughput Computing on CPU and GPU



 $\frac{T_{\rm GTX280}}{T_{\rm Core~i7}}$ 



#### Substantial gain in execution time (10x)!

before GPU	with GPU		
one year	one month plus a week		
one day	two hours and twinty four minutes		
one hour	six minutes		

#### GPU Programming today

Don't need to port the application to DirectX or OpenGL



- Propretary (only work on NVIDIA)
- Enhanced software support
- Several software libraries and examples

CUDA vs OpenCl - Open specification

- Work on NVIDIA and ATI video cards
- Aim at any computing device



Wednesday, July 11, 12

# Previously on parallel architectures

#### Superscalar processors



# Superscalar processors make the illusion of concurrent execution

Instruction from one thread arrive A hardware íssue unít decídes whích ínstructíons can execute símultaneously

 Front end

 isue unit

 Image: CPU

 Front end

 Image: CPU

 Image: CPU

# A program has instructions for several threads in memory



blue thread
red thread
green thread
yellow thread

Wednesday, July 11, 12

# Single threaded multicore



Wednesday, July 11, 12

### Single threaded multicore



## Single threaded multicore



Wednesday, July 11, 12

# Super-threadeding



# Multi-threadeding





Can execute instruction from more than 1 thread at a time





# Back to GPU architecture

# Streaming Processor (SP)



# Streaming Processor (SP)



# Streaming Multiprocessor (SM)

Streaming

Each SFU 4 FP multíply for sín, cosín





# Streaming Multiprocessor (SM)

Each SFU 4 FP multíply for sín, cosín



Multí-threaded can íssue several ínstructíons

Array of 8 (eight) SPS

# GPU Architecture (GT200)

#### Texture Processor Cluster 3 SM's

TPC (GT200)							
Geometry Controller							
SMC							
SM I cache MT issue C cache SP SP SP SP SP SP SP SP SP SP SFU SFU SFU SFU Shared Memory	SM I cache MT issue C cache SP SP SP SP SP SP SP SP SP SP SFU SFU SFU SFU Shared Memory	SM I cache MT issue C cache SP SP SP SP SP SP SP SP SP SP SFU SFU SFU SFU					
Texture Units Texture L1							

Wednesday, July 11, 12
#### GPU Architecture (GT200)

#### The beast

10 TPC's 3 SM's per TPC 8 SP's per SM

Total of 240 SP's



#### GPU Architecture (GT200)





## Schedule per group of 32 threads, called a **warp**

#### GPU Architecture (GT200)

# Each SM handles 32 warps simultaneously

#### $32 \times 32 = 1024$ threads per SM



## $1024 \times 30 = 30720$ simultaneous threads



Wednesday, July 11, 12



Massively parallel processor (GT200 - 30720 Threads)

- CPU send burst of threads to execute on the GPU

Use DMA to transfer from CPU DRAM to GPU DRAM

CPU becomes free after dispatching the threads

- Can do something useful meanwhile

Applications must be rewritten to cope with GPU



Wednesday, July 11, 12

#### Single threaded application



GT200 supports up to 30720 threads simultaneously!!!

## Definition of a single thread computing function (or **kernel**)



- 1- How to Compute the thread ID?
- 2- How do we copy data from CPU to GPU?
- 3- How to dispatch kernel on the device?
- 4- How to get results back when done?

Have support for operations on the Host (CPU) and Device (GPU)

Copy data from Host to Device
 Execute kernel on the device
 Wait for kernel to finish
 Copy data from Device to Host

Depends on the programming interface

mallocDeviceMemory copyFromHostToDevice computeKernel copyFromDeviceToHost



Wednesday, July 11, 12

- C extension
- Support for several platforms:
  - Linux
  - Windows
  - MacOS
- Need to install NVIDIA Driver, Toolbox and SDK

#### Provide several libraries



#### STL C++ Port to CUDA



Linear Algebra cuBLAS

**Requirements for Linux** 

- I NVIDIA CUDA aware card
- GCC installed
- Downloaded Toolkit, Driver, and SDK

Step-by-step installation:

- Exit GUI (Ctrl+Alt+Backspace)
- Install the CUDA Toolkit \$ ./cudatoolkit\_4.2.9\_linux\_64\_ubuntu11.04.run
- Install the driver \$ sudo ./devdriver\_4.2\_linux\_64\_295.41.run
- Restart GUI
  \$ sudo /etc/init.d/gdm start
- Install SDK
  - \$./gpucomputingsdk\_4.2.9\_linux.run
- Only the driver requires superuser priviledges



#### **Exploring SDK demos**

SDK has many toy applications:

\$ cd \$HOME/NVIDIA\_CUDA\_SDK

\$ make

- \$ make check
- \$ C/bin/linux/release/programName

}

#### Example

```
...
float *aHost, *bHost, *cHost;
...
global____void kernel(float *a, float *b, float *c){
    int i = threadidx.x;
    c[i] = a[i] + b[i];
}
```

```
int main(){
   float *aDev, *bDev, *cDev;
```

cudaMalloc(void \*aDev, 1024 \* sizeof(float)); cudaMemcpy(aDev, aHost, 1024 \* sizeof(float));

cudaMalloc(void \*bDev, 1024 \* sizeof(float)); cudaMemcpy(bDev, bHost, 1024 \* sizeof(float));

**kernel**<<<1, **1024**>>> (aDev, bDev, cDev);

cudaFree(aDev); cudaFree(bDev); cudaFree(cDev);

#### **Exercise Zero**

Device query.

#### Function directives

Kernel function has stric	t properties
must return void	
no static variables	
no recurrency	
no variable arguments	

	Execute on	Called from
device float DeviceFunc()	device	device
<b>global</b> void kernelFunc()	device	host
<pre>host float HostFunc()</pre>	host	host

can be used combined with \_\_device\_\_

Memory allocation

## cudaMalloc(...)

Allocate global memory

2 parameters: Pointer Number of bytes





Transfer data

### cudaMemcpy(...)

- 4 parameters: Source pointer Destination pointer Bytes to copy Transfer type
  - HostToHost HostToDevice DeviceToHost DeviceToDevice



Memory allocation

#### cudaFree(...)

Frees global memory I parameter: Pointer



#### Exercise I

Implementing the sum of two vectors using CUDA. You can assume that the array has a maximum of 1024 elements.



#### Thread indexing

Threads are organized in blocks

Blocks are organized in grids

Legacy from CG applications

Indexing may be in **3D**!!!



#### Grid

Block	Block	Block
(0,0)	(1,0)	(2,0)
Block	Block	Block
(0,1)	(I,I)	(2,1)
Block	Block	Block
(0,2)	(1,2)	(2,2)

#### Grid

Block (0,0)	Block (1,0)	Block (2,0)	Blo	ck		
Block	Block	Block	Thread	Thread	Thread	Thread
(0,1)	(I,I)	(2,1)	(0,0)	(1,0)	(2,0)	(3,0)
Block	Block	Block	Thread	Thread	Thread	Thread
(0,2)	(1,2)	(2,2)	(0,1)	(1,1)	(2,1)	(3,1)
			Thread (0,2)	Thread (1,2)	Thread (2,2)	Thread (3,2)

#### mapping threads

Block	Thread	Thread	Thread	Thread
	(0,0)	(1,0)	(0,0)	(1,0)
(0,0)	Thread	Thread	Thread	Thread
	(0,1)	(I,I)	(0,1)	(I,I)
	Thread	Thread	Thread	Thread
Block	(0,0)	(1,0)	(0,0)	(1,0)

Block (1,0)

Block (I,I)

dim3 Grid(2,2); dim3 Block(2,2); kernel<<<Grid,Block>>>(parameters);



#### How can we arrange 6 threads?

#### Block (0,0)

Thread	Thread	Thread	Thread	Thread	Thread
(0,0)	(1,0)	(2,0)	(3,0)	(4,0)	(5,0)





#### How can we arrange 6 threads?

#### Block (0,0)

Thread	Thread	Thread	Thread	Thread	Thread
(0,0)	(1,0)	(2,0)	(3,0)	(4,0)	(5,0)

Block (0,0)		Blc	ock (I	,0)	
Thread	Thread	Thread	Thread	Thread	Thread
(0,0)	(1,0)	(2,0)	(0,0)	(1,0)	(2,0)

#### How can we arrange 6 threads?

Block (0,0)		Block (1,0)		Bloc	< (2,0)
Thread (0,0)	Thread (1,0)	Thread (0,0)	Thread (1,0)	Thread (0,0)	Thread (1,0)

How can we arrange 6 threads?



#### Mapping on an unique grid



#### Mapping on an unique grid



idx = blockldx.x\*blockDim.x + threadIdx.x;

idy = blockIdx.y\*blockDim.y + threadIdx.y;

#### Mapping on an unique grid



(1,3)

(2,3)

(3,3)

(0,3)

idx = blockldx.x\*blockDim.x + threadIdx.x;

idy = blockldx.y\*blockDim.y + threadldx.y;

#### Get an unique thread index

Thread	Thread	Thread	Thread
(0,0)	(1,0)	(2,0)	(3,0)
Thread	Thread	Thread	Thread
(0,1)	(I,I)	(2,1)	(3,1)
Thread	Thread	Thread	Thread
(0,2)	(1,2)	(2,2)	(3,2)
Thread	Thread	Thread	Thread
(0,3)	(1,3)	(2,3)	(3,3)

k = idx + idy\*blockDim.x\*gridDim.x;
## CUDA Threads

#### Get an unique thread index

Thread	Thread	Thread	Thread
(0,0)	(1,0)	(2,0)	(3,0)
Thread	Thread	Thread	Thread
(0,1)	(I,I)	(2,1)	(3,1)
Thread	Thread	Thread	Thread
(0,2)	(1,2)	(2,2)	(3,2)
Thread	Thread	Thread	Thread
(0,3)	(1,3)	(2,3)	(3,3)



Thread	Thread	Thread	Thread
(0)	(1)	(2)	(3)
Thread	Thread	Thread	Thread
(4)	(5)	(6)	(7)
Thread	Thread	Thread	Thread
(8)	(9)	(10)	(11)

k = idx + idy\*blockDim.x\*gridDim.x;

#### Exercise II

Implementing the sum of two vectors using CUDA of a unlimited number of elements.

## **Exercise III** Filling the gaps of the third example.

### **Exercise IV** Implement a MPI/CUDA hybrid application.



Wednesday, July 11, 12

## GPU is good for...

loosely coupled threads (avoid synchronization)

computing bound applications

these architectures can not replace general purpose CPU

great insight for low power consumption and future architectures

#### CUDA Pros

## **CUDA** Cons

Support for several OS

A lot of documentation

Many libraries available

Great performance

NVIDIA propretary

# Architectures of the Future

Highly heterogeneous

Intel SandBridge AMD Fusion NVIDIA Tegra

Wednesday, July 11, 12

## Bibliography

Programming Massively Parallel Processors: A Handson Approach, David B. Kirk and Wen-Mei Hwu, Second Edition, Morgan Kaufmann, 2009



NVIDIA developer zone, <u>http://developer.nvidia.com/</u>