

## L6: Memory Hierarchy Optimization I, Data Placement

CS6963



## Administrative

- Projects mostly graded
  - Grades by Wednesday, possibly sooner
- Homework coming out this afternoon
  - Due Wednesday, Feb. 18

CS6963

2  
L6: Memory Hierarchy I

## Overview

- Where data can be stored
  - And how to get it there
- Some guidelines for where to store data
- High level description of how to write code to optimize for memory hierarchy
  - Fill in details Wednesday

CS6963

3  
L6: Memory Hierarchy I

## 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

4  
L6: Memory Hierarchy I

### Optimizing the Memory Hierarchy on GPUs

- 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.
  - **Align** data structures to address boundaries
  - **Coalesce** global memory accesses
  - **Avoid memory bank conflicts** to increase memory access parallelism

CS6963

5  
L6: Memory Hierarchy I



### Reuse and Locality

- Consider how data is accessed
  - **Data reuse:**
    - Same data used multiple times
    - Intrinsic in computation
  - **Data locality:**
    - Data is reused and is present in "fast memory"
    - Same data or same data transfer
- If a computation has reuse, what can we do to get locality?
  - Appropriate data placement and layout
  - Code reordering transformations

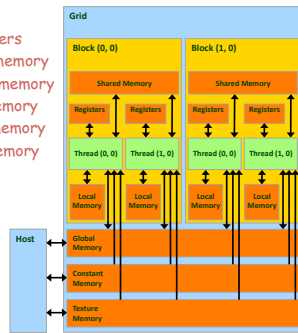
CS6963

6  
L6: Memory Hierarchy I



### Programmer's View: Memory Spaces

- 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**
  - Read only per-grid **texture memory**
- The host can read/write **global, constant, and texture memory**



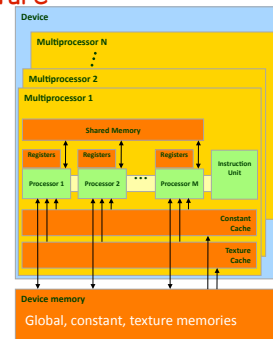
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

7  
L6: Memory Hierarchy



### Hardware Implementation: Memory Architecture

- The local, global, constant, and texture spaces are regions of device memory
- Each multiprocessor has:
  - A set of 32-bit **registers** per processor
  - **On-chip shared memory**
    - Where the shared memory space resides
  - A **read-only constant cache**
    - To speed up access to the constant memory space
  - A **read-only texture cache**
    - To speed up access to the texture memory space



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

8  
L6: Memory Hierarchy I



### Terminology Review

- device = GPU = set of multiprocessors
- Multiprocessor = set of processors & shared memory
- Kernel = GPU program
- Grid = array of thread blocks that execute a kernel
- Thread block = group of SIMD threads that execute a kernel and can communicate via shared memory

Memory	Location	Cached	Access	Who
Local	Off-chip	No	Read/write	One thread
Shared	On-chip	N/A - resident	Read/write	All threads in a block
Global	Off-chip	No	Read/write	All threads + host
Constant	Off-chip	Yes	Read	All threads + host
Texture	Off-chip	Yes	Read	All threads + host

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

9

L6: Memory Hierarchy I



### Access Times (REWRITE?)

- Register - dedicated HW - single cycle
- Constant and Texture caches - possibly single cycle, proportional to addresses accessed by warp
- Shared Memory - dedicated HW - single cycle
- Local Memory - DRAM, no cache - \*slow\*
- Global Memory - DRAM, no cache - \*slow\*
- Constant Memory - DRAM, cached, 1...10s...100s of cycles, depending on cache locality
- Texture Memory - DRAM, cached, 1...10s...100s of cycles, depending on cache locality
- Instruction Memory (invisible) - DRAM, cached

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

10

L5: Memory Hierarchy



### Data Placement: Conceptual

- Copies from host to device go to some part of global memory (possibly, constant or texture memory)
- How to use SP shared memory
  - Must construct or be copied from global memory by kernel program
- How to use constant or texture cache
  - Read-only "reused" data can be placed in constant & texture memory by host
- Local memory
  - Deals with capacity limitations of registers and shared memory
  - Eliminates worries about race conditions
  - ... but SLOW

CS6963

11

L6: Memory Hierarchy I



### Data Placement: Syntax

- Through type qualifiers
  - `__constant__`, `__shared__`, `__local__`, `__device__`
- Through `cudaMemcpy` calls
  - Flavor of call and symbolic constant designate where to copy
- Implicit default behavior
  - Device memory without qualifier is global memory
  - Host by default copies to global memory
  - Thread variables go into registers unless capacity exceeded, then local memory

CS6963

12

L6: Memory Hierarchy I



### Language Extensions: Variable Type Qualifiers

	Memory	Scope	Lifetime
<code>__device__ __local__ int LocalVar;</code>	local	thread	thread
<code>__device__ __shared__ int SharedVar;</code>	shared	block	block
<code>__device__ int GlobalVar;</code>	global	grid	application
<code>__device__ __constant__ int ConstantVar;</code>	constant	grid	application

- `__device__` is optional when used with `__local__`, `__shared__`, or `__constant__`
- Automatic variables without any qualifier reside in a register
  - Except arrays that reside in local memory

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

13  
L6: Memory Hierarchy



### 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;`

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

14  
L6: Memory Hierarchy I



### Rest of Today's Lecture

- Mechanics of how to place data in shared memory and constant memory
- Tiling transformation to reuse data within
  - Shared memory
  - Constant cache

15  
L6: Memory Hierarchy I



### Constant Memory Example

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

```
__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));
    ApplySignal<<<dimGrid,dimBlock>>(M);
}

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

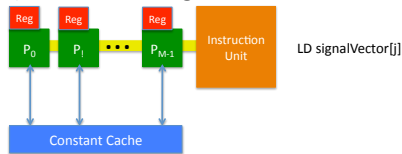
CS6963

16  
L6: Memory Hierarchy I



### 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



CS6963

17  
L6: Memory Hierarchy I

### Additional Detail

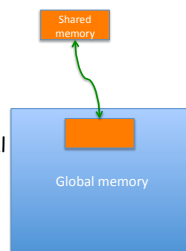
- Suppose each thread accesses different data from constant memory on same instruction
  - Reuse across threads?
    - Consider capacity of constant cache and locality
    - Code transformation needed? (later in lecture)
    - Cache latency proportional to number of accesses in a warp
  - No reuse?
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

CS6963

18  
L6: Memory Hierarchy I

### Now Let's Look at Shared Memory

- Common Programming Pattern (5.1.2 of CUDA manual)
  - Load data into shared memory
  - Synchronize (if necessary)
  - Operate on data in shared memory
  - Synchronize (if necessary)
  - Write intermediate results to global memory
  - Repeat until done



CS6963

19  
L6: Memory Hierarchy I

### Mechanics of Using Shared Memory

- `__shared__` type qualifier required
- Must be allocated from global/device function, or as "extern"
- Examples:

```
extern __shared__ float d_s_array[];    __global__ void compute2() {
/* a form of dynamic allocation */      __shared__ float d_s_array[M];
/* MEMSIZE is size of per-block */    /* create or copy from global memory */
/* shared memory*/                    d_s_array[j] = ...;
__host__ void outerCompute() {
  compute<<gs,bs,MEMSIZE>>();          /* write result back to global memory */
}                                       d_g_array[j] = d_s_array[j];
__global__ void compute() {
  d_s_array[i] = ...;                  }
}
```

CS6963

20  
L6: Memory Hierarchy I

## Matrix Transpose (from SDK)

```

_global__ void transpose(float *odata, float *idata, int width, int height)
{
    __shared__ float block[BLOCK_DIM][BLOCK_DIM];

    // 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];
}

```

odata and idata in global memory

Rearrange in shared memory and write back efficiently to global memory

CS6963

21  
L6: Memory Hierarchy I

## Recall Reuse and Locality

- Consider how data is accessed
  - **Data reuse:**
    - Same data used multiple times
    - Intrinsic in computation
  - **Data locality:**
    - Data is reused and is present in “fast memory”
    - Same data or same data transfer
- If a computation has reuse, what can we do to get locality?
  - Appropriate data placement and layout
  - Code reordering transformations

CS6963

22  
L6: Memory Hierarchy I

## Temporal Reuse in Sequential Code

- Same data used in distinct iterations  $I$  and  $I'$

```

for (i=1; i<N; i++)
    for (j=1; j<N; j++)
        A[j] = A[j] + A[j+1] + A[j-1];

```

- $A[j]$  has self-temporal reuse in loop  $i$

CS6963

23  
L6: Memory Hierarchy I

## Spatial Reuse (Ignore for now)

- Same data transfer (usually cache line) used in distinct iterations  $I$  and  $I'$

```

for (i=1; i<N; i++)
    for (j=1; j<N; j++)
        A[j] = A[j] + A[j+1] + A[j-1];

```

- $A[j]$  has self-spatial reuse in loop  $j$
- **Multi-dimensional array note:** C arrays are stored in row-major order

CS6963

24  
L6: Memory Hierarchy I

## Group Reuse

- Same data used by distinct references

```
for (i=1; i<N; i++)
  for (j=1; j<N; j++)
    A[j]= A[j]+A[j+1]+A[j-1];
```

- $A[j]$ ,  $A[j+1]$  and  $A[j-1]$  have group reuse (spatial and temporal) in loop  $j$

CS6963

25  
L6: Memory Hierarchy I

## Can Use Reordering Transformations!

- Analyze reuse in computation
- Apply loop reordering transformations to improve locality based on reuse
- With any loop reordering transformation, always ask
  - **Safety?** (doesn't reverse dependences)
  - **Profitability?** (improves locality)

CS6963

26  
L6: Memory Hierarchy I

26

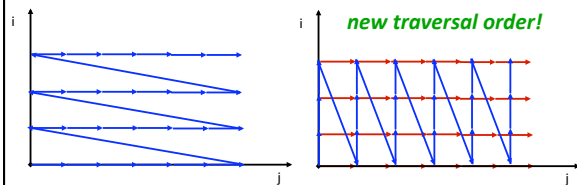


## Loop Permutation: A Reordering Transformation

Permute the order of the loops to modify the traversal order

```
for (i= 0; i<3; i++)
  for (j=0; j<6; j++)
    A[i][j+1]=A[i][j]+B[j];
```

```
for (j=0; j<6; j++)
  for (i= 0; i<3; i++)
    A[i][j+1]=A[i][j]+B[j];
```



Which one is better for row-major storage?

CS6963

27  
L6: Memory Hierarchy I

## Safety of Permutation

- **Intuition:** Cannot permute two loops  $i$  and  $j$  in a loop nest if doing so reverses the direction of any dependence.
- Loops  $i$  through  $j$  of an  $n$ -deep loop nest are *fully permutable* if for all dependences  $D$ , either
  - $(d_i, \dots, d_{j-1}) > 0$
  - or
  - for all  $k$ ,  $i \leq k \leq j$ ,  $d_k \geq 0$
- **Stated without proof:** Within the affine domain,  $n-1$  inner loops of  $n$ -deep loop nest can be transformed to be fully permutable.

CS6963

28  
L6: Memory Hierarchy I

### Simple Examples: 2-d Loop Nests

```
for (i= 0; i<3; i++)
for (j=0; j<6; j++)
  A[i][j+1]=A[i][j]+B[j];
```

```
for (i= 0; i<3; i++)
  for (j=0; j<6; j++)
    A[i+1][j-1]=A[i][j]
      +B[j];
```

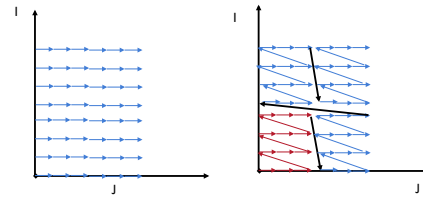
- Distance vectors
- Ok to permute?

CS6963

29  
L6: Memory Hierarchy I

### Tiling (Blocking): Another Loop Reordering Transformation

- Blocking reorders loop iterations to bring iterations that reuse data closer in time



CS6963

30  
L6: Memory Hierarchy I

### Tiling Example

```
for (j=1; j<M; j++)
  for (i=1; i<N; i++)
    D[i] = D[i] + B[j][i];
```

Strip  
mine

```
for (j=1; j<M; j++)
  for (ii=1; ii<N; ii+=s)
    for (i=ii, i<min(ii+s-1,N), i++)
      D[i] = D[i] +B[j][i];
```

Permute

```
for (ii=1; ii<N; ii+=s)
  for (j=1; j<M; j++)
    for (i=ii, i<min(ii+s-1,N), i++)
      D[i] = D[i] +B[j][i];
```

CS6963

31  
L6: Memory Hierarchy I

### Legality of Tiling

- Tiling = strip-mine and permutation
  - Strip-mine does not reorder iterations
  - Permutation must be legal
- OR
- strip size less than dependence distance

CS6963

32  
L6: Memory Hierarchy I



### A Few Words On Tiling

- Tiling can be used hierarchically to compute partial results on a block of data wherever there are capacity limitations
  - Between grids if data exceeds global memory capacity
  - Across thread blocks if shared data exceeds shared memory capacity
  - Within threads if data in constant cache exceeds cache capacity

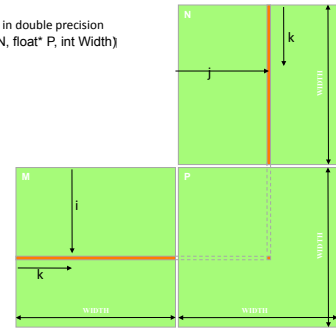
CS6963

33  
L6: Memory Hierarchy I



### Matrix Multiplication A Simple Host Version in C

```
// Matrix multiplication on the (CPU) host in double precision
void MatrixMulOnHost(float* M, float* N, float* P, int Width)
{
    for (int i = 0; i < Width; ++i)
        for (int j = 0; j < Width; ++j) {
            double sum = 0;
            for (int k = 0; k < Width; ++k) {
                double a = M[i * width + k];
                double b = N[k * width + j];
                sum += a * b;
            }
            P[j * Width + i] = sum;
        }
}
```



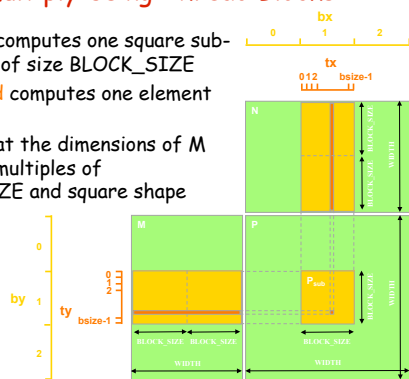
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009  
ECE498AL, University of Illinois, Urbana-Champaign

34  
L6: Memory Hierarchy I



### Tiled Multiply Using Thread Blocks

- One **block** computes one square sub-matrix  $P_{sub}$  of size `BLOCK_SIZE`
- One **thread** computes one element of  $P_{sub}$
- Assume that the dimensions of  $M$  and  $N$  are multiples of `BLOCK_SIZE` and square shape



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

35  
L6: Memory Hierarchy I



### Shared Memory Usage

- Assume each SMP has 16KB shared memory
  - Each Thread Block uses  $2 * 256 * 4B = 2KB$  of shared memory.
  - Can potentially have up to 8 Thread Blocks actively executing
  - For `BLOCK_SIZE = 16`, this allows up to  $8 * 512 = 4,096$  pending loads
    - In practice, there will probably be up to half of this due to scheduling to make use of SPs.
  - The next `BLOCK_SIZE = 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

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

36  
L6: Memory Hierarchy I



### First-order Size Considerations

- Each Thread Block should have a minimal of 192 threads
  - BLOCK\_SIZE of 16 gives  $16*16 = 256$  threads
- A minimal of 32 Thread Blocks
  - A  $1024*1024$  P Matrix gives  $64*64 = 4096$  Thread Blocks
- Each thread block performs  $2*256 = 512$  float loads from global memory for  $256 * (2*16) = 8,192$  mul/add operations.
  - Memory bandwidth no longer a limiting factor

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

37  
L6: Memory Hierarchy I



### CUDA Code - Kernel Execution Configuration

```
// Setup the execution configuration
dim3 dimBlock(BLOCK_SIZE, BLOCK_SIZE);
dim3 dimGrid(N.width / dimBlock.x,
             M.height / dimBlock.y);
```

For very large N and M dimensions, one will need to add another level of blocking and execute the second-level blocks sequentially.

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

38  
L6: Memory Hierarchy I



### CUDA Code - Kernel Overview

```
// Block index
int bx = blockIdx.x;
int by = blockIdx.y;
// Thread index
int tx = threadIdx.x;
int ty = threadIdx.y;

// Pvalue stores the element of the block sub-matrix
// that is computed by the thread
float Pvalue = 0;

// Loop over all the sub-matrices of M and N
// required to compute the block sub-matrix
for (int m = 0; m < M.width/BLOCK_SIZE; ++m) {
    code from the next few slides };
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

39  
L6: Memory Hierarchy I



### CUDA Code - Load Data to Shared Memory

```
// Get a pointer to the current sub-matrix Msub of M
Matrix Msub = GetSubMatrix(M, m, by);

// Get a pointer to the current sub-matrix Nsub of N
Matrix Nsub = GetSubMatrix(N, bx, m);

__shared__ float Ms[BLOCK_SIZE][BLOCK_SIZE];
__shared__ float Ns[BLOCK_SIZE][BLOCK_SIZE];

// each thread loads one element of the sub-matrix
Ms[ty][tx] = GetMatrixElement(Msub, tx, ty);

// each thread loads one element of the sub-matrix
Ns[ty][tx] = GetMatrixElement(Nsub, tx, ty);
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

40  
L6: Memory Hierarchy I



### CUDA Code - Compute Result

```
// Synchronize to make sure the sub-matrices are loaded
// before starting the computation
__syncthreads();

// each thread computes one element of the block sub-matrix
for (int k = 0; k < BLOCK_SIZE; ++k)
    Pvalue += Ms[ty][k] * Ns[k][tx];

// Synchronize to make sure that the preceding
// computation is done before loading two new
// sub-matrices of M and N in the next iteration
__syncthreads();
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

41  
L6: Memory Hierarchy I



### CUDA Code - Save Result

```
// Get a pointer to the block sub-matrix of P
Matrix Psub = GetSubMatrix(P, bx, by);

// Write the block sub-matrix to device memory;
// each thread writes one element
SetMatrixElement(Psub, tx, ty, Pvalue);
```

This code should run at about 45 GFLOPS

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
ECE 498AL, University of Illinois, Urbana-Champaign

42  
L6: Memory Hierarchy I



### Matrix Multiply in CUDA

- Imagine you want to compute extremely large matrices.
  - That don't fit in global memory
- This is where an additional level of tiling could be used, between grids

CS6963

43  
L6: Memory Hierarchy I



### Summary of Lecture

- How to place data in constant memory and shared memory
- Reordering transformations to improve locality
- Tiling transformation
- Matrix multiply example

CS6963

44  
L6: Memory Hierarchy I



### Next Time

- Complete this example
- Reasoning about reuse and locality
- Talk about projects and assign proposal