# CS 380 - GPU and GPGPU Programming
# Lecture 13: GPU Compute APIs, Pt. 3

Markus Hadwiger, KAUST

# Reading Assignment #7 (until Oct 20)

Read (required):

- Programming Massively Parallel Processors book, 3rd edition,
  Chapter 7 (*Parallel Patterns: Convolution*)

- PTX Instruction Set Architecture 7.4 (`https://docs.nvidia.com/cuda/pdf/ptx_isa_7.4.pdf`)
  Read Chapters 1 – 3; get an overview of Chapter 12;
  browse through the other chapters to get a feeling for what PTX looks like

- CUDA SASS, Chapter 4: `https://docs.nvidia.com/cuda/pdf/CUDA_Binary_Utilities.pdf`

Read (optional):

- Inline PTX Assembly in CUDA (CUDA SDK: `Inline_PTX_Assembly.pdf`)

- Dissecting GPU Architecture through Microbenchmarking:

    Volta: `https://arxiv.org/abs/1804.06826`
    Turing: `https://arxiv.org/abs/1903.07486`

`https://developer.download.nvidia.com/video/gputechconf/gtc/2019/presentation/`
   `s9839-discovering-the-turing-t4-gpu-architecture-with-microbenchmarks.pdf`

    Ampere: `https://www.nvidia.com/en-us/on-demand/session/gtcspring21-s33322/`

# Semester Project (proposal until Oct 25!)

- Choosing your own topic encouraged!
  (we will also suggest some topics)

  - Pick something that you think is really cool!

  - Can be completely graphics or completely computation, or both combined

  - Can be built on CS 380 frameworks, NVIDIA OpenGL SDK, CUDA SDK, ...

- Write short (1-2 pages) project proposal by end of Sep (announced later)

  - Talk to us before you start writing!
    (content and complexity should fit the lecture)

- Submit semester project with report (deadline: Dec 9)

- Present semester project (event in final exams week: Dec 13 (tentative))

# Semester Project Ideas (1)

Some ideas for topics

• Procedural shading with noise + marble etc. (GPU Gems 2, chapter 26)

• Procedural shading with noise + bump mapping (GPU Gems 2, chapter 26)

• Subdivision surfaces (GPU Gems 2, chapter 7)

• Ambient occlusion, screen space ambient occlusion

• Shadow mapping, hard shadows, soft shadows

• Deferred shading

• Particle system rendering + CUDA particle sort

• Advanced image filters: fast bilateral filtering, Gaussian kD trees

• Advanced image de-convolution (e.g., convex L1 optimization)

• PDE solvers (e.g., anisotropic diffusion filtering, 2D level set segmentation, 2D fluid flow)

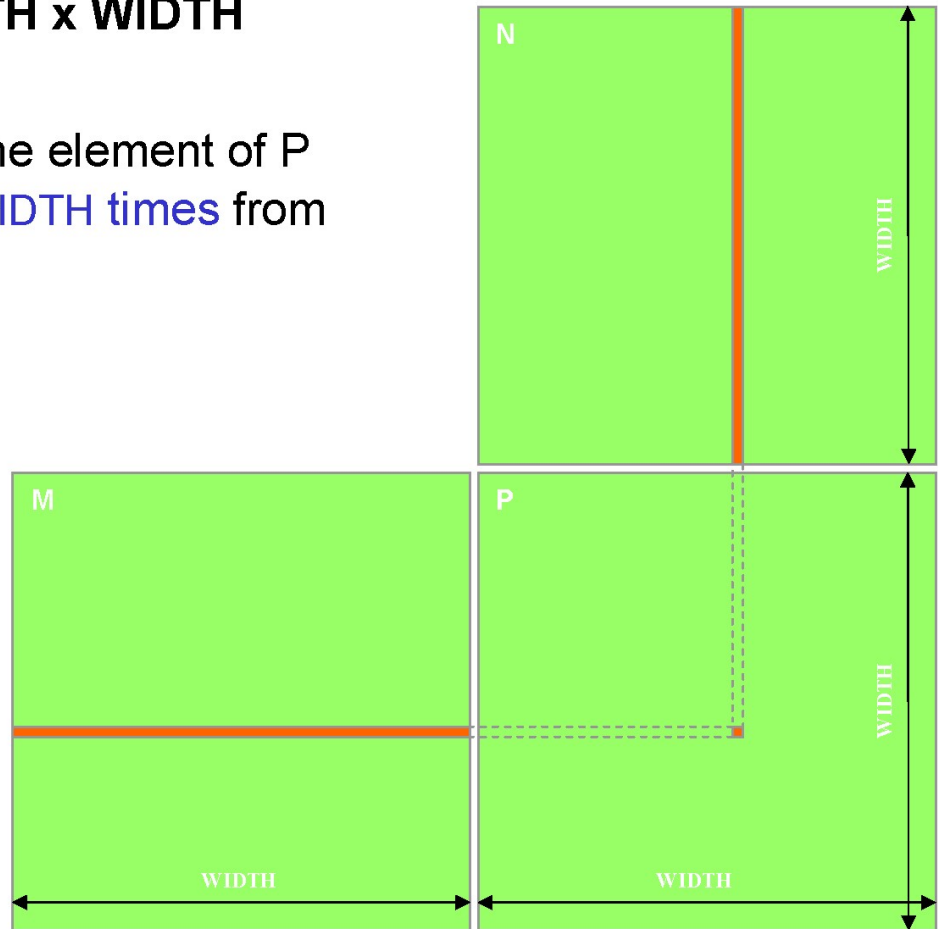# Semester Project Ideas (2)

Some ideas for topics

- Distance field computation (GPU Gems 3, chapter 34)

- Livewire ("intelligent scissors") segmentation in CUDA

- Linear systems solvers, matrix factorization (Cholesky, ...); with/without CUBLAS

- CUDA + matlab

- Fractals (Sierpinski, Koch, ...)

- Image compression

- Bilateral grid filtering for multichannel images

- Discrete wavelet transforms

- Fast histogram computations

- Terrain rendering from height map images; clipmaps or adaptive tesselation
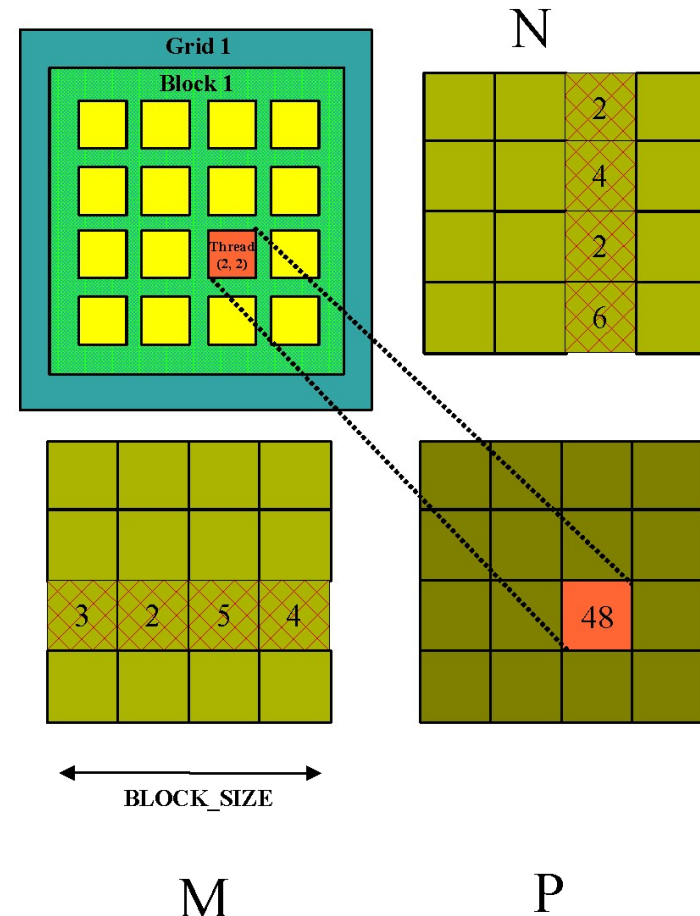
# Matrix-Matrix Multiplication

**P=MN**

Hendrik Lensch and Robert Strzodka

# Programming Model: Square Matrix Multiplication

- **P = M * N of size WIDTH x WIDTH**

- **Without tiling:**
  - One thread handles one element of P
  - M and N are loaded WIDTH times from global memory

# Multiply Using One Thread Block

- **One block of threads computes matrix P**
  - Each thread computes one element of P

- **Each thread**
  - Loads a row of matrix M
  - Loads a column of matrix N
  - Perform one multiply and addition for each pair of M and N elements
  - Compute to off-chip memory access ratio close to 1:1 (not very high)

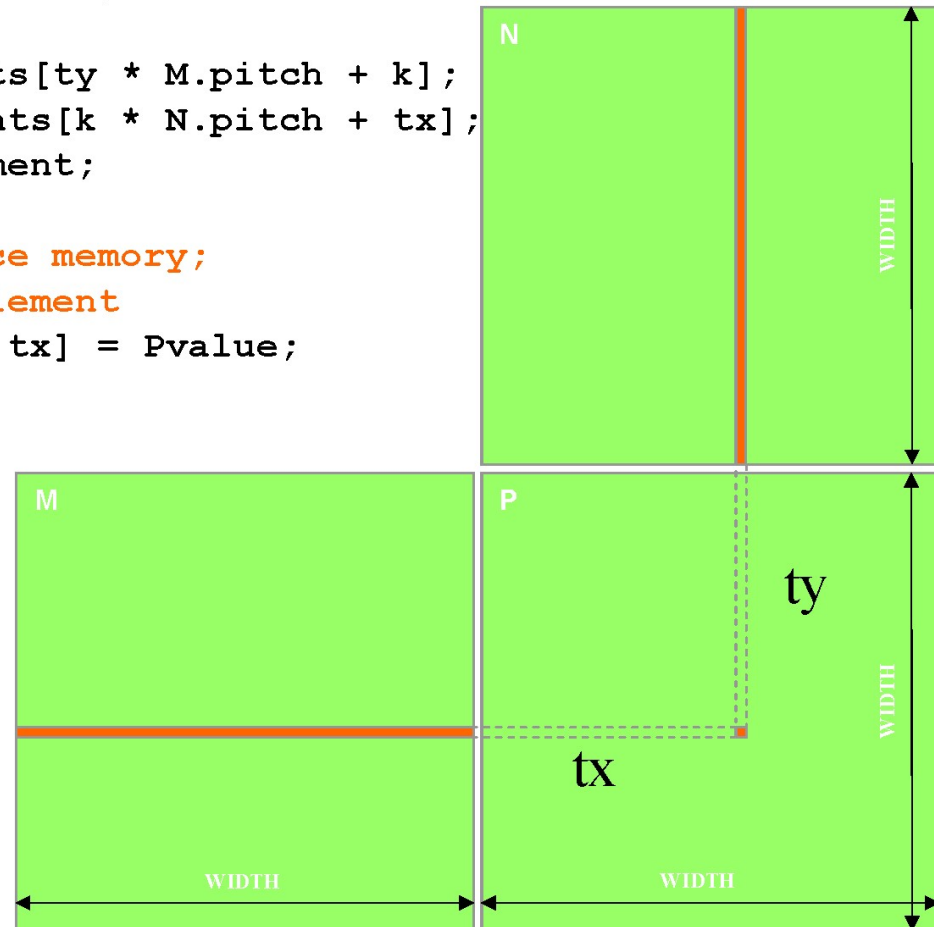- **Size of matrix limited by the number of threads allowed in a thread block**

N

Grid 1

Block 1

Thread (2, 2)

2

4

2

6

3 | 2 | 5 | 4

48

BLOCK_SIZE

M

P

# Matrix Multiplication
## Device-Side Kernel Function  (cont.)
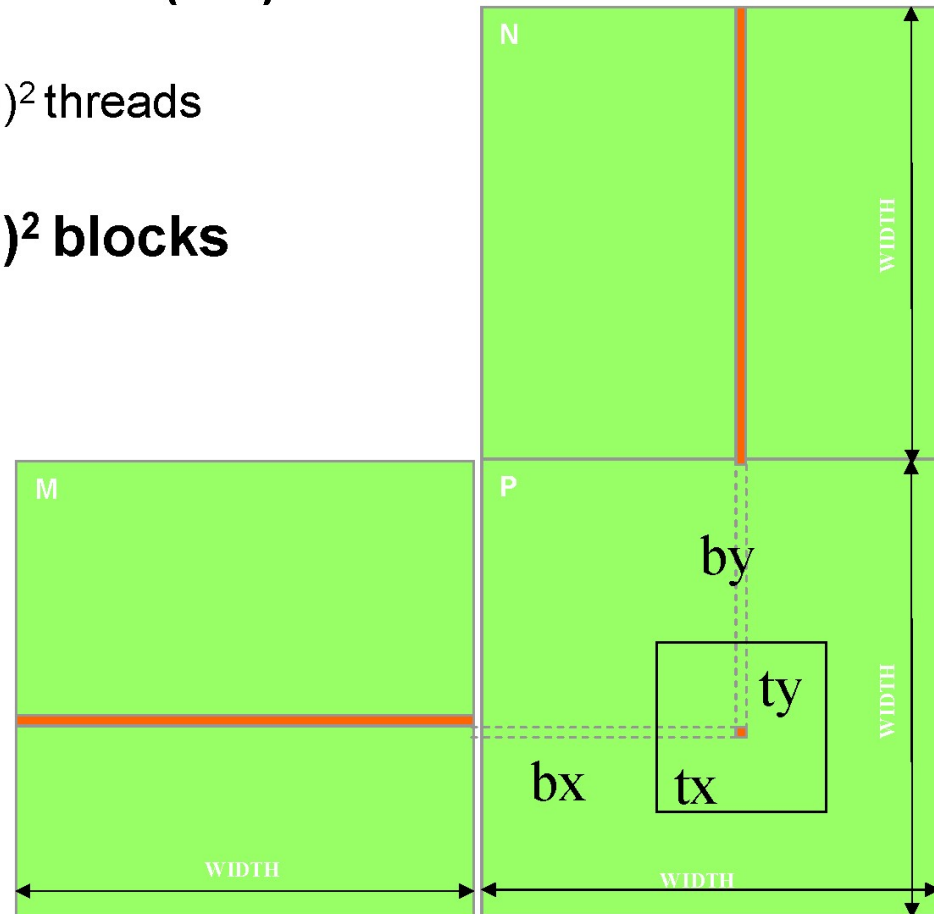
…
```
      for (int k = 0; k < M.width; ++k)
      {
         float Melement = M.elements[ty * M.pitch + k];
         float Nelement = Nd.elements[k * N.pitch + tx];
         Pvalue += Melement * Nelement;
      }
      // Write the matrix to device memory;
      // each thread writes one element
      P.elements[ty * blockDim.x+ tx] = Pvalue;
}
```
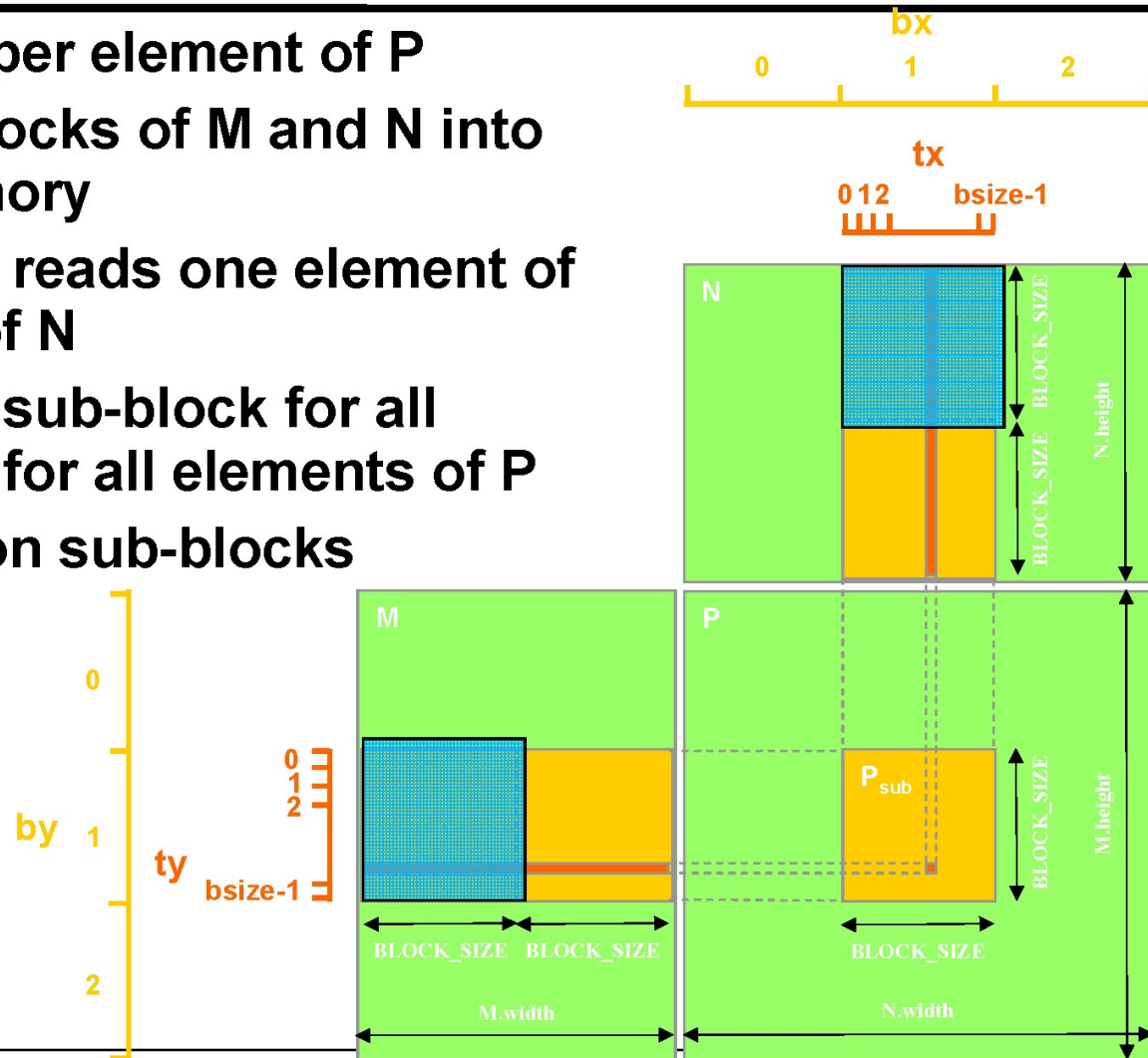
# Handling Arbitrary Sized Square Matrices

- **Have each 2D thread block to compute a (BLOCK_WIDTH)$^2$ sub-matrix (tile) of the result matrix**
  - Each has (BLOCK_WIDTH)$^2$ threads
- **Generate a 2D Grid of (WIDTH/BLOCK_WIDTH)$^2$ blocks**

You still need to put a loop around the kernel call for cases where WIDTH is greater than Max grid size!

N

M

P

by

ty

bx    tx

WIDTH

WIDTH

WIDTH

WIDTH

# Multiply Using Several Blocks - Idea

- One thread per element of P
- Load sub-blocks of M and N into shared memory
- Each thread reads one element of M and one of N
- Reuse each sub-block for all threads, i.e. for all elements of P
- Outer loop on sub-blocks
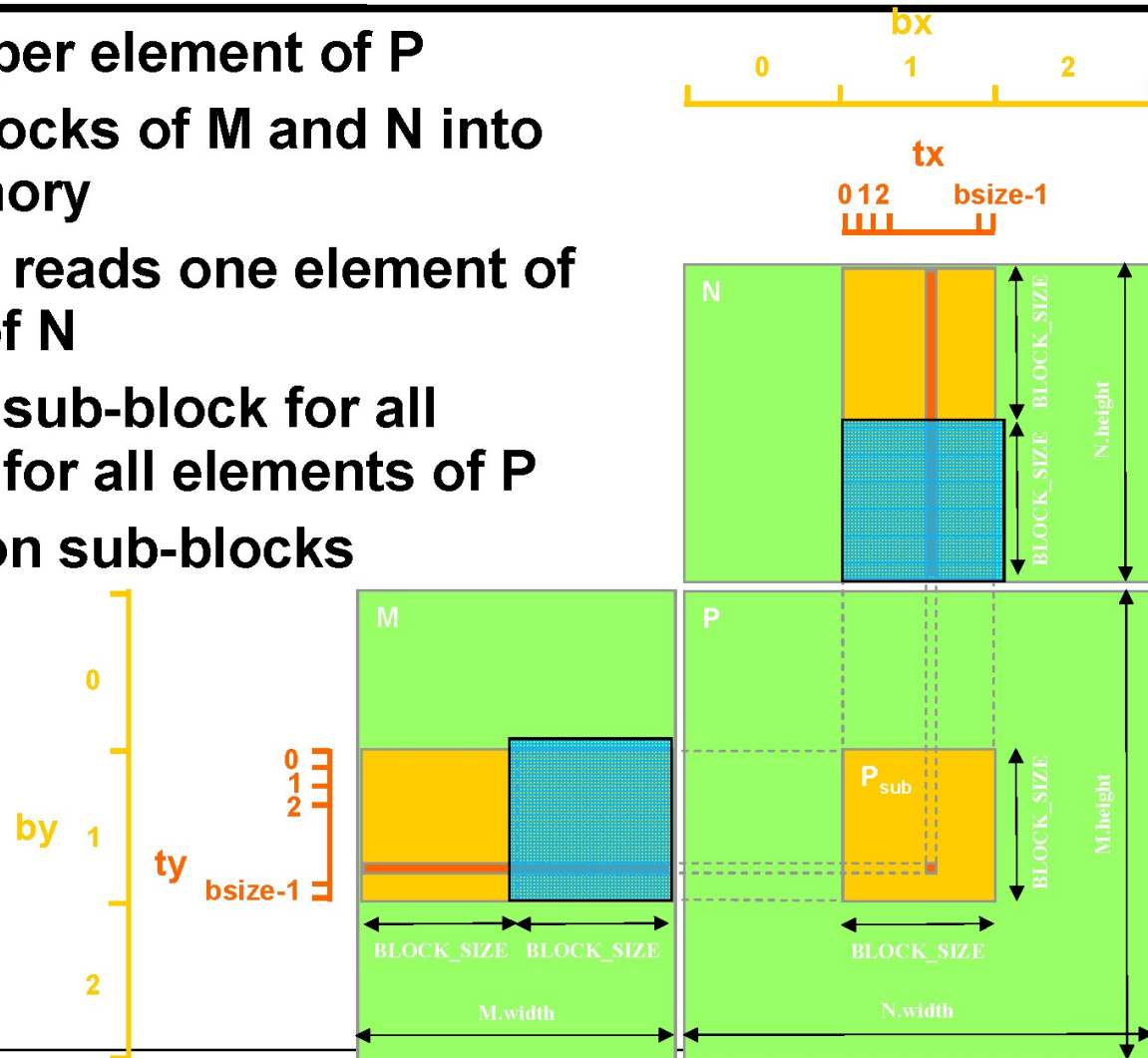
Hendrik Lensch and Robert Strzodka

# Multiply Using Several Blocks - Idea

- One thread per element of P
- Load sub-blocks of M and N into shared memory
- Each thread reads one element of M and one of N
- Reuse each sub-block for all threads, i.e. for all elements of P
- Outer loop on sub-blocks

Hendrik Lensch and Robert Strzodka

# Example: Matrix Multiplication (1)

- Copy matrices to device; invoke kernel; copy result matrix back to host

```
// Matrix multiplication - Host code
// Matrix dimensions are assumed to be multiples of BLOCK_SIZE
void MatMul(const Matrix A, const Matrix B, Matrix C)
{
    // Load A and B to device memory
    Matrix d_A;
    d_A.width = d_A.stride = A.width; d_A.height = A.height;
    size_t size = A.width * A.height * sizeof(float);
    cudaMalloc((void**)&d_A.elements, size);
    cudaMemcpy(d_A.elements, A.elements, size,
               cudaMemcpyHostToDevice);
    Matrix d_B;
    d_B.width = d_B.stride = B.width; d_B.height = B.height;
    size = B.width * B.height * sizeof(float);
    cudaMalloc((void**)&d_B.elements, size);
    cudaMemcpy(d_B.elements, B.elements, size,
               cudaMemcpyHostToDevice);
```

```
// Allocate C in device memory
Matrix d_C;
d_C.width = d_C.stride = C.width; d_C.height = C.height;
size = C.width * C.height * sizeof(float);
cudaMalloc((void**)&d_C.elements, size);

// Invoke kernel
dim3 dimBlock(BLOCK_SIZE, BLOCK_SIZE);
dim3 dimGrid(B.width / dimBlock.x, A.height / dimBlock.y);
MatMulKernel<<<dimGrid, dimBlock>>>(d_A, d_B, d_C);

// Read C from device memory
cudaMemcpy(C.elements, d_C.elements, size,
           cudaMemcpyDeviceToHost);

// Free device memory
cudaFree(d_A.elements);
cudaFree(d_B.elements);
cudaFree(d_C.elements);
}
```
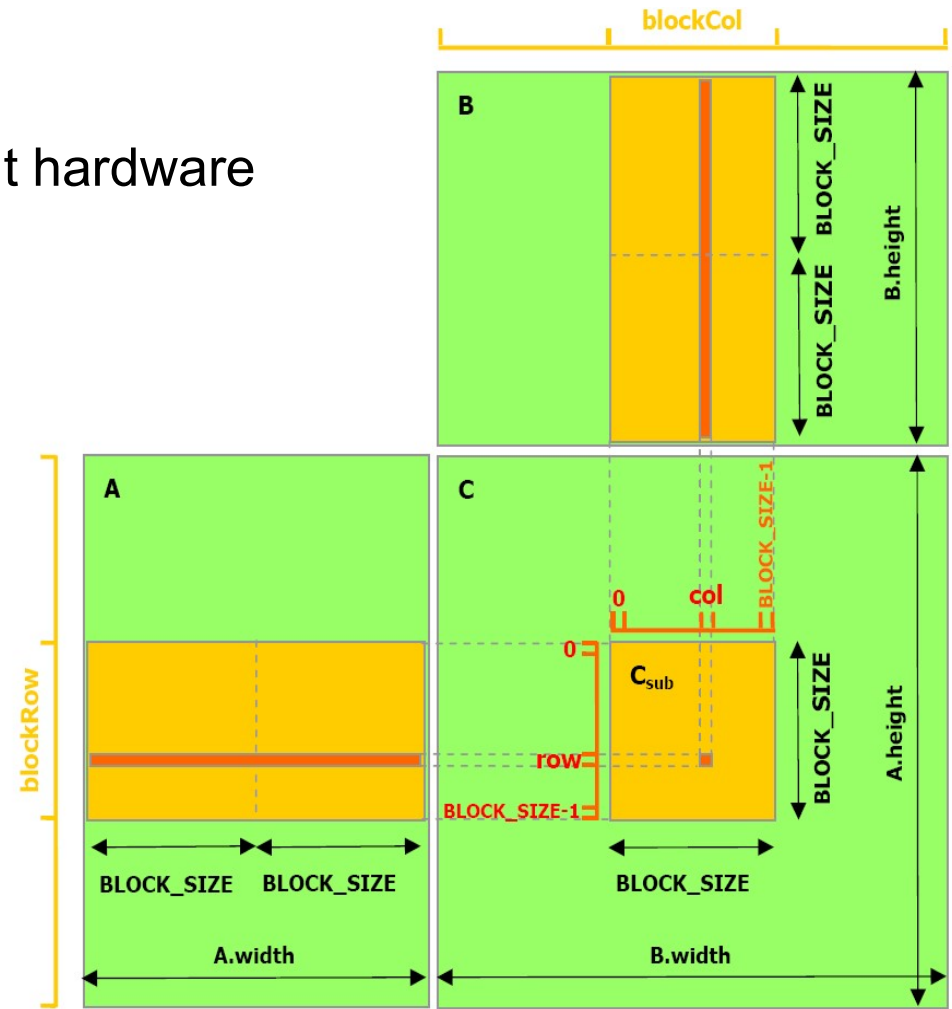
- Multiply matrix block-wise

- Set BLOCK_SIZE for efficient hardware use, e.g., to 16 on cc. 1.x or 16 or 32 on cc. 2.x +

- Maximize parallelism
  - Launch as many threads per block as block elements
  - Each thread fetches one element of block
  - Perform row * column dot products in parallel

```
__global__ void MatrixMul( float *matA, float *matB, float *matC, int w )
{
    __shared__ float blockA[ BLOCK_SIZE ][ BLOCK_SIZE ];
    __shared__ float blockB[ BLOCK_SIZE ][ BLOCK_SIZE ];

    int bx = blockIdx.x; int tx = threadIdx.x;
    int by = blockIdx.y; int ty = threadIdx.y;

    int col = bx * BLOCK_SIZE + tx;
    int row = by * BLOCK_SIZE + ty;

    float out = 0.0f;
    for ( int m = 0; m < w / BLOCK_SIZE; m++ ) {

        blockA[ ty ][ tx ] = matA[ row * w +   m * BLOCK_SIZE + tx       ];
        blockB[ ty ][ tx ] = matB[ col      + ( m * BLOCK_SIZE + ty ) * w ];
        __syncthreads();

        for ( int k = 0; k < BLOCK_SIZE; k++ ) {
            out += blockA[ ty ][ k ] * blockB[ k ][ tx ];
        }
        __syncthreads();
    }

    matC[ row * w + col ] = out;
}
```

16

Caveat: for brevity, this code assumes matrix sizes are a multiple of the block size (either because they really are, or because padding is used; otherwise guard code would need to be added)

# Memory Layout of a Matrix in C

# Memory Coalescing

- When accessing global memory, peak performance utilization occurs when all threads in a half warp (<span style="color:red">full warp on Fermi</span>) access continuous memory locations.

- <span style="color:red">Requirements relaxed on >=1.2 devices; L1 cache on Fermi!</span>

Not coalesced          coalesced



Md

Nd

Thread 1
Thread 2

WIDTH

WIDTH

18

# Memory Layout of a Matrix in C

Access direction in Kernel code



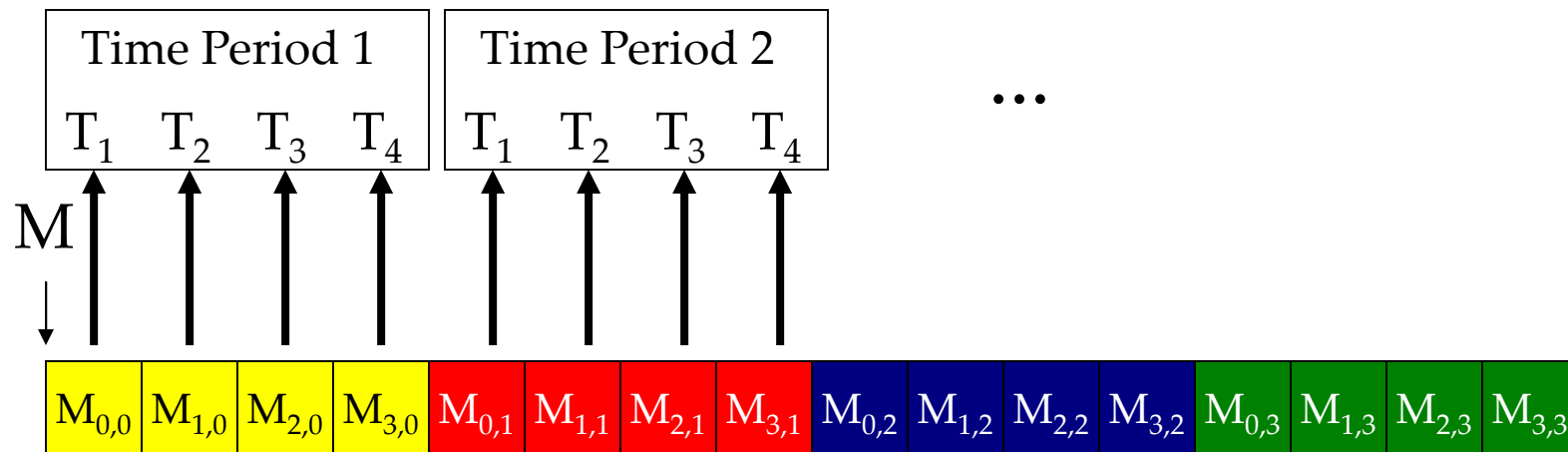| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ |
| $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ |
| $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ |
| $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |

| Time Period 1 | Time Period 2 |
| $T_1$  $T_2$  $T_3$  $T_4$ | $T_1$  $T_2$  $T_3$  $T_4$ |

...

M

| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ | $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ | $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ | $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |

# Memory Layout of a Matrix in C

Access direction in Kernel code

| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ |
|---|---|---|---|
| $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ |
| $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ |
| $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |

...

Time Period 2

$T_1$     $T_2$     $T_3$     $T_4$

Time Period 1

$T_1$     $T_2$     $T_3$     $T_4$

M

| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ | $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ | $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ | $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Thank you.