

• Purpose:	
<ul> <li>A chance to dig in deeper into a parallel programming model and explore concepts.</li> </ul>	
- Present results to work on communication of technical ideas	
<ul> <li>Write a non-trivial parallel program that combines two parallel programming languages/models. In some cases, just do two separate implementations.</li> </ul>	
- OpenMP + SSE-3	
- OpenMP + CUDA (but need to do this in separate parts of the code)	
- TBB + SSE-3	
- MPI + OpenMP	
- MPI + SSE-3	
- MPI + CUDA	
• Present results in a poster session on the last day of	
CIQSS CS4961 3	



1

## Outline

- Overview of the CUDA Programming Model for NVIDIA systems
- Presentation of basic syntax
- $\cdot$  Simple working examples
  - See http://www.cs.utah.edu/~mhall/cs6963s09
- Architecture
- Execution Model
- Heterogeneous Memory Hierarchy This lecture includes slides provided by: Wen-mei Hwu (UIUC) and David Kirk (NVIDIA) see http://courses.ecc.uiuc.edu/ece498/al1/

and Austin Robison (NVIDIA) 11/05/09

UNIVERSITY

_	Reading
•	David Kirk and Wen-mei Hwu manuscript (in progress) - http://www.toodoc.com/CUDA-textbook-by-David-Kirk- from-NVIDIA-and-Prof-Wen-mei-Hwu-pdf.html
•	CUDA 2.x Manual, particularly Chapters 2 and 4 (download from nvidia.com/cudazone)
•	Nice series from Dr. Dobbs Journal by Rob Farber - http://www.ddj.com/cpp/207200659
	11/05/09

<u>CUDA (Compute Unified Device Architecture)</u>
<ul> <li>Data-parallel programming interface to GPU</li> </ul>
<ul> <li>Data to be operated on is discretized into independent partition of memory</li> </ul>
<ul> <li>Each thread performs roughly same computation to different partition of data</li> </ul>
- When appropriate, easy to express and very efficient parallelization
<ul> <li>Programmer expresses</li> </ul>
- Thread programs to be launched on GPU, and how to launch - Data organization and movement between host and GPU - Synchronization, memory management, testing,
<ul> <li>CUDA is one of first to support <i>heterogeneous</i> architectures (more later in the semester)</li> </ul>
• CUDA environment
- Compiler, run-time utilities, libraries, emulation, performance



Compiled by nvcc

shared sdata;

\_device\_\_dfunc() { int ddata;

UNIVERSITY OF UTAH

compiler



## CUDA Programming Model: A Highly Multithreaded Coprocessor

- The GPU is viewed as a compute device that:
  - Is a coprocessor to the CPU or host
  - -Has its own DRAM (device memory)
  - Runs many threads in parallel
- Data-parallel portions of an application are executed on the device as kernels which run in parallel on many threads
- Differences between GPU and CPU threads •
  - GPU threads are extremely lightweight
    - Very little creation overhead
  - GPU needs 1000s of threads for full efficiency - Multi-core CPU needs only a few

11/05/09

UNIVERSITY



9

Known as a

cyclic data

distribution

UNIVERSITY OF UTAH

2





MAIN PROGRAM:	HOS	T FUNCTION:	
Initialization	Alloca	te memory on device for	
<ul> <li>Allocate memory on host for input and output</li> </ul>	Convi	Copy input to device	
<ul> <li>Assign random numbers to input array</li> </ul>	Set up	p grid/block	
Call host function	Call gl	lobal function	
Calculate final output from per-thread output	Сору о	device output to host	
Print result			
GLOBAL FUNCTION:		DEVICE FUNCTIO	
Thread scans subset of array elem	ients	Compare current eleme	
Call device function to compare wi	th "3"		
Compute local nacult		Return I if same, else U	

4

MAIN PROGRAM: Initialization • Allocate memory on host for input and output	#include <stdio.h> #define SIZE 16 #define BLOCKSIZE 4</stdio.h>
<ul> <li>Assign random numbers to input array</li> </ul>	int main(int argc, char **argv)
Call <i>host</i> function Calculate final output from per-thread output	{ int *in_array, *out_array;
rint result	}
11/05/09	

Initialization (OMIT) • Allocate memory on host for input and output • Assign random numbers to input array Calculate final output from per-thread output Print result Print result $d_{a} = \frac{1}{2} + \frac{1}{2} $	MAIN PROGRAM:	#include <stdio.h></stdio.h>
<ul> <li>Allocate memory on host for input and output</li> <li>Assign random numbers to input array</li> <li>Calculate final output from per-thread output</li> <li>Print result</li> <li>Assign random numbers to input array</li> <li>(int *in_array, *out_array; int sum = 0; /* initialization */ outer_compute(in_array, out_arra for (int i=0; i<blocksize; i++)="" {<br="">sum+=out_array[i]; } printf ("Result = %d\n",sum);</blocksize;></li> </ul>	Initialization (OMIT)	#define SIZE 16
<pre>* Assign random numbers to imput array Call host function Calculate final output from per-thread output Print result Print result int main(int argc, char **argv) { int *in_array, *out_array; int sum = 0; /* initialization */ outer_compute(in_array, out_arra for (int i=0; i<blocksize; i++)="" {<br="">sum+eout_array[i]; } printf ("Result = %d\n",sum);</blocksize;></pre>	Allocate memory on host for input and output	<pre>#define BLOCKSIZE 4host void outer_compute</pre>
Call host function { Calculate final output from per-thread output Print result /* initialization */ outer_compute(in_array, out_arra for (int i=0; i <blocksize; ("result='%d\n",sum);' i++)="" printf="" sum+="out_array[i];" td="" {="" }="" }<=""><td>input array</td><td>int main(int argc, char **argv)</td></blocksize;>	input array	int main(int argc, char **argv)
Calculate final output from per-thread output Print result int *in_array, *out_array; int sum = 0; /* initialization */ outer_compute(in_array, out_arra for (int i=0; i <blocksize; i++)="" {<br="">sum+=out_array[i]; } printf ("Result = %d\n",sum);</blocksize;>	Call host function	{
Print result /* initialization */ outer_compute(in_array, out_arra for (int i=0; i <blocksize; i++)="" {<br="">sum+=out_array[i]; } printf ("Result = %d\n",sum);</blocksize;>	Calculate final output from per-thread output	int *in_array, *out_array; int <i>s</i> um = 0;
outer_compute(in_array, out_arra for (int i=0; i <blocksize; i++)="" {<br="">sum+=out_array[i]; } printf ("Result = %d\n",sum);</blocksize;>	Print result	/* initialization */
for (int i=0; i <blocksize; i++)="" {<br="">sum+=out_array[i]; } printf ("Result = %d\n",sum);</blocksize;>		outer_compute(in_array, out_array)
sum+=out_array[i]; } printf ("Result = %d\n",sum);		for (int i=0; i <blocksize; i++)="" td="" {<=""></blocksize;>
} printf ("Result = %d\n",sum);		sum+=out_array[i];
printf ("Result = %d\n",sum);		}
		printf ("Result = %d\n",sum);



Host Function: Cop	<u>y Data To/From Host</u>		
HOST FUNCTION: - Allocate memory on device for copy of <i>input</i> and <i>output</i>	_host void outer_compute (int *h_in_array, int *h_out_array) { int *d_in_array, *d_out_array;	All cop	<b>)</b>
Copy input to <i>device</i> Set up grid/block Call <i>global</i> function Copy <i>device</i> output to host	cudaMalloc((void **) &d_in_array, SIZE*sizeof(int)); cudaMalloc((void **) &d_out_array, BLOCKSIZE*sizeof(int)); cudaMemcpy(d_in_array, h_in_array, SIZE*sizeof(int), cudaMemcpyHostToDevice); do computation	Col Se Cal Col	P) †    P)
11/05/09	cudaMemcpy(h_out_array.d_out_array, BOCKSIZE*Sizeof(in), cudaMemcpyDeviceToHost);		







<u>Another Example: A</u>	Adding Two Matrices	Closer Inspection of Computation and Data Partitioning
<pre>CPU C program void add_matrix_cpu(float *a, float *b, float *c, int N) { int i, j, index; for (i=0;i<n;i++) {<br="">for (j=0;j<n;j++) {<br="">index =i+j*N; c[index]=a[index]+b[index]; } } void main() {  add_matrix(a,b,c,N); }</n;j++)></n;i++)></pre>	CUDA C program global void add_matrix_gpu(float *a, float *b, float *c, int N) { int i =blockIdx.x*blockDim.x+threadIdx.x; int j=blockIdx.y*blockDim.y+threadIdx.y; int index =i+j*N; if( i < N && j <n) </n) 	<ul> <li>Define 2-d set of blocks, and 2-d set of threads per block         dim3 dimBlock(blocksize,blocksize);         dim3 dimGrid(N/dimBlock.x,N/dimBlock.y);         Each thread identifies what element of the matrix it operates on         int i=blockldx.x*blockDim.x+threadldx.x;         int j=blockldx.y*blockDim.y+threadldx.y;         int index =i+j*N;         if( i <n &&="" <="" <n)="" c[index]="a[index]+b[index];" j="" li=""> </n></li></ul>
Example source: Austin Robison, NVIDIA	U UNIVERSITY OF UTAH	11/05/09



















- Embarassingly/Pleasingly parallel computation

UNIVERSITY