# L6: Memory Hierarchy Optimization IV, Bandwidth Optimization

CS6963

---

# Administrative

- Next assignment available
  - Goals of assignment:
    - simple memory hierarchy management
    - block-thread decomposition tradeoff
  - Due Tuesday, Feb. 8, 5PM
  - Use handin program on CADE machines
    - "handin cs6963 lab2 <probfile>"
- Project proposals due Wednesday, March 9
  - Questions/discussion
- Mailing lists
  - cs6963s11-discussion@list.eng.utah.edu
    - Please use for all questions suitable for the whole class
    - Feel free to answer your classmates questions!
  - cs6963s1-teach@list.eng.utah.edu
    - Please use for questions to Sriram and me

CS6963    L5: Memory Hierarchy, IV    2

---

# Faculty Project Suggestions

- Mike Kirby:
  - Immersed boundary method
  - Spectral element library co-processing with CPU
  - Hybridized discontinuous Galerkin method
- Kris Sikorski:
  - Reliable algorithms for summation of large data sets
- Matt Might:
  - Containment analysis (program analysis)

---

# Project Proposal (due 3/9)

- Proposal Logistics:
  - Significant implementation, worth 55% of grade
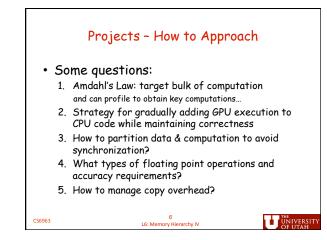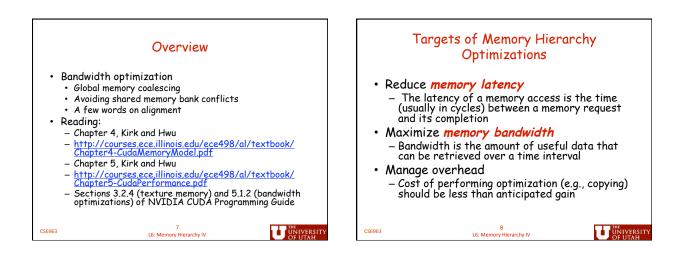  - Each person turns in the proposal (should be same as other team members)
- Proposal:
  - 3-4 page document (11pt, single-spaced)
  - Submit with handin program:
    "handin cs6963 prop <pdf-file>"

CS6963    4    L6: Memory Hierarchy IV

1

## Content of Proposal

I.   Team members: Name and a sentence on expertise for each member
II.  Problem description
   - What is the computation and why is it important?
   - Abstraction of computation: equations, graphic or pseudo-code, no more than 1 page
III. Suitability for GPU acceleration
   - Amdahl's Law: describe the inherent parallelism. Argue that it is close to 100% of computation. Use measurements from CPU execution of computation if possible.
   - Synchronization and Communication: Discuss what data structures may need to be protected by synchronization, or communication through host.
   - Copy Overhead: Discuss the data footprint and anticipated cost of copying to/from host memory.
IV.  Intellectual Challenges
   - Generally, what makes this computation worthy of a project?
   - Point to any difficulties you anticipate at present in achieving high speedup

CS6963                            5
                        L6: Memory Hierarchy IV                    THE UNIVERSITY OF UTAH

## Projects – How to Approach

- Some questions:
  1. Amdahl's Law: target bulk of computation and can profile to obtain key computations…
  2. Strategy for gradually adding GPU execution to CPU code while maintaining correctness
  3. How to partition data & computation to avoid synchronization?
  4. What types of floating point operations and accuracy requirements?
  5. How to manage copy overhead?

CS6963                            6
                        L6: Memory Hierarchy IV                    THE UNIVERSITY OF UTAH

## Overview

- Bandwidth optimization
  - Global memory coalescing
  - Avoiding shared memory bank conflicts
  - A few words on alignment
- Reading:
  – Chapter 4, Kirk and Hwu
  – http://courses.ece.illinois.edu/ece498/al/textbook/Chapter4-CudaMemoryModel.pdf
  – Chapter 5, Kirk and Hwu
  – http://courses.ece.illinois.edu/ece498/al/textbook/Chapter5-CudaPerformance.pdf
  – Sections 3.2.4 (texture memory) and 5.1.2 (bandwidth optimizations) of NVIDIA CUDA Programming Guide

CS6963                            7
                        L6: Memory Hierarchy IV                    THE UNIVERSITY OF UTAH

## Targets of Memory Hierarchy Optimizations

- Reduce *memory latency*
  – The latency of a memory access is the time (usually in cycles) between a memory request and its completion
- Maximize *memory bandwidth*
  – Bandwidth is the amount of useful data that can be retrieved over a time interval
- Manage overhead
  – Cost of performing optimization (e.g., copying) should be less than anticipated gain

CS6963                            8
                        L6: Memory Hierarchy IV                    THE UNIVERSITY OF UTAH

2

## Optimizing the Memory Hierarchy on GPUs, Overview

- Device memory access times non-uniform so *data placement* significantly affects performance.
  - But controlling data placement may require additional copying, so consider overhead.
- Optimizations to increase memory bandwidth. Idea: maximize utility of each memory access.
  - *Coalesce* global memory accesses
  - *Avoid memory bank conflicts* to increase memory access parallelism
  - *Align* data structures to address boundaries

CS6963     9     L6: Memory Hierarchy IV

---

## Data Location Impacts Latency of Memory Access

- Registers
  - Can load in current instruction cycle
- Constant or Texture Memory
  - In cache? Single address can be loaded for half-warp per cycle
  - O/W, global memory access
- Global memory
- Shared memory
  - Single cycle if accesses can be done in parallel

CS6963     10     L6: Memory Hierarchy IV

---

## Introduction to Memory System

- Recall execution model for a multiprocessor
  - **Scheduling unit:** A "warp" of threads is issued at a time (32 threads in current chips)
  - **Execution unit:** Each cycle, 8 "cores" or SPs are executing (32 cores in a Fermi)
  - **Memory unit:** Memory system scans a *"half warp"* or 16 threads for data to be loaded; (full warp for Fermi)

CS6963     11     L6: Memory Hierarchy IV

---

## Global Memory Accesses

- Each thread issues memory accesses to data types of varying sizes, perhaps as small as 1 byte entities
- Given an address to load or store, memory returns/updates "segments" of either 32 bytes, 64 bytes or 128 bytes
- Maximizing bandwidth:
  - Operate on an *entire* 128 byte segment for each memory transfer

CS6963     12     L6: Memory Hierarchy IV

## Understanding Global Memory Accesses

Memory protocol for compute capability 1.2 and 1.3* (CUDA Manual 5.1.2.1 and Appendix A.1)

- Start with memory request by smallest numbered thread. Find the memory segment that contains the address (32, 64 or 128 byte segment, depending on data type)
- Find other active threads requesting addresses within that segment and *coalesce*
- Reduce transaction size if possible
- Access memory and mark threads as "inactive"
- Repeat until all threads *in half-warp* are serviced

*Includes Tesla and GTX platforms as well as new Linux machines!

CS6963
13
L6: Memory Hierarchy IV

THE UNIVERSITY OF UTAH

## Protocol for most systems (including lab6 machines) even more restrictive

- For compute capability 1.0 and 1.1
  - Threads must access the words in a segment in sequence
  - The kth thread must access the kth word
  - *Alignment to the beginning of a segment becomes a very important optimization!*

CS6963
14
L6: Memory Hierarchy IV

THE UNIVERSITY OF UTAH

## Memory Layout of a Matrix in C



Consecutive threads will access different rows in memory.

Each thread will require a different memory operation.

Odd: But this is the RIGHT layout for a conventional multi-core!

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
ECE 498AL, University of Illinois, Urbana-Champaign
15
L6: Memory Hierarchy IV

THE UNIVERSITY OF UTAH

## Memory Layout of a Matrix in C



Each thread in a half-warp (assuming rows of 16 elements) will access consecutive memory locations.

GREAT! All accesses are coalesced.

With just a 4x4 block, we may need 4 separate memory operations to load data for a half-warp.

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
ECE 498AL, University of Illinois, Urbana-Champaign
16
L6: Memory Hierarchy IV

THE UNIVERSITY OF UTAH

## How to find out compute capability

See Appendix A.1 in NVIDIA CUDA Programming Guide to look up your device.

Also, recall "deviceQuery" in SDK to learn about features of installed device.

Older CADE machines are all Compute Capability 1.0 or 1.1.

Linux lab, most CADE machines and Tesla cluster are Compute Capability 1.2 and 1.3.

Fermi machines are 2.x.

CS6963                          17
                    L6: Memory Hierarchy IV

---

## Alignment

- Addresses accessed within a half-warp may need to be *aligned* to the beginning of a segment to enable coalescing
  – An aligned memory address is a multiple of the memory segment size
  – In compute 1.0 and 1.1 devices, address accessed by lowest numbered thread must be aligned to beginning of segment for coalescing
  – In future systems, sometimes alignment can reduce number of accesses

CS6963                          18
                    L6: Memory Hierarchy IV

---

## More on Alignment

- Objects allocated statically or by cudaMalloc begin at aligned addresses
  – But still need to think about index expressions
- May want to align structures

```
struct __align__(8) {          struct __align__(16) {
  float a;                       float a;
  float b;                       float b;
};                               float c;
                               };
```

CS6963                          19
                    L6: Memory Hierarchy IV

---

## What Can You Do to Improve Bandwidth to Global Memory?

- Think about spatial reuse and access patterns across threads
  – May need a different computation & data partitioning
  – May want to rearrange data in shared memory, even if no temporal reuse (transpose example)
  – Similar issues, but much better in future hardware generations

CS6963                          20
                    L6: Memory Hierarchy IV

5

2/2/11

## Bandwidth to Shared Memory: Parallel Memory Accesses

- Consider each thread accessing a different location in shared memory
- Bandwidth maximized if each one is able to proceed *in parallel*
- Hardware to support this
  - *Banked memory:* each bank can support an access on every memory cycle

CS6963
21
L6: Memory Hierarchy IV

## How addresses map to banks on G80

- Each bank has a bandwidth of 32 bits per clock cycle
- Successive 32-bit words are assigned to successive banks
- G80 has 16 banks
  - So bank = address % 16
  - Same as the size of a half-warp
    - No bank conflicts between different half-warps, only within a single half-warp

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
ECE 498AL, University of Illinois, Urbana-Champaign
22
L6: Memory Hierarchy IV

## Shared memory bank conflicts

- Shared memory is as fast as registers if there are no bank conflicts

- The fast case:
  - If all threads of a half-warp access different banks, there is no bank conflict
  - If all threads of a half-warp access the identical address, there is no bank conflict (broadcast)
- The slow case:
  - Bank Conflict: multiple threads in the same half-warp access the same bank
  - Must serialize the accesses
  - Cost = max # of simultaneous accesses to a single bank

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
ECE 498AL, University of Illinois, Urbana-Champaign
23
L6: Memory Hierarchy IV

## Bank Addressing Examples



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
ECE 498AL, University of Illinois, Urbana-Champaign
24
L6: Memory Hierarchy IV

6

## Bank Addressing Examples

- 2-way Bank Conflicts
  - Linear addressing stride == 2

- 8-way Bank Conflicts
  - Linear addressing stride == 8

25
L6: Memory Hierarchy IV

---

## Putting It Together: Global Memory Coalescing and Bank Conflicts

- Let's look at matrix transpose
- Simple goal: Replace A[i][j] with A[j][i]
- Any reuse of data?
- Do you think shared memory might be useful?

26
L6: Memory Hierarchy IV

---

## Matrix Transpose (from SDK)

```
_global__ void transpose(float *odata, float *idata, int width, int height)
{
```

*odata and idata in global memory*

```
    // read the element
    unsigned int xIndex = blockIdx.x * BLOCK_DIM + threadIdx.x;
    unsigned int yIndex = blockIdx.y * BLOCK_DIM + threadIdx.y;
    unsigned int index_in = yIndex * width + xIndex;
    temp = idata[index_in];

    // write the transposed element to global memory
    xIndex = blockIdx.y * BLOCK_DIM + threadIdx.x;
    yIndex = blockIdx.x * BLOCK_DIM + threadIdx.y;
    unsigned int index_out = yIndex * height + xIndex;
    odata[index_out] = temp;

}
```

CS6963

27
L6: Memory Hierarchy IV

---

## Coalesced Matrix Transpose

```
_global__ void transpose(float *odata, float *idata, int width, int height)
{
    __shared__ float block[BLOCK_D
```

*odata and idata in global memory*

```
    // read the matrix tile into shared memory
    unsigned int xIndex = blockIdx.x * BLOCK_DIM + threadIdx.x;
    unsigned int yIndex = blockIdx.y * BLOCK_DIM + threadIdx.y;
    unsigned int index_in = yIndex * width + xIndex;
    block[threadIdx.y][threadIdx.x] = idata[index_in];

    __syncthreads();

    // write the transposed matrix tile to global memory
    xIndex = blockIdx.y * BLOCK_DIM + threadIdx.x;
    yIndex = blockIdx.x * BLOCK_DIM + threadIdx.y;
    unsigned int index_out = yIndex * height + xIndex;
    odata[index_out] = block[threadIdx.x][threadIdx.y];
}
```

*Rearrange in shared memory and write back efficiently to global memory*

CS6963

28
L6: Memory Hierarchy IV

7

## Optimized Matrix Transpose (from SDK)

```
_global__ void transpose(float *odata, float *idata, int width, int height)
{
    __shared__ float block[BLOCK_DIM][          odata and idata in
                                               global memory

    // read the matrix tile into shared memory
    unsigned int xIndex = blockIdx.x * BLOCK_DIM + threadIdx.x;
    unsigned int yIndex = blockIdx.y * BLOCK_DIM + threadIdx.y;
    unsigned int index_in = yIndex * width + xIndex;
    block[threadIdx.y][threadIdx.x] = idata[index_in];

    __syncthreads();

    // write the transposed matrix tile to global memory       Rearrange in
    xIndex = blockIdx.y * BLOCK_DIM + threadIdx.x;            shared memory
    yIndex = blockIdx.x * BLOCK_DIM + threadIdx.y;            and write back
    unsigned int index_out = yIndex * height + xIndex;         efficiently to
    odata[index_out] = block[threadIdx.x][threadIdx.y];       global memory
}
```

CS6963
29
L6: Memory Hierarchy IV

---

## Further Optimization: Partition Camping

- A further optimization improves bank conflicts in global memory
  - But has not proven that useful in codes with additional computation
- Map blocks to different parts of chips

  int bid = blockIdx.x + gridDim.x*blockIdx.y;
  by = bid%gridDim.y;
  bx = ((bid/gridDim.y)+by)%gridDim.x;

30
L6: Memory Hierarchy IV

---

## Performance Results for Matrix Transpose (GTX280)



SDK-prev: all optimizations other than partition camping
CHiLL: generated by our compiler
SDK-new: includes partition camping

31
L6: Memory Hierarchy IV

---

## Summary of Lecture

- Completion of bandwidth optimizations
  - Global memory coalescing
  - Alignment
  - Shared memory bank conflicts
  - "Partitioning camping"
- Matrix transpose example

CS6963
32
L6: Memory Hierarchy IV

8

# Next Time

- A look at correctness
- Synchronization mechanisms

CS6963

33
L6: Memory Hierarchy IV

THE
UNIVERSITY
OF UTAH