### Plan



	< □ >	→ @ ▶ ★ 글 ▶ ★ 글 ▶ = 글	うくで		< □ >	★國▶★필▶★필▶ _ 필	うくで
(Moreno Maza)	CS4402-9535: Many-core Computing with Cl	UWO-CS4402-CS9535	1 / 83	(Moreno Maza)	CS4402-9535: Many-core Computing with CL	UWO-CS4402-CS9535	2 / 83
GPUs and CUDA	A: a Brief Introduction			GPUs and CU	IDA: a Brief Introduction		
Plan				GPUs			
i iuni				01 00			

### GPUs and CUDA: a Brief Introduction

- 2 CUDA Programming Model
- 3 CUDA Memory Model
- 4 CUDA Programming Basics
- 5 CUDA Hardware Implementation
- 6 CUDA Programming: Scheduling and Synchronization
- CUDA Tools
- 8 Sample Programs

- GPUs are massively multithreaded manycore chips:
  - NVIDIA Tesla products have up to 448 scalar processors with
  - over 12,000 concurrent threads in flight and
  - 1030.4 GFLOPS sustained performance (single precision).
- Users across science & engineering disciplines are achieving 100x or better speedups on GPUs.



# CUDA

## CUDA programming and memory models in a nutshell

- CUDA is a scalable parallel programming model **and** a software environment for parallel computing:
  - $\bullet\,$  Minimal extensions to familiar C/C++ environment
  - Heterogeneous serial-parallel programming model
- GPU Computing with CUDA brings data-parallel computing to the masses
  - as of 2008, over 46,000,000 (100,000,000, as of 2009) CUDA-capable GPUs sold,
  - a *developer kit* costs about \$400 (for 500 GFLOPS).
- Massively parallel computing has become a commodity technology!



(Moreno Maza) CS4402-9535: Many-core Co CUDA Programming Model	< □ > < □ > < □ > < ≥ > < ≥ > < ≥ > ≥       omputing with CL     UWO-CS4402-CS9535	୬	(Moreno Maza)	CS4402-9535: Many-core Computing CUDA Programming Model	d □ ► < d g with CL	→ < ≥ → < ≥ → < ≥ → UWO-CS4402-CS9535	୬୯.୯ 6 / 83
Plan			CUDA design g	goals			
<ol> <li>GPUs and CUDA: a Brief Introduction</li> <li>CUDA Programming Model</li> </ol>	ſ		<ul> <li>Enable hetero</li> <li>Scale to 100's</li> <li>Use C/C++ v</li> <li>Let programm</li> </ul>	ogeneous systems (i.e., CPU s of cores, 1000's of parallel with minimal extensions ners focus on parallel algorit	+GPU)   threads thms		
<ul> <li>3 CUDA Memory Model</li> <li>4 CUDA Programming Basics</li> <li>5 CUDA Hardware Implementation</li> </ul>			32 SP Cores	SM 128 SP Cores William CCEAN BP SP SP SP SP SP SP SP SP SP SP SP SP SP SP SP SP SP	SM Cette Cetter SP SP SP SP	SM I-Cache MT Issue	
<ul> <li>6 CUDA Programming: Scheduling and</li> <li>7 CUDA Tools</li> <li>8 Sample Programs</li> </ul>	Synchronization <□> <♂→ <≥> <≥> ≥	୬୯୯	GPU	Z40 SP		C-Cache SP SP SP SP SP SP SP SP SFU SFU SFU SFU Shared Memory	ল হাজ
(Moreno Maza) CS4402-9535: Many-core Co	omputing with CL UWO-CS4402-CS9535	7 / 83	(Moreno Maza)	CS4402-9535: Many-core Computing	g with Cl	UWO-CS4402-CS9535	8 / 83

#### CUDA Programming Model

## Heterogeneous programming (1/3)

- A CUDA program is a serial program with parallel kernels, all in C.
- The serial C code executes in a host (= CPU) thread
- The parallel kernel C code executes in many device threads across multiple GPU processing elements, called streaming processors (SP).



# Heterogeneous programming (3/3)

- The parallel code is written for a thread
  - Each thread is free to execute a unique code path
  - Built-in thread and block ID variables are used to map each thread to a specific data tile (more on this soon).
- Thus, each thread executes the same code on different data based on its thread and block ID.

Serial Code	Host	
Parallel Kernel KernelA (args);	Device	
Serial Code	Host	
Parallel Kernel KernelB (args);	Device	

# IDs and dimensions (1/2)

(Moreno Maza)

- A kernel is a grid of thread blocks.
- Each thread block has a 2-D ID, which is unique within the grid.
- Each thread has a 2-D ID, which is unique within its thread block.

CS4402-9535: Many-core Computing with Cl

- The dimensions are set at launch time by the host code
- IDs and dimension sizes are accessed via global variables in the **device code**: threadIdx, blockIdx, ..., blockDim, gridDim.
- Simplify memory addressing when processing multidimensional data

	G	id 1			
		Block (0, 0)	Bloc (1, 0	k )	Block (2, 0)
		Block (0, 1)	Bloc (1, 1	k )	Block (2, 1)
Block (	1, 1)				
Thread (0, 0)	Thread (1, 0)	Thread (2, 0)	Thread (3, 0)	Threa (4, 0	id )
Thread (0, 1)	Thread (1, 1)	Thread (2, 1)	Thread (3, 1)	Threa (4, 1	id )
Thread	Thread	Thread	Thread	Threa	d

3

11 / 83

・ 同 ト ・ ヨ ト ・ ヨ ト

UWO-CS4402-CS9535

10 / 8

# Heterogeneous programming (2/3)

- Thus, the parallel code (kernel) is launched and executed on a device by many threads.
- Threads are grouped into thread blocks (more on this soon).
- One kernel is executed at a time on the device.

Host

Device

Host

Device

CUDA Programming Model

CUDA Programming Model

• Many threads execute each kernel.

Serial Code

Parallel Kernel KernelA (args);

Serial Code

**Parallel Kernel** 

KernelB (args);

	Model	
IDs and dimensions $(2/2)$		Example: increment array elements $(1/2)$
Device	Block       Block       Block         (0, 0)       Block       Block         (1, 0)       Block       Block         (1, 1)       Block       Block         (2, 1)       Block       Block	Increment N-element vector a by scalar b Let's assume N=16, blockDim=4 -> 4 blocks int idx = blockDim.x * blockId.x + threadIdx.x;
(Moreno Maza) (Moreno Maza) (Moreno Maza) (Moreno Maza) (Moreno Maza) (Moreno Maza)	Thread Thread Thread $(4, 0)$ Thread Thread Thread $(4, 1)$ Thread Thread Thread $(4, 1)$ Thread Thread $(4, 2)$ Thread $(3, 2)$ Thread $(4, 2)$ Many-core Computing with CL Model Model	blockldx.x=0 blockDim.x=4 threadldx.x=0,1,2,3 idx=0,1,2,3 blockDim.x=4 threadldx.x=0,1,2,3 idx=4,5,6,7 blockldx.x=2 blockDim.x=4 threadldx.x=0,1,2,3 idx=4,5,6,7 blockldx.x=2 blockDim.x=4 threadldx.x=0,1,2,3 idx=8,9,10,11 blockldx.x=3 blockDim.x=4 threadldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=3 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0 threadldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0 threadldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=1 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=1 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0,1,2,3 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0,1,2,3 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0,1,2,3 blockldx.x=0,1,2,3 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0,1,2,3 idx=12,13,14,15 blockldx.x=0,
CPU program	CUDA program	<pre>// allocate host memory unsigned int numBytes = N * sizeof(float)</pre>
<pre>void increment_cpu(float *a, float b, int N) {     for (int idx = 0; idx<n; +="" a[idx]="a[idx]" b;="" idx++)="" pre="" }<=""></n;></pre>	global void increment_gpu(float *a, float b, int N) {     int idx = blockldx.x * blockDim.x + threadldx.x;     if( idx < N)         a[idx] = a[idx] + b; }	<pre>float* h_A = (float*) malloc(numBytes); // allocate device memory float* d_A = 0; cudaMalloc((void**)&amp;d_A, numbytes); // copy data from host to device cudaMemcpy(d_A, h_A, numBytes, cudaMemcpyHostToDevice);</pre>

(Moreno Maza)

CS4402-9535: Many-core Computing with Cl UWO-CS4402-CS9535 15 / 83

(Moreno Maza)

CS4402-9535: Many-core Computing with CL UWO-CS4402-CS9535 16 / 83

CUDA Programming Model

## Thread blocks (1/2)

- A Thread block is a group of threads that can:
  - Synchronize their execution
  - Communicate via shared memory
- Within a grid, thread blocks can run in any order:
  - Concurrently or sequentially
  - Facilitates scaling of the same code across many devices



# Thread blocks (2/2)

- Thus, within a grid, any possible interleaving of blocks must be valid.
- Thread blocks may coordinate but not synchronize
  - they may share pointers
  - they should not share locks (this can easily deadlock).
- The fact that thread blocks cannot synchronize gives scalability:
  - A kernel scales across any number of parallel cores
- However, within a thread bloc, threads in the same block may synchronize with barriers.
- That is, threads wait at the barrier until threads in the same block reach the barrier.

	< □ >	◆御 ▶ ◆ 臣 ▶ ◆ 臣 ▶ ○ 臣	- nac		< □ ▶	(本部)・ 大臣・ 大臣・ 三臣	• ৩৫৫
(Moreno Maza)	CS4402-9535: Many-core Computing with Cl	UWO-CS4402-CS9535	17 / 83	(Moreno Maza)	CS4402-9535: Many-core Computing with Cl	UWO-CS4402-CS9535	18 / 83
	CUDA Memory Model				CUDA Memory Model		
Plan				Memory hierarch	y (1/3)		

### 1 GPUs and CUDA: a Brief Introduction

2 CUDA Programming Model

### 3 CUDA Memory Model

- 4 CUDA Programming Basics
- 5 CUDA Hardware Implementation
- 6 CUDA Programming: Scheduling and Synchronization
- 7 CUDA Tools
- 8 Sample Programs

### Host (CPU) memory:

• Not directly accessible by CUDA threads



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ のへで

19 / 83

(Moreno Maza)

CS4402-9535: Many-core Computing with CL

#### CUDA Memory Model

Memory hierarchy (2/3)

### Global (on the device) memory:

- Also called device memory
- Accessible by all threads as well as host (CPU)
- Data lifetime = from allocation to deallocation



#### CUDA Memory Model

## Memory hierarchy (3/3)

### Shared memory:

- Each thread block has its own shared memory, which is accessible only by the threads within that block
- Data lifetime = block lifetime

#### Local storage:

- Each thread has its own local storage
- Data lifetime = thread lifetime



	< □ >	<li>&lt; ≥ &lt; ≥ &lt; ≥ &lt; ≥</li>	≜			< 🗆	★ ## ★ ★ ★ ★ ★ ★ ★ ★ ★	<u> ৯</u> ০৫
(Moreno Maza)	CS4402-9535: Many-core Computing with Cl	UWO-CS4402-CS9535	21 / 83	()	Moreno Maza)	CS4402-9535: Many-core Computing with CL	UWO-CS4402-CS9535	22 / 83
CUI	DA Programming Basics					CUDA Programming Basics		
				N/ .	1.11.1			

## Plan

- 1 GPUs and CUDA: a Brief Introduction
- 2 CUDA Programming Model
- 3 CUDA Memory Model
- 4 CUDA Programming Basics
- 5 CUDA Hardware Implementation
- 6 CUDA Programming: Scheduling and Synchronization
- 7 CUDA Tools
- 8 Sample Programs

# Vector addition on GPU (1/4)



(Moreno Maza)

23 / 83



### Vector addition on GPU (2/4)

### CUDA Programming Basics

## Vector addition on GPU (3/4)



### Variable Qualifiers (GPU code)

J code`	Launching	kernels	on	GPU
	Luunening	Renneis	011	01 0

device : • stored in global memory (not cached, high latency)	Launch parameters:
• accessible by all threads	<ul> <li>grid dimensions (up to 2D)</li> </ul>
<ul> <li>lifetime: application</li> </ul>	<ul> <li>thread-block dimensions (up to 3D)</li> </ul>
<pre>constant : • stored in global memory (cached)       • read-only for threads, written by host       • Lifetime: application</pre>	<ul> <li>shared memory: number of bytes per block</li> <li>for extern smem variables declared without size</li> <li>Optional, 0 by default</li> </ul>
<pre>shared : • stored in shared memory (latency comparable to registers)</pre>	<ul> <li>stream ID:</li> <li>Optional, 0 by default</li> </ul>
<ul><li>accessible by all threads in the same threadblock</li><li>lifetime: block lifetime</li></ul>	dim3 grid(16, 16);
<ul> <li>Unqualified variables: • scalars and built-in vector types are stored in registers</li> <li>• arrays are stored in device (= global) memory</li> </ul>	dim3 block(16,16); kernel<< <grid, 0="" 0,="" block,="">&gt;&gt;(); kernel&lt;&lt;&lt;32, 512&gt;&gt;&gt;();</grid,>

	∢ □ ▶	→ 御 ト ★ 国 ト ★ 国 ト → 国	1 9 A C		∢ □ ▶		ヨー つくで
(Moreno Maza)	CS4402-9535: Many-core Computing with Cl	UWO-CS4402-CS9535	29 / 83	(Moreno Maza)	CS4402-9535: Many-core Computing with CL	UWO-CS4402-CS9535	30 / 83
	CUDA Programming Basics				CUDA Programming Basics		
	Hearting / Delegas			Data Carica			

### GPU Memory Allocation / Release

	$\boldsymbol{c}$	1 A. 199
Data	CO	pies

Host (CPU) manages GPU memory:

- cudaMalloc (void \*\* pointer, size\_t nbytes)
- cudaMemset (void \* pointer, int value, size\_t count)
- cudaFree (void\* pointer)

```
int n = 1024;
int nbytes = 1024*sizeof(int);
int * d_a = 0;
cudaMalloc( (void**)&d_a, nbytes );
cudaMemset( d_a, 0, nbytes);
cudaFree(d_a);
```

- cudaMemcpy( void \*dst, void \*src, size\_t nbytes, enum cudaMemcpyKind direction);
  - returns after the copy is complete,

CUDA Programming Basics

- blocks the CPU thread,
- doesn't start copying until previous CUDA calls complete.
- enum cudaMemcpyKind
  - cudaMemcpyHostToDevice
  - cudaMemcpyDeviceToHost
  - cudaMemcpyDeviceToDevice
- Non-blocking memcopies are provided (more on this later)

**CUDA Programming Basics** 

#### CUDA Programming Basics

## Example kernel Source Code

## Kernel variations and output: what is in a?



### Kernel variations and utput: answers

# Code Walkthrough (1/4)

(Moreno Maza)



// walkthrough1.cu #include <stdio.h></stdio.h>
int main() { int dimx = 16; int num_bytes = dimx*sizeof(int);
int *d_a=0, *h_a=0; // device and host pointers

CS4402-9535: Many-core Computing with CL

UWO-CS4402-CS9535 36 / 83

▲□▶ ▲圖▶ ▲国▶ ▲国▶ - 国 - のへで

# Code Walkthrough (2/4)

# Code Walkthrough (3/4)

<pre>// walkthrough1.cu #include <stdio.h> int main() {     int dimx = 16;     int num_bytes = dimx*sizeof(int);     int *d_a=0, *h_a=0; // device and host pointers     h_a = (int*)malloc(num_bytes);     cudaMalloc( (void**)&amp;d_a, num_bytes );     if( 0==h_a    0==d_a )     {         printf("couldn't allocate memory\n");         return 1;     } </stdio.h></pre>	<pre>// walkthrough1.cu #include <stdio.h> int main() {     int dimx = 16;     int num_bytes = dimx*sizeof(int);     int *d_a=0, *h_a=0; // device and host pointers     h_a = (int*)malloc(num_bytes);     cudaMalloc( (void**)&amp;d_a, num_bytes );     if( 0==h_a    0==d_a )     {         printf("couldn't allocate memory\n");         return 1;     }     cudaMemset( d_a, 0, num_bytes );     cudaMemcpy( h_a, d_a, num_bytes,     cudaMemcpyDeviceToHost ); </stdio.h></pre>
Ioreno Maza) CS4402-9535: Many-core Computing with CU CUDA Programming Basics	(Moreno Maza) CS4402-9535: Many-core Computing with CU CUDA Programming Basics
Valkthrough (4/4)	Example: Shuffling Data
<pre>// walkthrough1.cu #include <stdio.h> int main() {     int dimx = 16;     int num_bytes = dimx*sizeof(int);     int num_bytes = dimx*sizeof(int);     int *d_a=0, *h_a=0; // device and host pointers     h_a = (int*)malloc(num_bytes);     cudaMalloc( (void**)&amp;d_a, num_bytes );     if( 0==h a    0==d a )</stdio.h></pre>	<pre>// Reorder values based on keys // Each thread moves one elementglobal void shuffle(int* prev_array, int*new_array, int* indices) {     int i = threadIdx.x + blockDim.x * blockIdx.x;     percential = preve error[indices[i]]; }</pre>

#### int main()

// Run grid of N/256 blocks of 256 threads each
shuffle<<< N/256, 256>>>(d\_old, d\_new, d\_ind);

Code

- ◆ □ ▶ ◆ ■ ▶ ◆ ■ ▶ ◆ ■ ● ● ● ● ●

cudaMemset( d\_a, 0, num\_bytes ); cudaMemcpy( h\_a, d\_a, num\_bytes, cudaMemcpyDeviceToHost );

for(int i=0; i<dimx; i++) printf("%d ", h\_a[i] );

printf("\n");

return 0;

free( h\_a ); cudaFree( d\_a ); {

}

CUDA Programming Basics	CUDA Programming Basics
Kernel with 2D Indexing $(1/2)$	Kernel with 2D Indexing $(2/2)$
<pre>global void kernel( int *a, int dimx, int dimy ) {     int ix = blockldx.x*blockDim.x + threadIdx.x;     int iy = blockldx.y*blockDim.y + threadIdx.y;     int idx = iy*dimx + ix;     a[idx] = a[idx]+1; }</pre>	<pre>int main() {     int dmx = 16;     int dmy = 16;     int dmy</pre>
(Woreno Waza) CS4402-9535: Wany-core Computing with CC UWO-CS4402-CS9535 41 / 83 CUDA Hardware Implementation	(Woreho Waza) CS4402-9535: Wany-core Computing with CC UWO-CS4402-CS9535 42 / 83 CUDA Hardware Implementation

### Plan

## Blocks Run on Multiprocessors

- 1 GPUs and CUDA: a Brief Introduction
- 2 CUDA Programming Model
- 3 CUDA Memory Model
- 4 CUDA Programming Basics
- 5 CUDA Hardware Implementation
- 6 CUDA Programming: Scheduling and Synchronization
- 7 CUDA Tools
- 8 Sample Programs



590

43 / 83

-

3

< 日 > < 同 > < 回 > < 回 > : < 回 > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □ > : < □

#### CUDA Hardware Implementation

### Streaming processors and multiprocessors

Streaming Processor



## Block Diagram for the G80 Family

- G80 (launched Nov 2006)
- 128 Thread Processors execute kernel threads
- Up to 12,288 parallel threads active



# Streaming Multiprocessor (1/2)

### • Processing elements:

- 8 scalar thread processors (SP)
- SM 32 GFLOPS peak at 1.35 GHz

CUDA Hardware Implementation

- 8192 32-bit registers (32KB)
- usual ops: float, int, branch, ...



# Streaming Multiprocessor (2/2)

- Hardware multithreading:
  - up to 8 blocks resident at once
  - up to 768 active threads in total

CUDA Hardware Implementation

#### • 16KB on-chip memory:

- low latency storage
- shared among threads of a block
- supports thread communication



Streaming Multiprocessor

#### CUDA Hardware Implementation

### Hardware Multithreading

#### • Hardware allocates resources to blocks:

- blocks need: thread slots, registers, shared memory
- blocks don't run until resources are available

### • Hardware schedules threads:

- hreads have their own registers
- any thread not waiting for something can run
- context switching is free every cycle

#### • Hardware relies on threads to hide latency:

• thus high parallelism is necessary for performance.

SM



## SIMT Thread Execution (1/3)

- At each clock cycle, a multiprocessor executes the same instruction on a group of threads called a warp
  - The number of threads in a warp is the warp size (32 on G80)
  - A half-warp is the first or second half of a warp.
- Within a warp, threads
  - share instruction fetch/dispatch
  - some become inactive when code path diverges
  - hardware automatically handles divergence
- Warps are the primitive unit of scheduling:
  - each active block is split nto warps in a well-defined way
  - threads within a warp are executed physically in parallel while warps and blocks are executed logically in parallel.



# SIMT Thread Execution (2/3)

### • SIMT execution is an implementation choice:

- sharing control logic leaves more space for ALUs
- largely invisible to programmer
- must be understodd for performance, not correctness
- As already mentioned, each multiprocessor processes batches of blocks, one batch after the other:
  - Active blocks = the blocks processed by one multiprocessor in one batch
  - Active threads = all the threads from the active blocks

# SIMT Thread Execution (3/3)

- The multiprocessor's registers and shared memory are split among the active threads
- Therefore, for a given kernel, the number of active blocks depends on:
  - The number of registers the kernel compiles to
  - How much shared memory the kernel requires
- If there cannot be at least one active block, the kernel fails to launch.



- GPUs and CUDA: a Brief Introduction
- **CUDA** Programming Model
- CUDA Memory Model
- **CUDA** Programming Basics
- CUDA Hardware Implementation

### CUDA Programming: Scheduling and Synchronization

- CUDA Tools
- Sample Programs

# Thread Synchronization Function

- void \_\_syncthreads();
- Synchronizes all threads in a block:
  - once all threads have reached this point, execution resumes normally.
  - this is used to avoid hazards when accessing shared memory.
- Should be used in conditional code only if the condition is uniform across the entire thread block.

	< □ >		) ଏଦ		< 🗆	<ul> <li>&lt; @ &gt; &lt; E &gt; &lt; E &gt; &lt; E</li> </ul>	৩৫৫
(Moreno Maza)	CS4402-9535: Many-core Computing with Cl	UWO-CS4402-CS9535	53 / 83	(Moreno Maza)	CS4402-9535: Many-core Computing with CL	UWO-CS4402-CS9535	54 / 83
CUDA Programming: Scheduling	and Synchronization			CUDA Programming: Schedulin	g and Synchronization		
GPU Atomic Intege	er Operations			Host Synchronizati	ion		

590

55 / 83

-

- Atomic operations on integers in global memory:
  - associative operations on signed/unsigned ints, such as
  - add, min, max, . and, or, xor.
  - they have names like atomicAdd, atomicMin, atomicAnd, ....
- Requires hardware with 1.1 compute capability
- Should be used only when strictly necessary: non-locking mechanisms should be prefered for performance consideration.

- All kernel launches are asynchronous
  - control returns to CPU immediately
  - kernel starts executing once all previous CUDA calls have completed
- Memcopies are synchronous
  - control returns to CPU once the copy is complete
  - copy starts once all previous CUDA calls have completed
- cudaThreadSynchronize()
  - host code execution resumes when all previous CUDA calls complete
- Asynchronous CUDA calls provide:
  - non-blocking memcopies (more on this later)
  - ability to overlap memcopies and kernel execution

(日) (周) (王) (王) (王)

3

200

### Example host code (recall)

#### // allocate host memory

unsigned int numBytes = N \* sizeof(float)
float\* h\_A = (float\*) malloc(numBytes);

#### // allocate device memory

float\* d\_A = 0; cudaMalloc((void\*\*)&d\_A, numbytes);

#### // copy data from host to device

cudaMemcpy(d\_A, h\_A, numBytes, cudaMemcpyHostToDevice);

#### // execute the kernel

increment\_gpu<<< N/blockSize, blockSize>>>(d\_A, b, N);

// copy data from device back to host
cudaMemcpy(h\_A, d\_A, numBytes, cudaMemcpyDeviceToHost);

### // free device memory

#### CUDA Programming: Scheduling and Synchronization

### **Device Management**

- CPU can query and select GPU devices:
  - cudaGetDeviceCount( int\* count )
  - cudaSetDevice( int device )
  - cudaGetDevice( int \*current\_device )
  - cudaGetDeviceProperties( cudaDeviceProp\* prop, int device )
  - cudaChooseDevice( int \*device, cudaDeviceProp\* prop )
- Multi-GPU setup:
  - device 0 is used by default
  - ${\scriptstyle \bullet}\,$  one CPU thread can control one GPU
  - multiple CPU threads can control the same GPU but their calls are serialized by the driver.
  - CUDA resources allocated by a CPU thread can be consumed only by CUDA calls from the same CPU thread.

cudaFree (d_A) ;	
(Moreno Maza) CS4402-9535: Many-core Computing with CL UWO-CS4402-CS9535 57 / 83	(Moreno Maza) CS4402-9535: Many-core Computing with CL UWO-CS4402-CS9535 58 / 83 CUDA Programming: Scheduling and Synchronization
CUDA Error Reporting to CPU	CUDA Event API
<ul> <li>All CUDA calls return error code: <ul> <li>except for kernel launches</li> <li>the error code type is cudaError_t</li> </ul> </li> <li>cudaError_t cudaGetLastError(void): <ul> <li>returns the code for the last error ( <i>no error</i> has also a code)</li> </ul> </li> <li>char* cudaGetErrorString(cudaError_t code): <ul> <li>returns a null-terminted character string describing the error</li> </ul> </li> <li>printf(%s\n, cudaGetErrorString( cudaGetLastError() ) );</li> </ul>	<ul> <li>Events are inserted (recorded) into CUDA call streams</li> <li>Usage scenarios: <ul> <li>measure elapsed time for CUDA calls (clock cycle precision)</li> <li>query the status of an asynchronous CUDA call</li> <li>block CPU until CUDA calls prior to the event are completed</li> </ul> </li> <li>cudaEvent_t start, stop;</li> <li>cudaEventCreate(&amp;start);</li> <li>cudaEventRecord(start, 0);</li> <li>kernel&lt;&lt;<grid, block="">&gt;&gt;();</grid,></li> <li>cudaEventSynchronize(stop);</li> <li>float et;</li> <li>cudaEventElapsedTime(&amp;et, start, stop);</li> <li>cudaEventDestroy(start); cudaEventDestroy(stop);</li> </ul>
(Marana Mara) C\$4402.0525: Many care Computing with CL LIWO C\$4402.050535 50 / 83	(Morono Mazza) CS4402 0525: Many core Computing with Cl UWO CS4402 CS0535 60 / 83

	CUDA Tools		
Plan	The nvcc compiler		
1 GPUs and CUDA: a Brief Introduction	<ul> <li>Any source file containing CUDA language extensions must be compiled with nvcc:</li> </ul>		
2 CUDA Programming Model	<ul> <li>NVCC separates code running on the host from code running on the device.</li> </ul>		
3 CUDA Memory Model	<ul> <li>Two-stage compilation:</li> </ul>		
CUDA Programming Basics	<ul> <li>First generates Parallel Thread eXecution code (PTX)</li> <li>Then produces Device-specific binary object</li> </ul>		
5 CUDA Hardware Implementation	• NVCC is a compiler driver:		
6 CUDA Programming: Scheduling and Synchronization	<ul> <li>Works by invoking all the necessary tools and compilers like cudacc, g++,</li> </ul>		
7 CUDA Tools	• An executable with CUDA code requires:		
8 Sample Programs	<ul> <li>the CUDA core library (cuda)</li> <li>the CUDA runtime library (cudart)</li> </ul>		
オロマネ 御マネ キャー・ボーン ひんの	・ロッ 《聖》 《世》 、 聞 、 うくろ		
(Moreno Maza) CS4402-9535: Many-core Computing with CL UWO-CS4402-CS9535 61 / 83 CUDA Tools	(Moreno Maza) CS4402-9535: Many-core Computing with CL UWO-CS4402-CS9535 62 / 83 CUDA Tools		
Compiling CUDA code	PTX Example (SAXPY code)		

## Compiling CUDA code



(Moreno Maza)

(Moreno Maza)

CS4402-9535: Many-core Computing with CL

64 / 83

#### CUDA Tools

## Debugging CUDA code

## Developing a CUDA program

- An executable compiled in **device emulation mode** (nvcc -deviceemu) runs completely on the host using the CUDA runtime:
  - no need of any device and CUDA driver
  - each device thread is emulated with a host thread
- However, the device emulation mode has several pitfalls:
  - emulated device threads execute sequentially, so simultaneous accesses of the same memory location by multiple threads potentially produce different results.
  - results of floating-point computations will slightly differ because of different compiler outputs, different instruction sets. etc.
  - dereferencing device pointers on the host may produce correct results in device emulation mode while generating errors in device execution mode
- In fact in the latest version of nvcc the device emulation mode is no longer supported!

Decompose the targeted application according to the many-core programming model of CUDA:

CUDA Tools

- such a program alternates serial code and vectorized code
- such that the parallel code has enough work and enough parallelism
- Write serial C code for each targeted CUDA kernel
- Sor each targeted CUDA kernel, carefully decompose the work into thread blocks:
  - this implies mapping the thred blocks to the data
  - leading to potentially delicate index caculation:
  - proving them mathematically often prevents from painful debugging!
- Verify each kernel against its C counterpart
- Debugging may lead to further decompose a kernel into smaller kernels.

	< □ >	<b>∂</b> + <b>∂</b> + <b>∂</b> + <b>∂</b> + <b>∂</b> = <b>∂</b>	5000		< □ >	<ul> <li>&lt; □ &gt; &lt; □ &gt; &lt; □ &gt; &lt; □</li> </ul>	. nac
(Moreno Maza)	CS4402-9535: Many-core Computing with Cl	UWO-CS4402-CS9535	65 / 83	(Moreno Maza)	CS4402-9535: Many-core Computing with Cl	UWO-CS4402-CS9535	66 / 83
	Sample Programs				Sample Programs		
Plan				Matrix multiplicati	on (1/16)		

- GPUs and CUDA: a Brief Introduction
- **CUDA** Programming Model
- CUDA Memory Model
- **CUDA** Programming Basics
- **CUDA** Hardware Implementation
- CUDA Programming: Scheduling and Synchronization
- CUDA Tools

### 8 Sample Programs

- The goals of this e8 xample are:
  - Understanding how to write a kernel for a non-toy example
  - Understanding how to map work (and data) to the thread blocks
  - Understanding the importance of using shared memory
- We start by writing a naive kernel for matrix multiplication which does not use shared memory.
- Then we analyze the performance of this kernel and realize that it is limited by the global memory latency.
- Finally, we present a more efficient kernel, which takes advantage of a tile decomposition and makes use of shared memory.

< ロ > < 同 > < 回 > < 回 > :

500

67 / 83

-

### Matrix multiplication (2/16)

# Matrix multiplication (3/16)

- Consider multiplying two rectangular matrices A and B with respective formats  $m \times n$  and  $n \times p$ . Define  $C = A \times B$ .
- Principle: each thread computes an element of *C* through a 2D kernel.



```
__global__ void mat_mul(float *a, float *b,
                        float *ab, int width)
{
  // calculate the row & col index of the 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 col of b
 for(int k = 0; k < width; ++k)
    result += a[row*width+k] * b[k*width+col];
  ab[row*width+col] = result;
}
```

Sample Programs



# Matrix multiplication (4/16)

- Analyze the previous CUDA kernel for multiplying two rectangular matrices A and B with respective formats  $m \times n$  and  $n \times p$ . Define  $C = A \times B$ .
- Each element of *C* is computed by one thread:
  - then each row of A is read p times and
  - each column of B is read m times, thus
  - 2 m n p reads in total for 2 m n p flops.
- Let t be an integer dividing m and p. We decompose C into  $t \times t$ tiles. If tiles are computed one after another, then:
  - (m/t)(tn)(p/t) slots are read in A
  - (p/t)(t n)(m/t) slots are read in A, thus
  - 2mnp/t reads in total for 2mnp flops.
- For a CUDA implementation, t = 16 such that each tile is computed by one thread block.

# Matrix multiplication (5/16)

- The previous explanation can be adapted to a particular GPU architecture, so as to estimate the performance of the first (naive) kernel.
- The first kernel has a global memory access to flop ratio (GMAC) of 8 Bytes / 2 ops, that is, 4 B/op.
- Suppose using a GeForce GTX 260, which has 805 GFLOPS peak performance.
- In order to reach **peak fp performance** we would need a memory bandwidth of  $GMAC \times Peak FLOPS = 3.2 TB/s$ .
- Unfortunately, we only have 112 GB/s of actual memory bandwidth (BW) on a GeForce GTX 260.
- Therefore an upper bound on the performance of our implementation is BW / GMAC = 28 GFLOPS.

71 / 83

(Moreno Maza)

#### Sample Programs

## Matrix multiplication (6/16)

# Matrix multiplication (7/16)

- The picture below illustrates our second kernel
- Each thread block computes a tile in C, which is obtained as a dot product of tile-vector of A by a tile-vector of B.
- Tile size is chosen in order to maximize data locality.



### • So a thread block computes a $t \times t$ tile of C.

Sample Programs

- Each element in that tile is a dot-prouct of a row from A and a column from B.
- We view each of these dot-products as a sum of small dot products:

$$c_{i,j} = \sum_{k=o}^{t-1} a_{i,k} b_{k,j} + \sum_{k=t}^{2t-1} a_{i,k} b_{k,j} + \cdots \sum_{k=n-1-t}^{n-1} a_{i,k} b_{k,j}$$

- Therefore we fix  $\ell$  and then compute  $\sum_{k=\ell t}^{(\ell+1)t-1} a_{i,k} b_{k,j}$  for all i,j in the working thread block.
- We do this for  $\ell = 0, 1, ..., (n/t 1)$ .
- This allows us to store the working tiles of A and B in shared memory.



# Matrix multiplication (8/16)

- We assume that A, B, C are stored in row-major layout.
- Observe that for computing a tile in C our kernel code does need to know the number of rows in A.
- It just needs to know the width (number of columns) of A and B.
- The following code fragments are taken from Example 2.

#### #define BLOCK\_SIZE 16

```
template <typename T>
__global__ void matrix_mul_ker(T* C, const T *A, const T *B,
```

```
size_t wa, size_t wb)
```

```
// Block index; WARNING: should be at most 2<sup>16</sup> - 1
int bx = blockIdx.x; int by = blockIdx.y;
```

```
// Thread index
```

```
int tx = threadIdx.x; int ty = threadIdx.y;
```

# Matrix multiplication (9/16)

- We need the position in \*A of the first element of the first working tile from A; we call it aBegin.
- We will need also the position in \*A of the last element of the last working tile from A; we call it aEnd.
- Moreover, we will need the offset between two consecutive working tiles of A; we call it aStep.

int aBegin = wa \* BLOCK\_SIZE \* by;

int aEnd = aBegin + wa - 1;

```
int aStep = BLOCK_SIZE;
```

(Moreno Maza)

э.

CS4402-9535: Many-core Computing with Cl

3

75 / 83

CS4402-9535: Many-core Computing with CL

Sample Programs	Sample Programs		
Matrix multiplication $(10/16)$	Matrix multiplication $(11/16)$		
• Similarly for <i>B</i> we have bBegin and bStep.	• The main loop starts by copying the working tiles of A and B to shared memory.		
<ul> <li>We will not need a bEnd since once we are done with a row of A, we are also done with a column of B.</li> </ul>	<pre>for(int a = aBegin, b = bBegin; a &lt;= aEnd; a += aStep, b += bSt</pre>		
<ul> <li>Finally, we initially the accumulator of the working thread; we call it Csub.</li> </ul>	snared int As[BLUCK_SIZE][BLUCK_SIZE];		
int bBegin = BLOCK_SIZE * bx;	shared int Bs[BLOCK_SIZE][BLOCK_SIZE];		
<pre>int bStep = BLOCK_SIZE * wb;</pre>	<pre>// Load the tiles from global memory to shared memory // each thread loads one element of each tile </pre>		
<pre>int Csub = 0;</pre>	As[ty][tx] = A[a + wa * ty + tx]; Bs[ty][tx] = B[b + wb * ty + tx];		
	<pre>// synchronize to make sure the matrices are loadedsyncthreads();</pre>		
くちゃ (川) (二) (二) (二) (二) (二) (二) (二) (二) (二) (二	▲□> ▲圖> ▲国> ▲国> ■ シンタの		
(Moreno Maza) CS4402-9535: Many-core Computing with CL UWO-CS4402-CS9535 77 / 83 Sample Programs	(Moreno Maza) CS4402-9535: Many-core Computing with CL UWO-CS4402-CS9535 78 / 83 Sample Programs		

```
Matrix multiplication (12/16) Matrix
```

# Matrix multiplication (13/16)

• Compute a small "dot-product" for each element in the working tile of *C*.

```
// Multiply the two tiles together
// each thread computes one element of the tile of C
for(int k = 0; k < BLOCK_SIZE; ++k) {
    Csub += As[ty][k] * Bs[k][tx];
}
// synchronize to make sure that the preceding comput:
// done before loading two new tiles of A dnd B in the
__syncthreads();
}</pre>
```

• Once computed, the working tile of C is written to global memory.

```
// Write the working tile of $C$ to global memory;
// each thread writes one element
int c = wb * BLOCK_SIZE * by + BLOCK_SIZE * bx;
C[c + wb * ty + tx] = Csub;
```

・ 「 ト ・ ヨ ト ・ ヨ ト

-

(레이 에트이 에트이

500

э

#### Sample Programs

### Matrix multiplication (14/16)

### Sample Programs

# Matrix multiplication (15/16)

- Experimentation performed on a GT200.
- **Tiling** and using **shared memory** were clearly worth the effort.



• Each thread block should have many threads:

• TILE\_WIDTH = 16 implies  $16 \times 16 = 256$  threads

### • There should be many thread blocks:

- A  $1024 \times 1024$  matrix would require 4096 thread blocks.
- Since one streaming multiprocessor (SM) can handle 768 threads, each SM will process 3 thread blocks, leading it full occupancy.

《曰》 《圖》 《臣》 《臣》

- Each thread block performs  $2 \times 256$  reads of a 4-byte float while performing  $256 \times (2 \times 16) = 8,192$  fp ops:
  - Memory bandwidth is no longer limiting factor

Sample Programs

Matrix multiplication (16/16)

(Moreno Maza)

• Effective use of different memory resources reduces the number of accesses to global memory

CS4402-9535: Many-core Computing with Cl

- But these resources are finite!
- The more memory locations each thread requires, the fewer threads an SM can accommodate.

Resource	Per GT200 SM	Full Occupancy on GT200	
Registers	16384	<= 16384 / 768 threads = <b>21 per thread</b>	
shared Memory	16KB	<= 16KB / 8 blocks = 2KB per block	