CS/EE 217 GPU Architecture and Parallel Programming

Lectures 4 and 5: Memory Model and Locality

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2013

Objective

- To learn to efficiently use the important levels of the CUDA memory hierarchy
 - Registers, shared memory, global memory
 - Tiled algorithms and barrier synchronization

The Von-Neumann Model



Going back to the program

• Every instruction needs to be fetched from memory, decoded, then executed.

- The decode stage typically accesses register file

- Instructions come in three flavors: Operate, Data transfer, and Program Control Flow.
- An example instruction cycle is the following:

Fetch | Decode | Execute | Memory

Operate Instructions

- Example of an operate instruction: ADD R1, R2, R3
- Instruction cycle for an operate instruction: Fetch | Decode | Execute | Memory

Memory Access Instructions

- Examples of memory access instruction: LDR R1, R2, #2 STR R1, R2, #2
- Instruction cycle for an operate instruction: Fetch | Decode | Execute | Memory

Registers vs Memory

- Registers are "free"
 - No additional memory access instruction
 - Very fast to use, however, there are very few of them
- Memory is expensive (slow), but very large



Programmer View of CUDA Memories

• Each thread can:

- Read/write per-thread
 registers (~1 cycle)
- Read/write per-block
 shared memory (~5
 cycles)
- Read/write per-grid
 global memory (~500
 cycles)
- Read/only per-grid
 constant memory (~5
 cycles with caching)



Shared Memory in CUDA

- A special type of memory whose contents are explicitly declared and used in the source code
 - Located in the processor
 - Accessed at much higher speed (in both latency and throughput)
 - Still accessed by memory access instructions
 - Commonly referred to as scratchpad memory in computer architecture

CUDA Variable Type Qualifiers

	Variable decl	aration	Memory	Scope	Lifetime
		<pre>int LocalVar;</pre>	register	thread	thread
device	_shared	<pre>int SharedVar;</pre>	shared	block	block
device		<pre>int GlobalVar;</pre>	global	grid	application
device	_constant	<pre>int ConstantVar;</pre>	constant	grid	application

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



__global__ void MatrixMulKernel(float* d_M, float* d_N, float* d_P, int Width)

1. _____shared____float ds_M[TILE_WIDTH][TILE_WIDTH];
2. _____shared____float ds_N[TILE_WIDTH][TILE_WIDTH];

A Common Programming Strategy

- Global memory resides in device memory (DRAM)
 slow access
- So, a profitable way of performing computation on the device is to tile input 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 memory-level parallelism
 - Performing the computation on the subset from shared memory; each thread can efficiently multi-pass over any data element
 - Copying results from shared memory to global memory

Matrix-Matrix Multiplication using Shared Memory

Base Matrix Multiplication Kernel

_global___ void MatrixMulKernel(float* d_M, float* d_N, float* d_P, 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 += d_M[Row*Width+k]* d_N[k*Width+Col];</pre>

```
d P[Row*Width+Col] = Pvalue;
```

How about performance on Fermi?

- All threads access global memory for their input matrix elements
 - Two memory accesses (8 bytes) per floating point multiply add
 - 4B/s of memory bandwidth/ FLOPS
 - 4*1,000 = 4,000 GB/s required to achieve peak FLOP rating
 - 150 GB/s limits the code at 37.5 GFLOPS
- The actual code runs at about 25 GFLOPS
- Need to drastically cut down memory accesses to get closer to the peak 1,000 GFLOPS



Shared Memory Blocking Basic Idea



17

Basic Concept of Blocking/Tiling

- In a congested traffic system, significant reduction of vehicles can greatly improve the delay seen by all vehicles
 - Carpooling for commuters
 - Blocking/Tiling for global memory accesses
 - drivers = threads,
 - cars = data





Some computations are more challenging to block/tile than others.

- Some carpools may be easier than others
 - More efficient if neighbors are also classmates or coworkers
 - Some vehicles may be more suitable for carpooling
- Similar variations exist in blocking/tiling





Carpools need synchronization.

• Good – when people have similar schedule



Same with Blocking/Tiling

Good – when threads have similar access timing



Bad – when threads have very different timing

Outline of Technique

- Identify a block/tile of global memory content that are accessed by multiple threads
- Load the block/tile from global memory into onchip memory
- Have the multiple threads to access their data from the on-chip memory
- Move on to the next block/tile

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 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



Loading a Tile

- All threads in a block participate
 - Each thread loads one Md element and one Nd element in based tiled code
- Assign the loaded element to each thread such that the accesses within each warp is coalesced (more later).



SM

$M_{0.0}$	M_	$\mathbf{M}_{\mathbf{a}}$	\mathbf{M}_{a}	\mathbf{M}_{aa}	M_{01}
0,0	0,1	0,2	0,5	0,0	0,1
M _{1,0}	M _{1,1}	M _{1,2}	M _{1,3}	M _{1,0}	M _{1,1}
M _{2,0}	M _{2,1}	M _{2,2}	M _{2,3}		
M _{3,0}	M _{3,1}	M _{3,2}	M _{3,3}		

P _{0,0}	P _{0,1}	P _{0,2}	P _{0,3}
P _{1,0}	P _{1,1}	P _{1,2}	P _{1,3}
P _{2,0}	P _{2,1}	P _{2,2}	P _{2,3}
P _{3,0}	P _{3,1}	P _{3,2}	P _{3,3}

SM

N _{0,0}	N _{0,1}	N _{0,2}	N _{0,3}
N _{1,0}	N _{1,1}	N _{1,2}	N _{1,3}
N _{2,0}	N _{2,1}	N _{2,2}	N _{2,3}
N _{3,0}	N _{3,1}	N _{3,2}	N _{3,3}

M _{0,0}	M _{0,1}	M _{0,2}	M _{0,3}
M _{1,0}	M _{1,1}	M _{1,2}	M _{1,3}
M _{2,0}	M _{2,1}	M _{2,2}	M _{2,3}



N _{0,0}	N _{0,1}	N _{0,2}	N _{0,3}
N _{1,0}	N _{1,1}	N _{1,2}	N _{1,3}
N _{2,0}	N _{2,1}	N _{2,2}	N _{2,3}
N _{3,0}	N _{3,1}	N _{3,2}	N _{3,3}

M _{0,0}	M _{0,1}	M _{0,2}	M _{0,3}
M _{1,0}	M _{1,1}	M _{1,2}	M _{1,3}
M _{2,0}	M _{2,1}	M _{2,2}	M _{2,3}





SM



P _{0,0}	P _{0,1}	P _{0,2}	P _{0,3}
P _{1,0}	P _{1,1}	P _{1,2}	P _{1,3}
P _{2,0}	P _{2,1}	P _{2,2}	P _{2,3}

N _{0,0}	N _{0,1}	N _{0,2}	N _{0,3}
N _{1,0}	N _{1,1}	N _{1,2}	N _{1,3}
N _{2,0}	N _{2,1}	N _{2,2}	N _{2,3}
N _{3,0}	N _{3,1}	N _{3,2}	N _{3,3}

M _{0,0}	M _{0,1}	M _{0,2}	M _{0,3}
M _{1,0}	M _{1,1}	M _{1,2}	M _{1,3}
M _{2,0}	M _{2,1}	M _{2,2}	M _{2,3}



Barrier Synchronization

- An API function call in CUDA
 - __syncthreads()
- Best used to coordinate tiled algorithms
 To oncure that all elements of a tile are leaded
 - To ensure that all elements of a tile are loaded
 - To ensure that all elements of a tile are consumed



Figure 4.11 An example execution timing of barrier synchronization.

Loading an Input Tile , bx , 2



Loading an Input Tile ²

E WIDTH-1

m

k

Accessing tile 1 in 2D indexing: M[Row][1*TILE_WIDTH+tx] N[1*TILE_WIDTH+ty][Col]



Loading Input Tile m

However, M and N are dynamically allocated and can only use 1D indexing:

M[Row][m*TILE_WIDTH+tx] M[Row*Width + m*TILE_WIDTH + tx]

N[m*TILE_WIDTH+ty][Col] N[(m*TILE_WIDTH+ty) * Width + Col]

Row

d_M

m*TILE_WIDTH



Tiled Matrix Multiplication Kernel

```
global void MatrixMulKernel(float* d_M, float* d_N, float* d_P, int Width)
{
1.
    shared float ds M[TILE WIDTH][TILE WIDTH];
      shared float ds N[TILE_WIDTH][TILE_WIDTH];
    int bx = blockIdx.x; int by = blockIdx.y;
3.
   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;
6. int Col = bx * TILE WIDTH + tx;
   float Pvalue = 0;
7.
  Loop over the Md and Nd tiles required to compute the Pd element
//
     for (int m = 0; m < Width/TILE WIDTH; ++m) {</pre>
8.
  Coolaborative loading of Md and Nd tiles into shared memory
//
9.
       ds M[ty][tx] = d M[Row*Width + m*TILE WIDTH+tx];
10.
     ds N[ty][tx] = d N[Col+(m*TILE WIDTH+ty)*Width];
      syncthreads();
11.
12.
      for (int k = 0; k < TILE WIDTH; ++k)
         Pvalue += ds M[ty][k] * ds N[k][tx];
13.
      synchthreads();
14.
15.}
16.
      d P[Row*Width+Col] = Pvalue;
```

Compare with the Base Kernel

_global___ void MatrixMulKernel(float* d_M, float* d_N, float* d_P, 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 += d_M[Row*Width+k]* d_N[k*Width+Col];</pre>

```
d P[Row*Width+Col] = Pvalue;
```

First-order Size Considerations

- Each thread block should have many threads
 - TILE_WIDTH of 16 gives 16*16 = 256 threads
 - TILE_WIDTH of 32 gives 32*32 = 1024 threads
- For 16, each block performs 2*256 = 512 float loads from global memory for 256 * (2*16) = 8,192 mul/add operations.
- For 32, each block performs 2*1024 = 2048 float loads from global memory for 1024 * (2*32) = 65,536 mul/add operations

Shared Memory and Threading

- Each SM in Fermi has 16KB or 48KB shared memory*
 - SM size is implementation dependent!
 - For TILE_WIDTH = 16, each thread block uses 2*256*4B = 2KB of 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 next TILE_WIDTH 32 would lead to 2*32*32*4B= 8KB shared memory usage per thread block, allowing 2 or 6 thread blocks active at the same time
- Using 16x16 tiling, we reduce the accesses to the global memory by a factor of 16
 - The 150GB/s bandwidth can now support (150/4)*16 = 600GFLOPS!

*Configurable vs L1, total 64KB

Boundary conditions

- What to do if the matrix size is not a multiple of width?
 - Tricky problem, lets work through an example
- Too many boundary checks can cause control divergence and overhead
- Something that you have to work through for lab
 2

- I'll start a piazza discussion on the topic

Device Query

• Number of devices in the system

int dev_count; cudaGetDeviceCount(&dev_count);

Capability of devices

cudaDeviceProp dev_prop; for (i = 0; i < dev_count; i++) {

cudaGetDeviceProperties(&dev_prop, i);

// decide if device has sufficient resources and capabilities

- cudaDeviceProp is a built-in C structure type
 - dev_prop.dev_prop.maxThreadsPerBlock
 - dev_prop.sharedMemoryPerBlock

Summary- Typical Structure of a CUDA Program

- Global variables declaration
 - __host__
 - __device__... _global__, __constant__, __texture__
- Function prototypes
 - __global__ void kernelOne(...)
 - float handyFunction(...)
- Main ()
 - allocate memory space on the device cudaMalloc(&d_GlbIVarPtr, bytes)
 - transfer data from host to device cudaMemCpy(d_GlbIVarPtr, h_GI...)
 - 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 auto, __shared__
 - automatic variables transparently assigned to registers
 - syncthreads()...
- Other functions
 - float handyFunction(int inVar...);

repeat

needed

as

ANY MORE QUESTIONS? READ CHAPTER 5!