CSCI 4830/7000 Advanced Computer Graphics Spring 2011



CUDA

- "Compute Unified Device Architecture"
- General purpose programming model
 - User kicks off batches of threads on the GPU
 - GPU = dedicated super-threaded, massively data parallel coprocessor
- Targeted software stack
 - Compute oriented drivers, language, and tools
- Driver for loading computation programs into GPU
 - Standalone Driver Optimized for computation
 - Interface designed for compute graphics-free API
 - Data sharing with OpenGL buffer objects
 - Guaranteed maximum download & readback spee
 - Explicit GPU memory management

Parallel Computing on a GPU

- 8-series GPUs deliver 25 to 200+ GFLC on compiled parallel C applications
 - Available in laptops, desktops, and clusters
- GPU parallelism is doubling every year
- Programming model scales transparently
- Programmable in C with CUDA tools
- Multithreaded SPMD model uses application data parallelism and thread parallelism

Tesla D870



GeForce 8800



Block IDs and Thread IDs



Courtesy: NDVIA

Matrix Multiplication Using Multiple Blocks

TILE WIDTH

- Break-up Pd into tiles
- Each block calculates one tile
 - Each thread calculates one element
 - Block size equal tile size

tv

0

2

by



A Small Example



A Small Example: Multiplication



Revised Matrix Multiplication Kernel using Multiple Blocks __global_void MatrixMulKernel(float* Md, float* Nd, float* Pd, int Width)

// Calculate the row index of the Pd element and M
int Row = blockIdx.y*TILE_WIDTH + threadIdx.y;
// Calculate the column idenx of Pd and N
int Col = blockIdx.x*TILE WIDTH + threadIdx.x;

float Pvalue = 0; // each thread computes one element of the block sub-matrix for (int k = 0; k < Width; ++k) Pvalue += Md[Row*Width+k] * Nd[k*Width+Col];

```
Pd[Row*Width+Col] = Pvalue;
}
```

Revised Step 5: Kernel Invocation (Host-side Code)

// Setup the execution configuration dim3 dimGrid(Width/TILE_WIDTH, Width/TILE_WIDTH); dim3 dimBlock(TILE_WIDTH, TILE_WIDTH);

// Launch the device computation threads!
MatrixMulKernel<<<dimGrid, dimBlock>>>(Md, Nd, Pd, Width);

CUDA Thread Block

- All threads in a block execute the same kernel program (SPMD)
- Programmer declares block:
 - Block size 1 to **512** concurrent threads
 - Block shape 1D, 2D, or 3D
 - Block dimensions in threads
- Threads have thread id numbers within block
 - Thread program uses thread id to select work and address shared data
- Threads in the same block share data and synchronize while doing their share of the work
- Threads in different blocks cannot cooperate
 - Each block can execute in any order relative to other blocs!

CUDA Thread Block



Courtesy: John Nickolls, NVIDIA

Transparent Scalability

- Hardware is free to assigns blocks to any processor at any time
 - A kernel scales across any number of parallel processors



G80 CUDA mode – A Review

- Processors execute computing threads
- New operating mode/HW interface for computing





- Or 128 (threads/block) * 6 blocks, etc.
- Threads run concurrently
 - SM maintains thread/block id #s
 - SM manages/schedules thread execution

G80 Example: Thread Scheduling

- Each Block is executed as 32thread Warps
 - An implementation decision, not part of the CUDA programming model
 - Warps are scheduling units in SM
- If 3 blocks are assigned to an SM and each block has 256 threads, how many Warps are there in an SM?
 - Each Block is divided into 256/32 = 8 Warps
 - There are 8 * 3 = 24 Warps





G80 Example: Thread Scheduling (Cont.)

- SM implements zero-overhead warp scheduling
 - Warps whose next instruction has its operands 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

G80 Block Granularity Considerations

- For Matrix Multiplication using multiple blocks, should I use 8X8, 16X16 or 32X32 blocks?
 - For 8X8, we have 64 threads per Block. Since each SM can take up to 768 threads, there are 12 Blocks. However, each SM can only take up to 8 Blocks, only 512 threads will go into each SM!
 - For 16X16, we have 256 threads per Block. Since each SM can take up to 768 threads, it can take up to 3 Blocks and achieve full capacity unless other resource considerations overrule.
 - For 32X32, we have 1024 threads per Block. Not even one can fit into an SM!

Example 49: Matrix Multiply

Download CUDA from

http://developer.nvidia.com/object/cuda_3_0_download s.html

- Default install is in /usr/local/cuda
- Usage: ex49 <Bw> <Bn>
 - Bw*Bw threads per block
 - Bn*Bn blocks
 - Matrix size Bw*Bn x Bw*Bn
- Order n³ problem

Matrix Multiply Threads × Blocks = 960 CUDA Performance: GeForce GTX 480M vs. Core i7 980X 3.33GHz



Matrix Multiply 16x16 Threads per Block CUDA Performance: GeForce 9400M vs. Core 2 Duo 2.53GHz Block



Matrix Multiply 16x16 Threads per Block CUDA Performance: GeForce GTX 480M vs. Core i7 980X 3.33GHz



Matrix Multiply 32x32 Threads per Block CUDA Performance: GeForce GTX 480M vs. Core i7 980X 3.33GHz



Application Programming Interface

- The API is an extension to the C programming language
- It consists of:
 - Language extensions
 - To target portions of the code for execution on the device
 - A runtime library split into:
 - A common component providing built-in vector types and a subset of the C runtime library in both host and device codes
 - A host component to control and access one or more devices from the host
 - A device component providing device-specific functions

Language Extensions: Built-in Variables

- dim3 gridDim;
 - Dimensions of the grid in blocks (gridDim.z unused)
- dim3 blockDim;
 - Dimensions of the block in threads
- dim3 blockIdx;
 - Block index within the grid
- dim3 threadIdx;
 - Thread index within the block

Common Runtime Component: Mathematical Functions

- pow, sqrt, cbrt, hypot
- exp, exp2, expm1
- log, log2, log10, log1p
- sin, cos, tan, asin, acos, atan, atan2
- sinh, cosh, tanh, asinh, acosh, atanh
- ceil, floor, trunc, round
- Etc.
 - When executed on the host, a given function uses the C runtime implementation if available
 - These functions are only supported for scalar types, not vector types

Device Runtime Component: Mathematical Functions

Some mathematical functions (e.g. sin(x)) have a less accurate, but faster device-only version (e.g. sin(x))



– **___exp**

Host Runtime Component

- Provides functions to deal with:
 - Device management (including multi-device systems)
 - Memory management
 - Error handling
- Initializes the first time a runtime function is called
- A host thread can invoke device code on only one device
 - Multiple host threads required to run on multiple devices

Device Runtime Component: Synchronization Function • void ______syncthreads ();

- Synchronizes all threads in a block
- Once all threads have reached this point, execution resumes normally
- Used to avoid RAW / WAR / WAW hazards when accessing shared or global memory
- Allowed in conditional constructs only if the conditional is uniform across the entire thread block

G80 Implementation of CUDA Memories

- Each thread can:
 - Read/write per-thread registers
 - Read/write per-thread local memory
 - Read/write per-block
 shared memory
 - Read/write per-grid
 global memory
 - Read/only per-grid
 constant memory



CUDA Variable Type Qualifiers

Variable declaration	Memory	Scope	Lifetime
	local	thread	thread
	shared	block	block
device int GlobalVar;	global	grid	application
	constant	grid	application

- <u>device</u> is optional when used with
 <u>local</u>, <u>shared</u>, or <u>constant</u>
- Automatic variables without any qualifier reside in a register
 - Except arrays that reside in local memory

Where to Declare Variables?



Variable Type Restrictions

- Pointers can only point to memory allocated or declared in global memory:
 - Allocated in the host and passed to the kernel:

__global___void KernelFunc(float* ptr)
- Obtained as the address of a global variable:
float* ptr = &GlobalVar;

A Common Programming Strategy

- Global memory resides in device memory (DRAM) - much slower access than shared memory
- So, a profitable way of performing computation on the device is to tile data to take advantage of fast shared memory:
 - Partition data into subsets that fit into shared memory
 - Handle each data subset with one thread block by:
 - Loading the subset from global memory to shared memory, using multiple threads to exploit memorylevel parallelism
 - Performing the computation on the subset from shared

A Common Programming Strategy (Cont.)

- Constant memory also resides in device memory (DRAM) - much slower access than shared memory
 - But... cached!
 - Highly efficient access for read-only data
- Carefully divide data according to access patterns
 - R/Only \rightarrow constant memory (very fast if in cache)
 - R/W shared within Block \rightarrow shared memory (very fast)
 - R/Wtwitthinference of the thread have here and the difference of the second of the s
 - R/W inputs/results \rightarrow global memory (very slow)

GPU Atomic Integer Operations

- Atomic operations on integers in global memory:
 - Associative operations on signed/unsigned ints
 - add, sub, min, max, …
 - and, or, xor
 - Increment, decrement
 - Exchange, compare and swap
- Requires hardware with compute capability 1.1 and above.

Review: Matrix Multiplication Kernel using Multiple Blocks __global_void MatrixMulKernel(float* Md, float* Nd, float* Pd, int Width)

// Calculate the row index of the Pd element and M
int Row = blockIdx.y*TILE_WIDTH + threadIdx.y;
// Calculate the column idenx of Pd and N
int Col = blockIdx.x*TILE WIDTH + threadIdx.x;

float Pvalue = 0; // each thread computes one element of the block sub-matrix for (int k = 0; k < Width; ++k) Pvalue += Md[Row*Width+k] * Nd[k*Width+Col];

```
Pd[Row*Width+Col] = Pvalue;
}
```

How about performance on G80?

- All threads access global memory for their input matrix elements
 - Two memory accesses (8 b)tes) per floating point multiply-add
 - 4B/s of memory bandwidth/FLOPS
 - 4*346.5 = 1386 GB/s required to achieve peak FLOP rating
 - 86.4 GB/s limits the code at 21.6 GFLOPS
- The actual code runs at about 15 GFLOPS
- Need to drastically cut down memory accesses to get closer to the peak 346.5 GFLOPS



Idea: Use Shared Memory to reuse global memory data

- Each input element is read by Width threads.
- Load each element into Shared Memory and have several threads use the local version to reduce the memory bandwidth
 - Tiled algorithms



Tiled Multiply

 Break up the execution of the kernel into phases so that the data accesses in each phase is focused on one subset (tile) of Md and Nd

0

2

by

tv

TILE WIDTH

TILE WIDT

TILE WIDTH



A Small Example



Every Md and Nd Element is used exactly twice in generating a 2X2 tile of P

	P _{0,0}	P _{1,0}	P _{0,1}	P _{1,1}
1	thread _{0,0}	thread _{1,0}	thread _{0,1}	thread _{1,1}
	NA * NI		NA * NI	
	IVI _{0,0} " IN _{0,0}	M _{0,0} " N _{1,0}	IM _{0,1} " IN _{0,0}	M _{0,1} " N _{1,0}
Access order	M _{1,0} N _{0,1}	N _{1,1}	M _{1,1} * N _{0,1}	M _{1,1} * N _{1,1}
	M _{2,0} * N _{0,2}	M _{2,0} * N _{1,2}	M _{2,1} * N _{0,2}	M _{2,1} * N _{1,2}
	M _{3,0} * N _{0,3}	M _{3,0} * N _{1,3}	M _{3,1} * N _{0,3}	M _{3,1} * N _{1,3}

Breaking Md and Nd into Tiles

- Break up the inner product loop of each thread into phases
- At the beginning of each phase, load the Md and Nd elements that everyone needs during the phase into shared memory
- Everyone access the Md and Nd elements from the shared memory during the phase

	$ \begin{bmatrix} 1 & Nd_{0,0} & Nd_{1,0} \\ 1 & Nd_{0,1} & Nd_{1,1} \\ 1 & Nd_{0,2} & Nd_{1,2} \\ 1 & Nd_{0,3} & Nd_{1,3} \end{bmatrix} $
Md _{0,0} Md _{1,0} Md _{2,0} Md _{3,0}	$\begin{array}{c c} Pd_{8,0} & Pd_{2,0} \\ Pd_{8,0} & Pd_{2,0} \\ \end{array} Pd_{3,0} \end{array}$
$\mathrm{Md}_{_{0,1}}\mathrm{Md}_{_{1,1}}\mathrm{Md}_{_{2,1}}\mathrm{Md}_{_{3,1}}$	$Pd_{0,1} Pd_{1,1} Pd_{2,1} Pd_{3,1}$
	$\begin{array}{ c c c c } Pd_{0,2} & Pd_{1,2} & Pd_{2,2} & Pd_{3,2} \end{array}$
	$Pd_{0,3} Pd_{1,3} Pd_{2,3} Pd_{3,3}$

Each phase of a Thread Block uses one tile from Md and one from Nd

	Phase 1			Phase 2		
T _{0,0}	Md _{₀,₀} ↓ Mds _{₀,₀}	Nd _{₀,₀} ↓ Nds _{₀,₀}	$PValue_{0,0} += Mds_{0,0}*Nds_{0,0} + Mds_{1,0}*Nds_{0,1}$	Md _{2,0} ↓ Mds _{0,0}	Nd _{0,2} ↓ Nds _{0,0}	PValue _{0,0} += Mds _{0,0} *Nds _{0,0} + Mds _{1,0} *Nds _{0,1}
Τ _{1,0}	Md _{1,0} ↓ Mds _{1,0}	Nd _{1,0} ↓ Nds _{1,0}	PValue _{1,0} += Mds _{0,0} *Nds _{1,0} + Mds _{1,0} *Nds _{1,1}	Md _{3,0} ↓ Mds _{1,0}	Nd _{1,2} ↓ Nds _{1,0}	PValue _{1,0} += Mds _{0,0} *Nds _{1,0} + Mds _{1,0} *Nds _{1,1}
T _{0,1}	Md _{₀,1} ↓ Mds _{₀,1}	Nd _{0,1} ↓ Nds _{0,1}	PoValue _{0,1} += Mds _{0,1} *Nds _{0,0} + Mds _{1,1} *Nds _{0,1}	Md _{2,1} ↓ Mds _{0,1}	Nd _{₀,3} ↓ Nds _{₀,1}	PdValue _{0,1} += Mds _{0,1} *Nds _{0,0} + Mds _{1,1} *Nds _{0,1}
T _{1,1}	Md _{1,1} ↓ Mds _{1,1}	Nd _{1,1} ↓ Nds _{1,1}	PdValue _{1,1} += Mds _{0,1} *Nds _{1,0} + Mds _{1,1} *Nds _{1,1}	Md _{3,1} ↓ Mds _{1,1}	Nd _{1,3} ↓ Nds _{1,1}	PdValue _{1,1} += Mds _{0,1} *Nds _{1,0} + Mds _{1,1} *Nds _{1,1}

Kernel

```
global void MatrixMulKernel(float* Md, float* Nd, float* Pd, int Width)
{
1.
    shared float Mds[TILE WIDTH][TILE WIDTH];
2.
    shared float Nds[TILE WIDTH][TILE WIDTH];
   int bx = blockIdx.x; int by = blockIdx.y;
3.
4. int tx = threadIdx.x; int ty = threadIdx.y;
// Identify the row and column of the Pd element to work on
5. int Row = by * TILE WIDTH + ty;
    int Col = bx * TILE WIDTH + tx;
6.
    float Pvalue = 0;
7.
// Loop over the Md and Nd tiles required to compute the Pd element
8. for (int m = 0; m < Width/TILE WIDTH; ++m) {
   Collaborative loading of Md and Nd tiles into shared memory
//
9.
        Mds[ty][tx] = Md[Row*Width + (m*TILE WIDTH + tx)];
        Nds[ty][tx] = Nd[(m*TILE WIDTH + ty)*Width + Col];
٠
        ____syncthreads();
٠
     for (int k = 0; k < TILE WIDTH; ++k)
12.
          Pvalue += Mds[ty][k] * Nds[k][tx];
14.
       syncthreads();
15. Pd[Row*Width + Col] = Pvalue;
}
```

CUDA Code – Kernel Execution Configuration

First-order Size Considerations in G80

- Each thread block should have many threads
 TILE WIDTH of 16 gives 16*16 = 256 threads
- There should be many thread blocks
 - A 1024*1024 Pd gives 64*64 = 4096 Thread Blocks
 - TILE_WIDTH of 16 gives each SM 3 blocks, 768 threads (full capacity)
- Each thread block perform 2*256 = 512 float loads from global memory for 256 * (2*16) = 8,192 mul/add operations.
 - Memory bandwidth no longer a limiting factor

Tiled Multiply

by

k.

TILE WIDTH

m

- Each block computes one square sub-matrix Pd_{sub} of size TILE_WIDTH
- Each thread computes one element of Pd_{sub}

tv

TILE WIDTH

0

2

by



G80 Shared Memory and Threading

- Each SM in G80 has 16KB shared memory
 - SM size is implementation dependent!
 - For TILE_WIDTH = 16, each thread block uses 2*256*4B = 2KB of shared memory.
 - The shared memory can potentially have up to 8 Thread Blocks actively executing
 - This allows up to 8*512 = 4,096 pending loads. (2 per thread, 256 threads per block)
 - The threading model limits the number of thread blocks to 3 so shared memory is not the limiting factor here
 - The next TILE_WIDTH 32 would lead to 2*32*32*4B= 8KB shared memory usage per thread block, allowing only up to two thread blocks active at the same time
- Using 16x16 tiling, we reduce the accesses to the global memory by a factor of 16
 - The 86.4B/s bandwidth can now support (86.4/4)*16 = 347.6 GFLOPS!

```
illed Matrix Multiplication
 _global___ void MatrixMulKernel(float* Md, Floatend, Int Width)
{
1.
    shared float Mds[TILE WIDTH][TILE WIDTH];
   shared float Nds[TILE WIDTH][TILE WIDTH];
3.
   int bx = blockIdx.x; int by = blockIdx.y;
   int tx = threadIdx.x; int ty = threadIdx.y;
4.
// Identify the row and column of the Pd element to work on
5. int Row = by * TILE WIDTH + ty;
   int Col = bx * TILE WIDTH + tx;
6.
   float Pvalue = 0;
7.
// Loop over the Md and Nd tiles required to compute the Pd element
8. for (int m = 0; m < Width/TILE WIDTH; ++m) {
// Collaborative loading of Md and Nd tiles into shared memory
       Mds[ty][tx] = Md[Row*Width + (m*TILE WIDTH + tx)];
9.
       Nds[ty][tx] = Nd[(m*TILE WIDTH + ty)*Width + Col];
•
       syncthreads();
٠
    for (int k = 0; k < TILE WIDTH; ++k)
12.
         Pvalue += Mds[ty][k] * Nds[k][tx];
•
      syncthreads();
14.
15. Pd[Row*Width + Col] = Pvalue;
}
```

Tiling Size Effects



Summary- Typical Structure of a CUDA Program

- Global variables declaration
 - <u>host</u>
 - ____device__... _global__, __constant__, __texture__
- Function prototypes
 - __global__ void kernelOne(...)
 - float handyFunction(...)
- Main ()
 - allocate memory space on the device cudaMalloc(&d_GlblVarPtr, bytes)
 - transfer data from host to device cudaMemCpy(d_GlbIVarPtr, h_Gl...)
 - execution configuration setup
 - kernel call kernelOne<<<execution configuration>>>(args...);
 - transfer results from device to host cudaMemCpy(h_GlblVarPtr,...)
 - optional: compare against golden (host computed) solution
- Kernel void kernelOne(type args,...)
 - variables declaration _local_, _shared_
 - automatic variables transparently assigned to registers or local memory
 - syncthreads()...
- Other functions
 - float handyFunction(int inVar...);