

## L5: Memory Hierarchy Optimization III, Data Placement, cont. and Memory Bandwidth Optimizations

CS6235

L5: Memory Hierarchy, III

1



## Administrative

- Next assignment available
  - Next three slides
  - Goals of assignment:
    - simple memory hierarchy management
    - block-thread decomposition tradeoff
  - Due Tuesday, Feb. 7, 5PM
  - Use handin program on CADE machines
    - "handin CS6235 lab2 <probfile>"

CS6235

L5: Memory Hierarchy, III

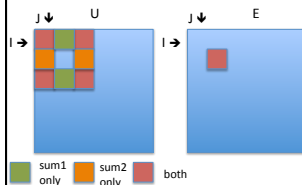
2



## Assignment 2: Memory Hierarchy Optimization Due Tuesday, February 7 at 5PM

Sobel edge detection:

Find the boundaries of the image where there is significant difference as compared to neighboring "pixels" and replace values to find edges



```
for (i = 1; i < ImageNRows - 1; i++)
  for (j = 1; j < ImageNCols - 1; j++)
  {
    sum1 = u[i-1][j+1] - u[i-1][j-1]
          + 2 * u[i][j+1] - 2 * u[i][j-1]
          + u[i+1][j+1] - u[i+1][j-1];
    sum2 = u[i-1][j-1] + 2 * u[i-1][j] + u[i-1][j+1]
          - u[i+1][j-1] - 2 * u[i+1][j] - u[i+1][j+1];

    magnitude = sum1*sum1 + sum2*sum2;
    if (magnitude > THRESHOLD)
      e[i][j] = 255;
    else
      e[i][j] = 0;
  }
```

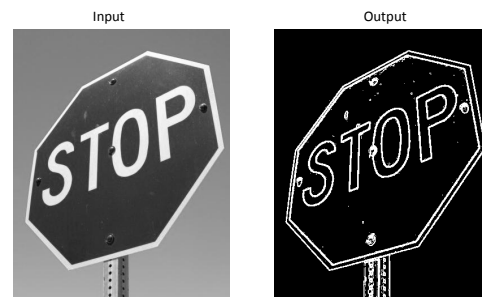
CS6963

3

L4: Memory Hierarchy I



## Example



## General Approach

### 0. Provided

- Input file
- Sample output file
- CPU implementation

### 1. Structure

- Compare CPU version and GPU version output [[compareInt](#)]
- Time performance of two GPU versions (see 2 & 3 below) [[EventRecord](#)]

### 2. GPU version 1 (partial credit if correct)

implementation using global memory

### 3. GPU version 2 (highest points to best performing versions)

use memory hierarchy optimizations from previous, current and Monday's lecture

### 4. Extra credit: Try two different block / thread decompositions. What happens if you use more threads versus more blocks? What if you do more work per thread? Explain your choices in a README file.

Handin using the following on CADE machines, where probfile includes all files

"handin cs6235 lab2 <probfile>"



## Overview of Lecture

- Review: Tiling for computation partitioning and fixed capacity storage
  - Now for constant memory and registers
- Quick look at texture memory
- First bandwidth optimization: global memory coalescing
- Reading:
  - Chapter 5, Kirk and Hwu book
  - Or, <http://courses.ece.illinois.edu/ece498/al/textbook/Chapter4-CudaMemoryModel.pdf>

CS6235

L5: Memory Hierarchy, III

6



## 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

CS6235

L5: Memory Hierarchy, III

7



## Constant Memory Example

- Signal recognition example:
  - Apply input signal (a vector) to a set of precomputed transform matrices
  - Compute  $M_1V, M_2V, \dots, M_nV$

```
__constant__ float d_signalVector[M];
__device__ float R[N][M];

__host__ void outerApplySignal () {
    float *h_inputSignal;
    dim3 dimGrid(N);
    dim3 dimBlock(M);
    cudaMemcpyToSymbol (d_signalVector,
        h_inputSignal, M*sizeof(float));
    // input matrix is in d_mat
    ApplySignal<<<dimGrid,dimBlock>>>
        (d_mat, M);
}

__global__ void ApplySignal (float * d_mat,
    int M) {
    float result = 0.0; /* register */
    for (j=0; j<M; j++)
        result += d_mat[blockIdx.x][threadIdx.x][j] *
            d_signalVector[j];
    R[blockIdx.x][threadIdx.x] = result;
}
```

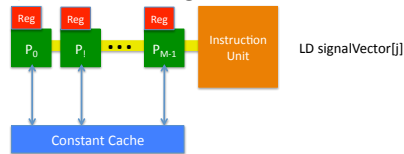
CS6235

L3: Memory Hierarchy, 1



### More on Constant Cache

- Example from previous slide
  - All threads in a block accessing same element of signal vector
  - Brought into cache for first access, then latency equivalent to a register access



CS6235

L3: Memory Hierarchy, 1



### Additional Detail

- Suppose each thread accesses different data from constant memory on same instruction
  - Reuse across threads?
    - Consider capacity of constant cache and locality
    - Code transformation needed? -- tile for constant memory, constant cache
    - Cache latency proportional to number of accesses in a warp
  - No reuse?
    - Should not be in constant memory.

CS6235

L3: Memory Hierarchy, 1



### "Tiling" for Registers

- A similar technique can be used to map data to registers
- Unroll-and-jam
  - Unroll outer loops in a nest and fuse together resulting inner loops
  - Equivalent to "strip-mine" followed by permutation and unrolling
- Fusion safe if dependences are not reversed
- Scalar replacement
  - May be followed by replacing array references with scalar variables to help compiler identify register opportunities

CS6235

L5: Memory Hierarchy, III



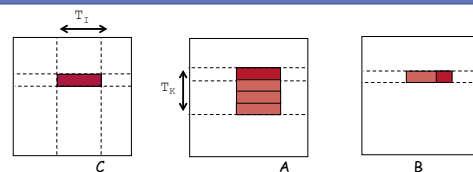
### Unroll-and-jam for matrix multiply

#### Tiling inner loops I and K (+permutation)

```

for (K = 0; K < N; K += TK)
  for (I = 0; I < N; I += TI)
    for (J = 0; J < N; J++)
      for (KK = K; KK < min(K + TK, N); KK++)
        for (II = I; II < min(I + TI, N); II++)
          P[J][II] = P[J][II] + M[KK][II] * N[J][KK];

```



L5: Memory Hierarchy, III

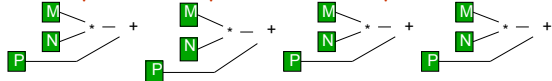
12



### First, Apply Unroll-and-Jam Unroll II loop, $T_I = 4$ (equiv. to unroll&jam)

```
for (K = 0; K < N; K += T_K)
  for (I = 0; I < N; I += 4)
    for (J = 0; J < N; J++)
      for (KK = K; KK < min(K + T_K, N); KK++)
        P[J][I] = P[J][I] + M[KK][I] * N[J][KK];
        P[J][I+1] = P[J][I+1] + M[KK][I+1] * N[J][KK];
        P[J][I+2] = P[J][I+2] + M[KK][I+2] * N[J][KK];
        P[J][I+3] = P[J][I+3] + M[KK][I+3] * N[J][KK];
```

Now parallel computations are exposed



LS: Memory Hierarchy, III

13



### Now can expose registers using scalar replacement (or simply unroll kk loop)

Scalar Replacement: Replace accesses to P with scalars

```
for (K = 0; K < N; K += T_K)
  for (I = 0; I < N; I += 4)
    for (J = 0; J < N; J++) {
      P0 = P[J][I]; P1 = P[J][I+1]; P2 = P[J][I+2]; P3 = P[J][I+3];
      for (KK = K; KK < min(K + T_K, N); KK++) {
        P0 = P0 + M[KK][I] * N[J][KK];
        P1 = P1 + M[KK][I+1] * N[J][KK];
        P2 = P2 + M[KK][I+2] * N[J][KK];
        P3 = P3 + M[KK][I+3] * N[J][KK];
      }
      P[J][I] = P0; P[J][I+1] = P1; P[J][I+2] = P2; P[J][I+3] = P3;
    }
```

Now P accesses can be mapped to "named registers"

LS: Memory Hierarchy, III

14



### Overview of Texture Memory

- Recall, texture cache of read-only data
- Special protocol for allocating and copying to GPU
  - texture<Type, Dim, ReadMode> texRef;
    - Dim: 1, 2 or 3D objects
- Special protocol for accesses (macros)
  - tex2D(<name>, dim1, dim2);
- In full glory can also apply functions to textures
- Writing possible, but unsafe if followed by read in same kernel

CS6235

15  
LS: Memory Hierarchy III



### Using Texture Memory (simpleTexture project from SDK)

```
cudaMalloc((void**) &d_data, size);
cudaChannelFormatDesc channelDesc = cudaCreateChannelDesc(32, 0, 0, 0,
  cudaChannelFormatKindFloat);
cudaArray* cu_array;
cudaMallocArray(&cu_array, &channelDesc, width, height);
cudaMemcpyToArray(cu_array, 0, 0, h_data, size, cudaMemcpyHostToDevice);
// set texture parameters
tex.addressMode[0] = tex.addressMode[1] = cudaAddressModeWrap;
tex.filterMode = cudaFilterModeLinear; tex.normalized = true;
cudaBindTextureToArray(tex, cu_array, channelDesc);
// execute the kernel
transformKernel<<< dimGrid, dimBlock, 0 >>>(d_data, width, height, angle);
```

Kernel function:  
// declare texture reference for 2D float texture  
texture<float, 2, cudaReadModeElementType> tex;

... = tex2D(tex, i, j);

CS6235

16  
LS: Memory Hierarchy III



## When to use Texture (and Surface) Memory

(From 5.3 of CUDA manual) Reading device memory through texture or surface fetching present some benefits that can make it an advantageous alternative to reading device memory from global or constant memory:

- If memory reads to global or constant memory will not be coalesced, higher bandwidth can be achieved providing that there is locality in the texture fetches or surface reads (this is less likely for devices of compute capability 2.x given that global memory reads are cached on these devices);
- Addressing calculations are performed outside the kernel by dedicated units;
- Packed data may be broadcast to separate variables in a single operation;
- 8-bit and 16-bit integer input data may be optionally converted to 32-bit floating-point values in the range [0.0, 1.0] or [-1.0, 1.0] (see Section 3.2.4.1.1).

LS: Memory Hierarchy, III

17



## Memory Bandwidth Optimization

- Goal is to maximize utility of data for each data transfer from global memory
- Memory system will "coalesce" accesses for a collection of consecutive threads if they are within an aligned 128 byte portion of memory (from half-warp or warp)
- Implications for programming:
  - Desirable to have consecutive threads in tx dimension accessing consecutive data in memory
  - Significant performance impact, but Fermi data cache makes it slightly less important

LS: Memory Hierarchy, III

18



## Introduction to Global Memory Bandwidth: Understanding Global Memory Accesses

### Memory protocol for compute capability 1.2\* (CUDA Manual 5.1.2.1)

- 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

CS6235

LS: Memory Hierarchy, III  
LS: Memory Hierarchy II

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

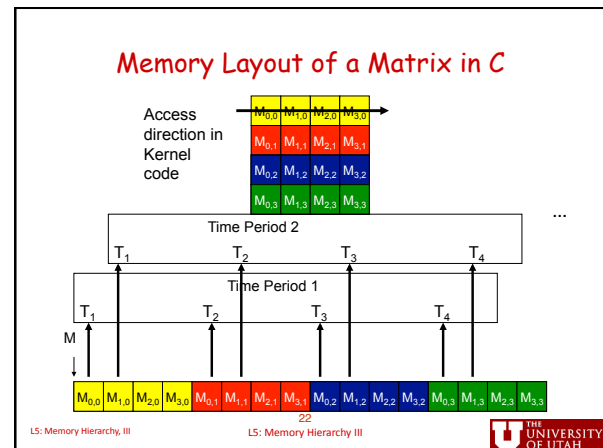
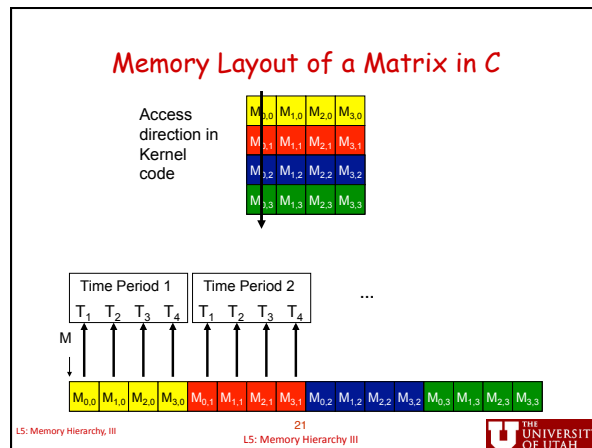
- For compute capability 1.0 and 1.1
  - Threads must access the words in a segment in sequence
  - The kth thread must access the kth word

CS6235

LS: Memory Hierarchy, III

20





### Summary of Lecture

- Tiling transformation
  - For computation partitioning
  - For limited capacity in shared memory
  - For registers
- Matrix multiply example
- Unroll-and-jam for registers
- Bandwidth optimization
  - Global memory coalescing

CS6235

LS: Memory Hierarchy, III

23

THE UNIVERSITY OF UTAH

### Next Time

- Complete bandwidth optimizations
  - Shared memory bank conflicts
  - Bank conflicts in global memory (briefly)

CS6235

LS: Memory Hierarchy, III

24

THE UNIVERSITY OF UTAH