

## L4: Memory Hierarchy Optimization II, Locality and Data Placement

CS6235

L3: Memory Hierarchy, 1



## Administrative

- Assignment due Friday, Jan. 18, 5 PM
  - Use handin program on CADE machines
    - "handin CS6235 lab1 <probfile>"
- Mailing list
  - [CS6235@list.eng.utah.edu](mailto:CS6235@list.eng.utah.edu)
    - Please use for all questions suitable for the whole class
    - Feel free to answer your classmates questions!

CS6235

L3: Memory Hierarchy, 1



## Overview of Lecture

- More on tiling for shared memory and constant memory
- Reading:
  - Chapter 5, Kirk and Hwu book
  - Or, <http://courses.ece.illinois.edu/ece498/al/textbook/Chapter4-CudaMemoryModel.pdf>

CS6235

L3: Memory Hierarchy, 1



## Review: 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

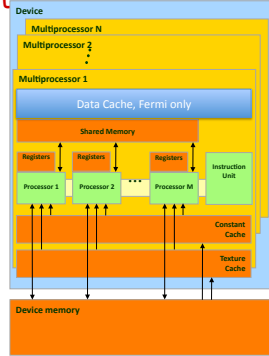
CS6235

L3: Memory Hierarchy, 1



## Hardware Implementation: Memory Architecture

- The local, global, constant, and texture spaces are regions of device memory (DRAM)
- 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
  - Data cache (Fermi only)**



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
EEL 458AL, University of Illinois, Urbana-Champaign

L3: Memory Hierarchy, 1



## Review: 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

CS6235

L3: Memory Hierarchy, 1



## Memory Hierarchy Example: Matrix vector multiply

```
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
    a[i] += c[j][i] * b[j];
  }
}
```

Remember to:

- Consider correctness of parallelization strategy (next week)
- Exploit locality in shared memory and registers

L3: Memory Hierarchy, 1



## Resulting CUDA code (Automatically Generated by our Research Compiler)

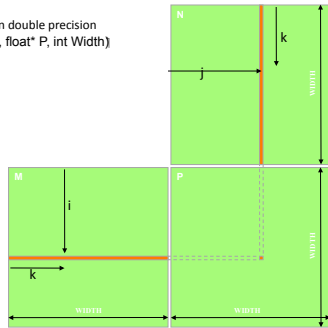
```
__global__ mv_GPU(float* a, float* b, float** c) {
  int bx = blockIdx.x; int tx = threadIdx.x;
  __shared__ float bcpy[32];
  double acpy = a[tx + 32 * bx];
  for (k = 0; k < 32; k++) {
    bcpy[tx] = b[32 * k + tx];
    __syncthreads();
    //this loop is actually fully unrolled
    for (j = 32 * k; j <= 32 * k + 32; j++) {
      acpy = acpy + c[j][32 * bx + tx] * bcpy[j];
    }
    __syncthreads();
  }
  a[tx + 32 * bx] = acpy;
}
```

L3: Memory Hierarchy, 1



### Ch. 4: 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[i * Width + j] = sum;
        }
}
```



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007  
EET-458AM, University of Illinois, Urbana-Champaign L3: Memory Hierarchy, 1



### Discussion (Simplified Code)

<pre>for (int i = 0; i &lt; Width; ++i)     for (int j = 0; j &lt; Width; ++j) {         double sum = 0;         for (int k = 0; k &lt; Width; ++k) {             double a = M[i * width + k];             double b = N[k * width + j];             sum += a * b;         }         P[i * Width + j] = sum;     } }</pre>	<pre>for (int i = 0; i &lt; Width; ++i)     for (int j = 0; j &lt; Width; ++j) {         double sum = 0;         for (int k = 0; k &lt; Width; ++k) {             sum += M[i][k] * N[k][j];         }         P[i][j] = sum;     } }</pre>
---	--

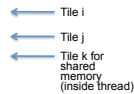
L4: Memory Hierarchy, II

10



### Let's Look at This Code

```
for (int i = 0; i < Width; ++i)
    for (int j = 0; j < Width; ++j) {
        double sum = 0;
        for (int k = 0; k < Width; ++k) {
            sum += M[i][k] * N[k][j];
        }
        P[i][j] = sum;
    }
}
```



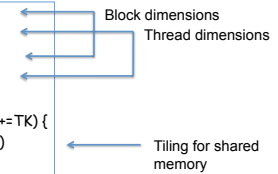
L4: Memory Hierarchy, II

11



### Strip-Mined Code

```
for (int ii = 0; ii < Width; ii+=TI)
    for (int i=ii; i<ii+TI; i++)
        for (int jj=0; jj<Width; jj+=TJ)
            for (int j = jj; j < jj+TJ; j++) {
                double sum = 0;
                for (int kk = 0; kk < Width; kk+=TK) {
                    for (int k = kk; k < kk+TK; k++)
                        sum += M[i][k] * N[k][j];
                }
                P[i][j] = sum;
            }
}
```



L4: Memory Hierarchy, II

12



### But this code doesn't match CUDA Constraints

```

for (int ii = 0; ii < Width; ii+=TI)
  for (int i=ii; i<ii+TI; i++)
    for (int jj=0; jj<Width; jj+=TJ)
      for (int j = jj; j < jj+TJ; j++){
        double sum = 0;
        for (int kk = 0; kk < Width; kk+=TK) {
          for (int k = kk; k < kk+TK; k++)
            sum += M[i][k] * N[k][j];
        }
        P[i][j] = sum;
      }
  
```

Block dimensions – must be unit stride

Thread dimensions – must be unit stride

Tiling for shared memory

Can we fix this?

L4: Memory Hierarchy, II 13 THE UNIVERSITY OF UTAH

### Unit Stride Tiling - Reflect Stride in Subscript Expressions

```

for (int ii = 0; ii < Width/TI; ii++)
  for (int i=0; i<TI; i++)
    for (int jj=0; jj<Width/TJ; jj++)
      for (int j = 0; j < TJ; j++){
        double sum = 0;
        for (int kk = 0; kk < Width; kk+=TK) {
          for (int k = kk; k < kk+TK; k++)
            sum += M[ii*TI+i][k] * N[k][jj*TJ+j];
        }
        P[ii*TI+i][jj*TJ+j] = sum;
      }
  
```

Block dimensions – must be unit stride

Thread dimensions – must be unit stride

Tiling for shared memory, no need to change

L4: Memory Hierarchy, II 14 THE UNIVERSITY OF UTAH

### What Does this Look Like in CUDA

```

#define TI 32
#define TJ 32
dim3 dimGrid(Width/TI, Width/TJ);
dim3 dimBlock(TI, TJ);
matMult<<<dimGrid, dimBlock>>>(M, N, P);

__global__ matMult(float *M, float *N, float *P) {
  ii = blockIdx.y; jj = blockIdx.x;
  i = threadIdx.y; j = threadIdx.x;
  double sum = 0;
  for (int kk = 0; kk < Width; kk+=TK) {
    for (int k = kk; k < kk+TK; k++) {
      sum += M[(ii*TI+i)*Width+k] *
        N[k*Width+jj*TJ+j];
    }
  }
  P[(ii*TI+i)*Width+jj*TJ+j] = sum;
}
  
```

Block and thread loops disappear

Tiling for shared memory, next slides

Array accesses to global memory are "linearized"

L4: Memory Hierarchy, II 15 THE UNIVERSITY OF UTAH

### What Does this Look Like in CUDA

```

#define TI 32
#define TJ 32
dim3 dimGrid(Width/TI, Width/TJ);
dim3 dimBlock(TI, TJ);
matMult<<<dimGrid, dimBlock>>>(M, N, P);
__global__ matMult(float *M, float *N, float *P) {
  ii = blockIdx.y; jj = blockIdx.x;
  i = threadIdx.y; j = threadIdx.x;
  double sum = 0;
  for (int kk = 0; kk < Width; kk+=TK) {
    Mds[j][i] = M[(ii*TI+i)*Width+kk*TK+j];
    Nds[j][i] = N[(kk*TK+i)*Width+jj*TJ+j];
    __syncthreads();
    for (int k = 0; k < TK; k++) {
      sum += Mds[j][k] * Nds[k][i];
    }
    __syncthreads();
  }
  P[(ii*TI+i)*Width+jj*TJ+j] = sum;
}
  
```

Block and thread loops disappear

Tiling for shared memory

Array accesses to global memory are "linearized"

L4: Memory Hierarchy, II 16 THE UNIVERSITY OF UTAH

## Derivation of code in text

- $TI = TJ = TK = \text{"TILE\_WIDTH"}$
- All matrices square, Width  $\times$  Width
- Copies of M and N in shared memory
  - TILE\_WIDTH  $\times$  TILE\_WIDTH
- "Linearized" 2-d array accesses:
  - $a[i][j]$  is equivalent to  $a[i*\text{Row} + j]$
- Each SM computes a "tile" of output matrix P from a block of consecutive rows of M and a block of consecutive columns of N
  - dim3 Grid (Width/TILE\_WIDTH, Width/TILE\_WIDTH);
  - dim3 Block (TILE\_WIDTH, TILE\_WIDTH)
- Then, location  $P[i][j]$  corresponds to
  - $P[\text{by}*\text{TILE\_WIDTH}+\text{ty}][\text{bx}*\text{TILE\_WIDTH}+\text{tx}]$  or
  - $P[\text{Row}][\text{Col}]$

L5: Memory Hierarchy, III

17



## Final Code (from text, p. 87)

```

__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];
3 & 4.  int bx = blockIdx.x; int by = blockIdx.y; int tx = threadIdx.x; int ty = threadIdx.y;
//Identify the row and column of the Pd element to work on
5 & 6.  int Row = by * TILE_WIDTH + ty;  int Col = bx * TILE_WIDTH + tx;
7.  float Pvalue = 0;
// Loop over the Md and Nd tiles required to compute the Pd element
8.  for (int m=0; m < Width / TILE_WIDTH; ++m) {
// Collaborative (parallel) loading of Md and Nd tiles into shared memory
9.  Mds [ty] [bx] = Md [Row*Width + (m*TILE_WIDTH + tx)];
10. Nds [ty] [bx] = Nd [(m*TILE_WIDTH + ty)*Width + Col];
11.  __syncthreads(); // make sure all threads have completed copy before calculation
12.  for (int k = 0; k < TILE_WIDTH; ++k) // Update Pvalue for TKxTK tiles in Mds and Nds
13.  Pvalue += Mds [ty] [k] * Nds [k] [bx];
14.  __syncthreads(); // make sure calculation complete before copying next tile
} // m loop
15.  Pd [Row*Width + Col] = Pvalue;
}

```

L5: Memory Hierarchy, III

18



## Summary of Lecture

- Matrix-matrix multiply example

CS6235

L3: Memory Hierarchy, 1

