# The CUDA Programming Model

# CUDA

- "Compute Unified Device Architecture"
  - General purpose programming model
  - User kicks off batches of threads on the GPU
- GPU = dedicated super-threaded, massively data parallel co-processor
- Targeted software stack
  - Compute oriented drivers, language, and tools
- Driver for loading computation programs into GPU
  - Standalone Driver - Optimized for computation
  - Interface designed for compute – graphics-free API
  - Guaranteed maximum download & readback speeds
  - Explicit GPU memory management

# CUDA Devices and Threads

- A compute device
  - Is a coprocessor to the CPU or host
  - Has its own DRAM (device memory)
  - Runs many threads in parallel
  - Is typically a GPU but can also be another type of parallel processing device
- Data-parallel portions of an application are expressed as device kernels which run on many threads
- Differences between GPU and CPU threads
  - GPU threads are extremely lightweight
    - Very little creation overhead
  - GPU needs 1000s of threads for full efficiency
    - Multi-core CPU needs only a few

# Extended C

- Declspecs
  - global, device, shared, local, constant

- Keywords
  - threadIdx, blockIdx

- Intrinsics
  - __syncthreads

- Runtime API
  - Memory, symbol, execution management

- Function launch

```c
__device__ float filter[N];

__global__ void convolve (float *image)  {

  __shared__ float region[M];

  ...

  region[threadIdx] = image[i];

  __syncthreads();

  ...

  image[j] = result;
}

// Allocate GPU memory
void *myimage = cudaMalloc(bytes);



// 100 blocks, 10 threads per block
convolve<<<100, 10>>> (myimage);
```
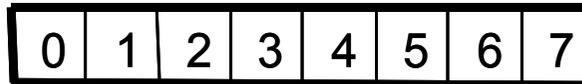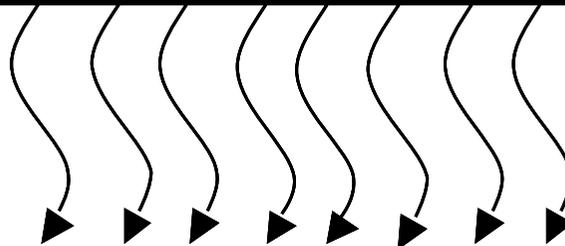
# Arrays of Parallel Threads

- A CUDA kernel is executed by an array of threads
  - All threads run the same code (SPMD)
  - Each thread has an ID that it uses to compute memory addresses and make control decisions
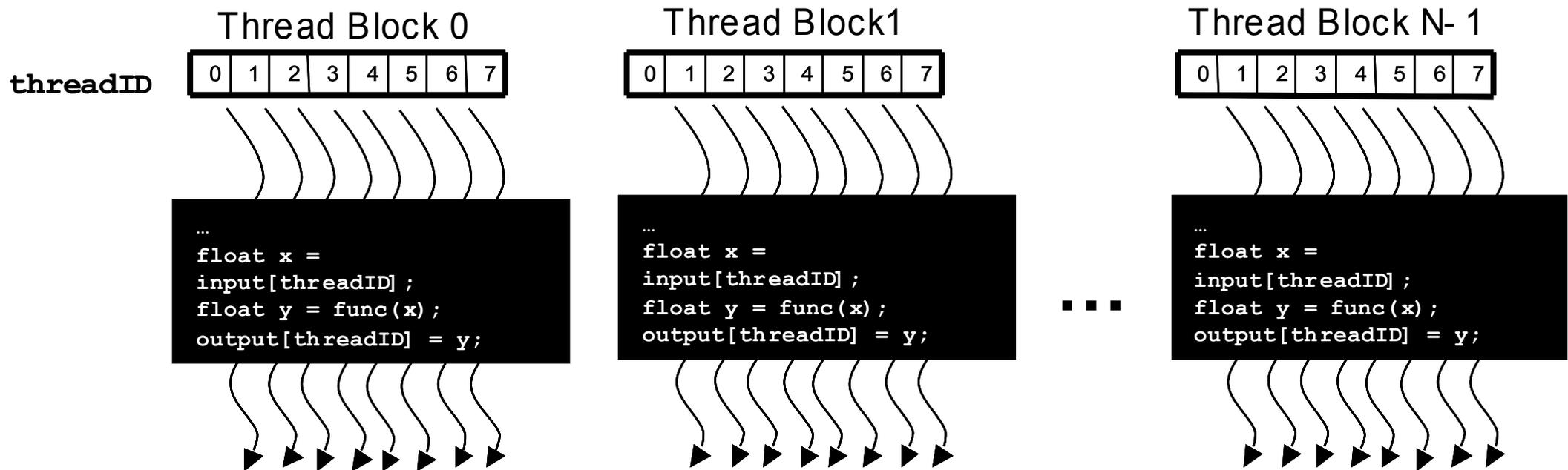
**threadID**

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|

```
…
float x = input[threadID ];
float y = func(x);
output[threadID ] = y;
…
```
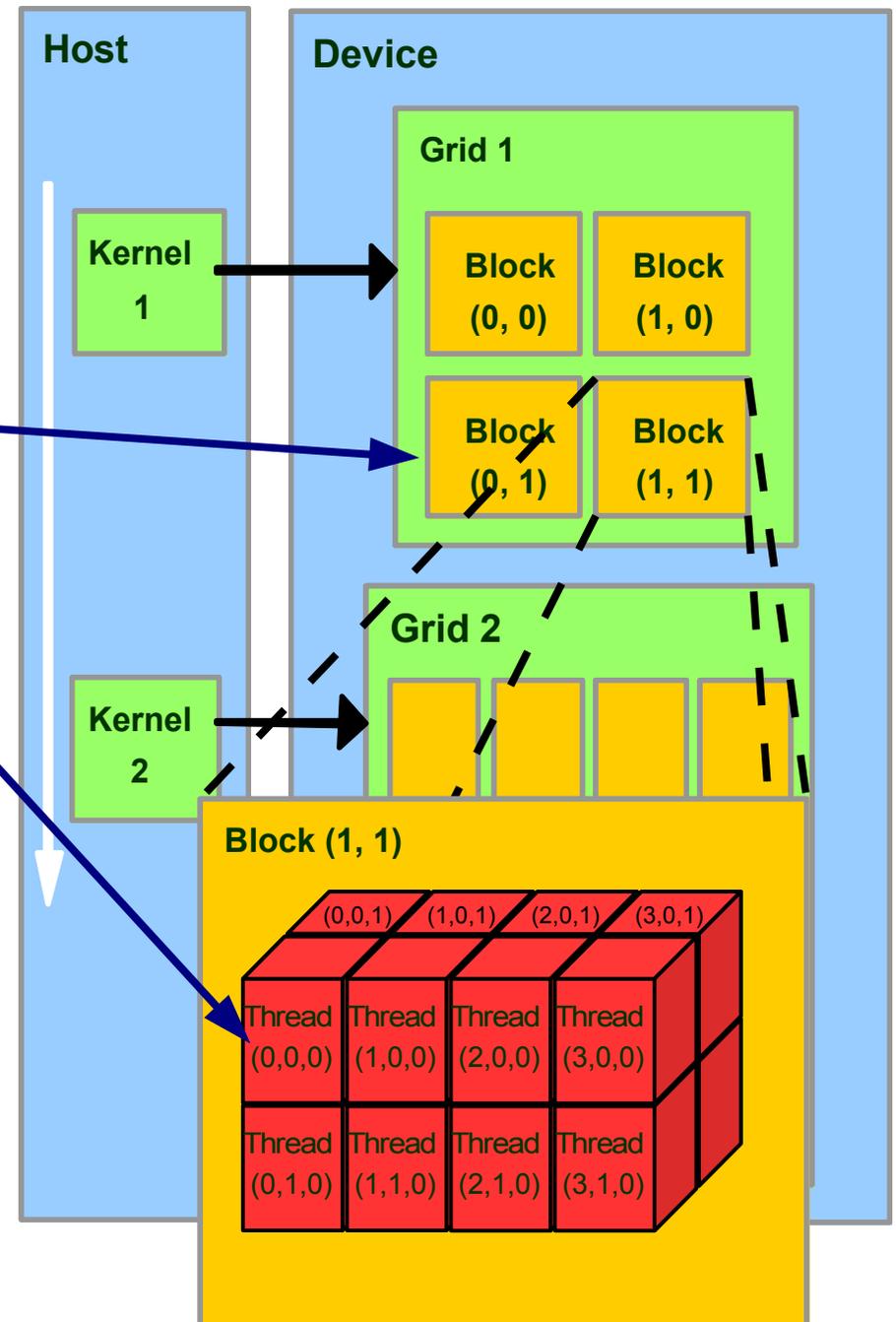
# Thread Blocks: Scalable Cooperation

- Divide monolithic thread array into multiple blocks
- Threads within a block cooperate via shared memory, atomic operations and barrier synchronization
- Threads in different blocks cannot cooperate

Thread Block 0

threadID | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

```
…
float x =
input[threadID];
float y = func(x);
output[threadID] = y;
```

Thread Block1

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

```
…
float x =
input[threadID];
float y = func(x);
output[threadID] = y;
```

Thread Block N- 1

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

```
…
float x =
input[threadID];
float y = func(x);
output[threadID] = y;
```
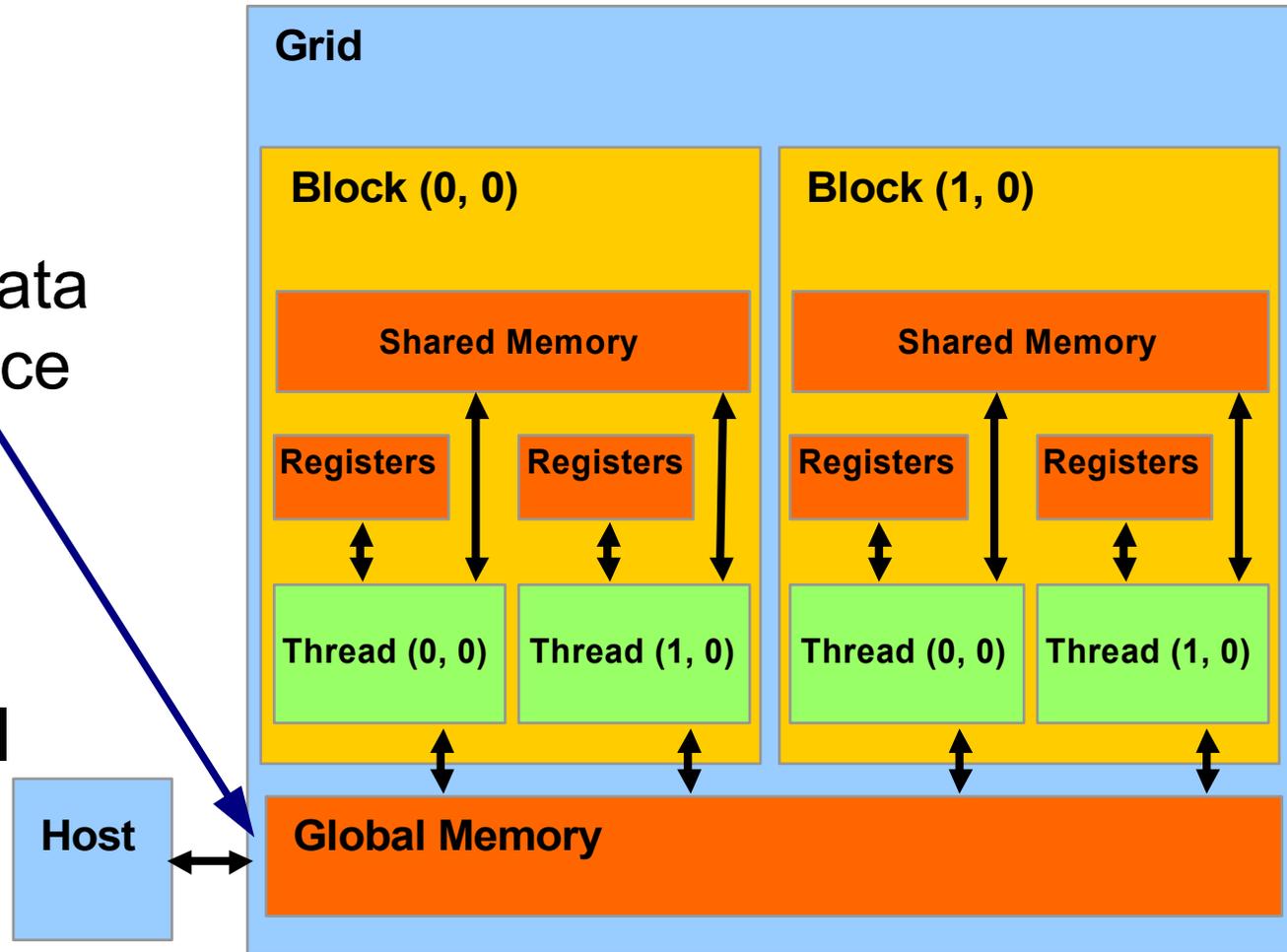
# Block IDs and Thread IDs

- Each thread uses IDs to decide what data to work on
  - Block ID: 1D or 2D
  - Thread ID: 1D, 2D, or 3D

- Simplifies memory addressing when processing multidimensional data
  - Image processing
  - Solving PDEs on volumes
  - …

**Host**

**Device**

**Grid 1**

**Kernel 1**

Block (0, 0)  Block (1, 0)

Block (0, 1)  Block (1, 1)

**Grid 2**

**Kernel 2**

**Block (1, 1)**

(0,0,1) (1,0,1) (2,0,1) (3,0,1)

Thread (0,0,0)  Thread (1,0,0)  Thread (2,0,0)  Thread (3,0,0)

Thread (0,1,0)  Thread (1,1,0)  Thread (2,1,0)  Thread (3,1,0)

# CUDA Memory Model Overview

- Global memory
  - Main means of communicating R/W Data between host and device
  - Contents visible to all threads
  - Long latency access
- We will focus on global memory for now
  - Constant and texture memory will come later

# CUDA API Highlights: Easy and Lightweight

- The API is an extension to the ANSI C programming language
  - Low learning curve
- The hardware is designed to enable lightweight runtime and driver
  - High performance

# CUDA Device Memory Allocation

- cudaMalloc()
  - Allocates object in the device <u>Global Memory</u>
  - Requires two parameters
    - **Address of a pointer** to the allocated object
    - **Size of** the allocated object
- cudaFree()
  - Frees object from device Global Memory
    - Pointer to freed object

# CUDA Device Memory Allocation (cont.)
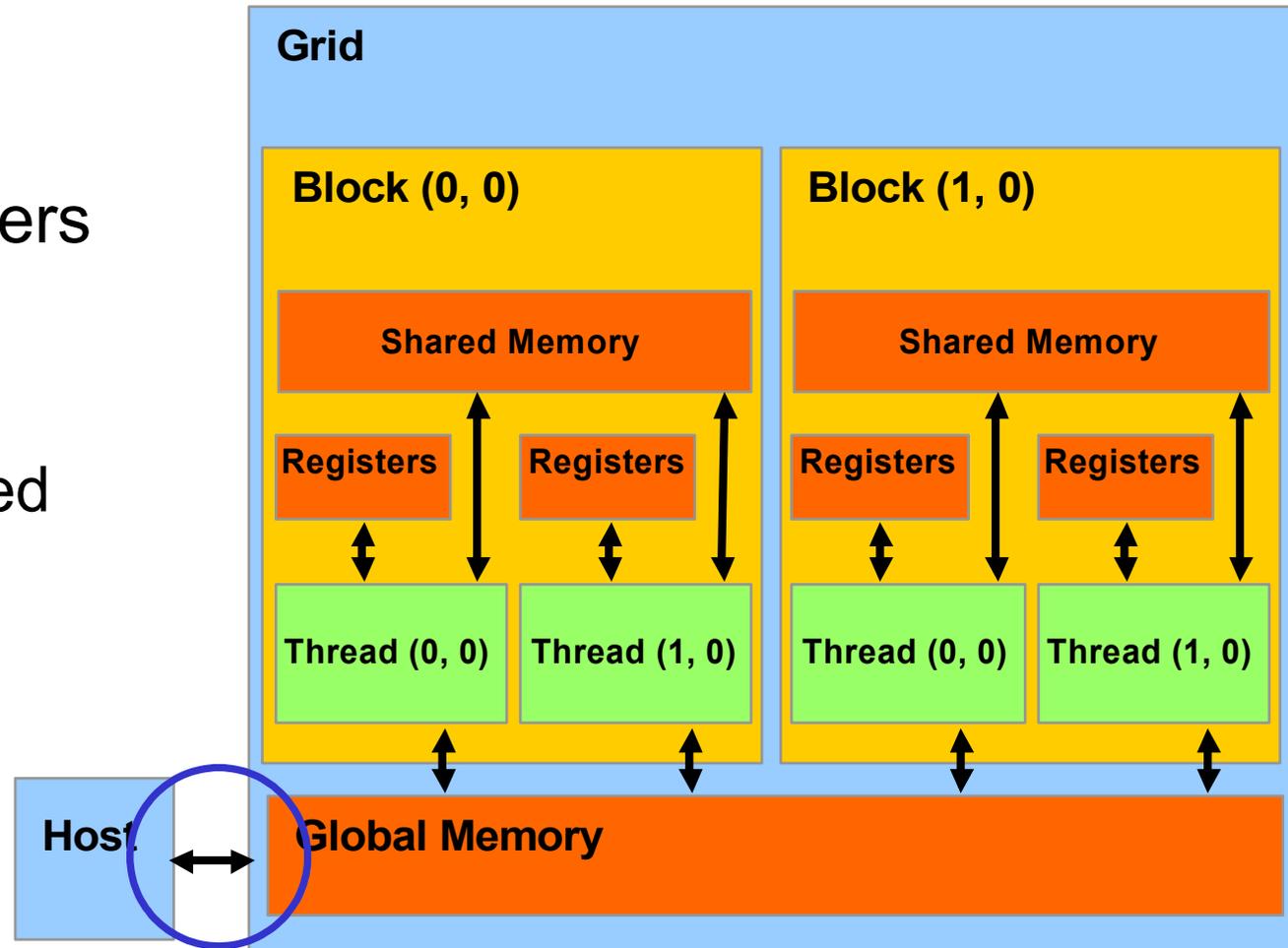
- Code example:
  - Allocate a  64 * 64 single precision float array
  - Attach the allocated storage to Md
  - "d" is often used to indicate a device data structure

```
TILE_WIDTH = 64;
float* Md
int size = TILE_WIDTH * TILE_WIDTH * sizeof(float);

cudaMalloc((void**)&Md, size);
//...
cudaFree(Md);
```

# CUDA Host-Device Data Transfer

- cudaMemcpy()
  - memory data transfer
  - Requires four parameters
    - Pointer to destination
    - Pointer to source
    - Number of bytes copied
    - Type of transfer
      - Host to Host
      - Host to Device
      - Device to Host
      - Device to Device
- Asynchronous transfer

**Grid**

**Block (0, 0)**

Shared Memory

Registers   Registers

Thread (0, 0)   Thread (1, 0)

**Block (1, 0)**

Shared Memory

Registers   Registers

Thread (0, 0)   Thread (1, 0)

**Host**

**Global Memory**

# CUDA Host-Device Data Transfer(cont.)

- Code example:
  - Transfer a  64 * 64 single precision float array
  - M is in host memory and Md is in device memory
  - cudaMemcpyHostToDevice and cudaMemcpyDeviceToHost are symbolic constants

```
cudaMemcpy(Md, M, size, cudaMemcpyHostToDevice);

cudaMemcpy(M, Md, size, cudaMemcpyDeviceToHost);
```

# CUDA Function Declarations

|  | Executed on the: | Only callable from the: |
|---|---|---|
| `__device__` `float DeviceFunc()` | device | device |
| `__global__` `void  KernelFunc()` | device | host |
| `__host__`   `float HostFunc()` | host | host |

- **`__global__`** defines a kernel function
  - Must return void

# CUDA Function Declarations (cont.)

- **`__device__`** functions cannot have their address taken
- For functions executed on the device:
  - No recursion
  - No static variable declarations inside the function
  - No variable number of arguments

# Calling a Kernel Function – Thread Creation

- A kernel function must be called with an execution configuration:

```
__global__  void KernelFunc(...);
dim3    DimGrid(100, 50);      // 5000 thread blocks
dim3    DimBlock(4, 8, 8);     // 256 threads per block
size_t SharedMemBytes = 64; // 64 bytes of shared memory
KernelFunc<<< DimGrid, DimBlock, SharedMemBytes >>>(...);
```

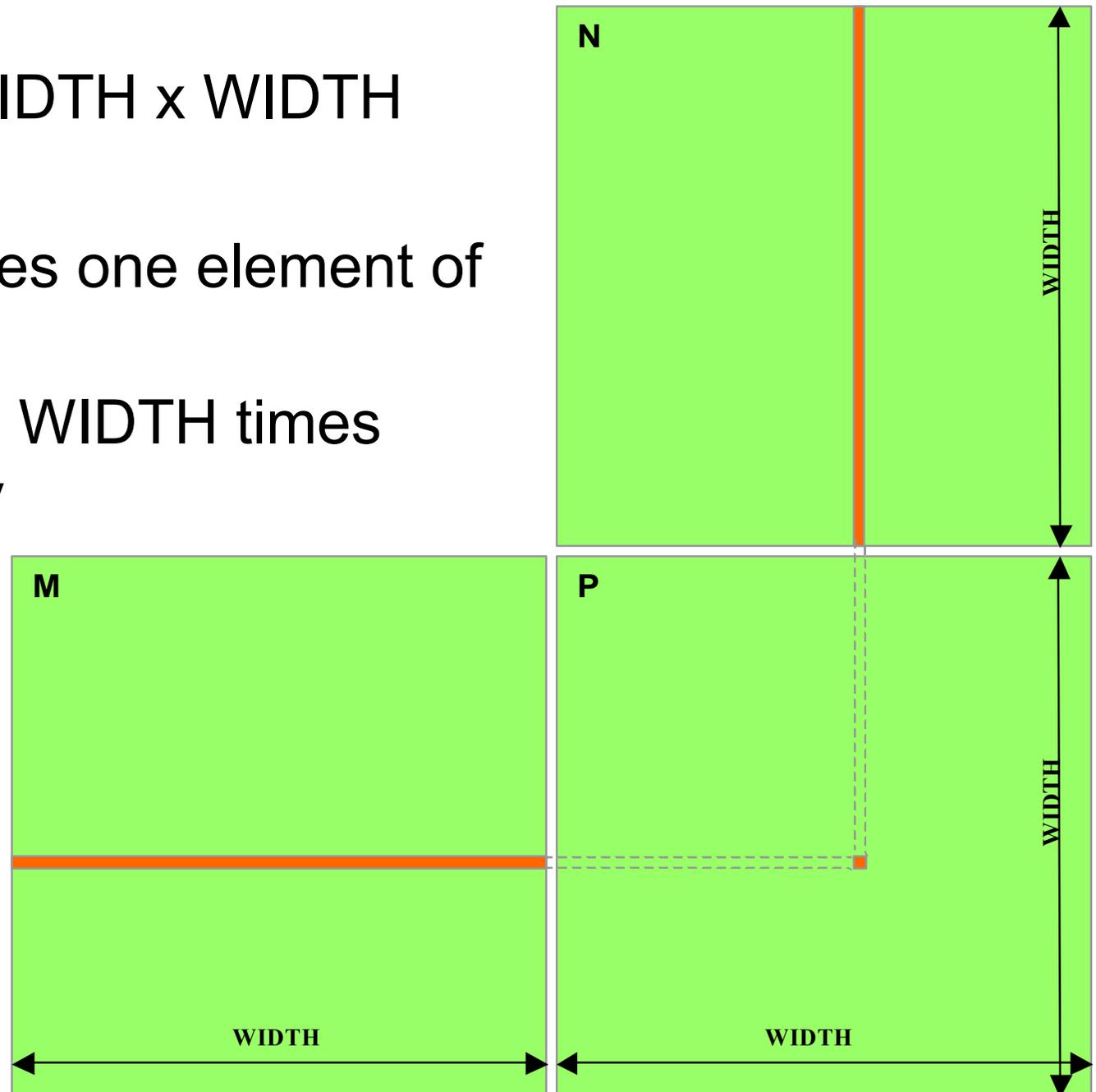- Any call to a kernel function is asynchronous, explicit synch needed for blocking

# A Simple Running Example: Matrix Multiplication

- A simple matrix multiplication example that illustrates the basic features of memory and thread management in CUDA programs
  - Leave shared memory usage until later
  - Local, register usage
  - Thread ID usage
  - Memory data transfer API between host and device
  - Assume square matrix for simplicity

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

# Memory Layout of a Matrix in C

# Step 1: Matrix Multiplication: A Simple Host Version in C

```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;
        }
}
```



N

M

P

WIDTH

WIDTH

WIDTH

WIDTH

# Step 2: Input Matrix Data Transfer (Host-side Code)

```
void MatrixMulOnDevice(float* M, float* N, float* P, int Width)
{
    int size = Width * Width * sizeof(float);
    float *Md, *Nd, *Pd;
    …
// 1. Allocate and Load M, N to device memory
    cudaMalloc(&Md, size);
    cudaMemcpy(Md, M, size, cudaMemcpyHostToDevice);

    cudaMalloc(&Nd, size);
    cudaMemcpy(Nd, N, size, cudaMemcpyHostToDevice);

    // Allocate P on the device
    cudaMalloc(&Pd, size);
```

# Step 3: Output Matrix Data Transfer (Host-side Code)
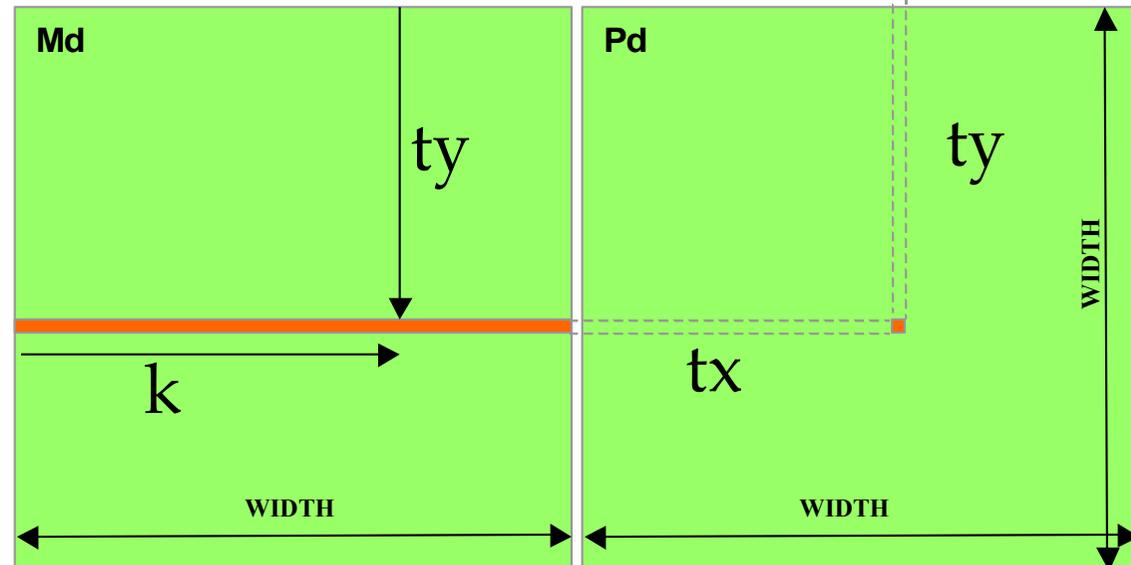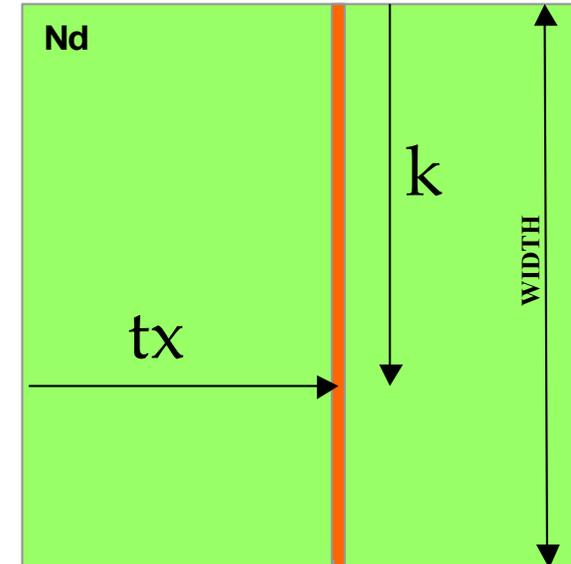
```
    // 2. Kernel invocation code – to be shown later
    …

    // 3. Read P from the device
    cudaMemcpy(P, Pd, size, cudaMemcpyDeviceToHost);

    // Free device matrices
    cudaFree(Md); cudaFree(Nd); cudaFree (Pd);
}
```

# Step 4: Kernel Function

```
// Matrix multiplication kernel – per thread code

__global__ void MatrixMulKernel
                (float* Md, float* Nd, float* Pd, int Width)
{

    // Pvalue is used to store the element of the matrix
    // that is computed by the thread
    float Pvalue = 0;
    for (int k = 0; k < Width; ++k) {
        float Melement = Md[threadIdx.y*Width+k];
        float Nelement = Nd[k*Width+threadIdx.x];
        Pvalue += Melement * Nelement;
    }


    Pd[threadIdx.y*Width+threadIdx.x] = Pvalue;
}
```

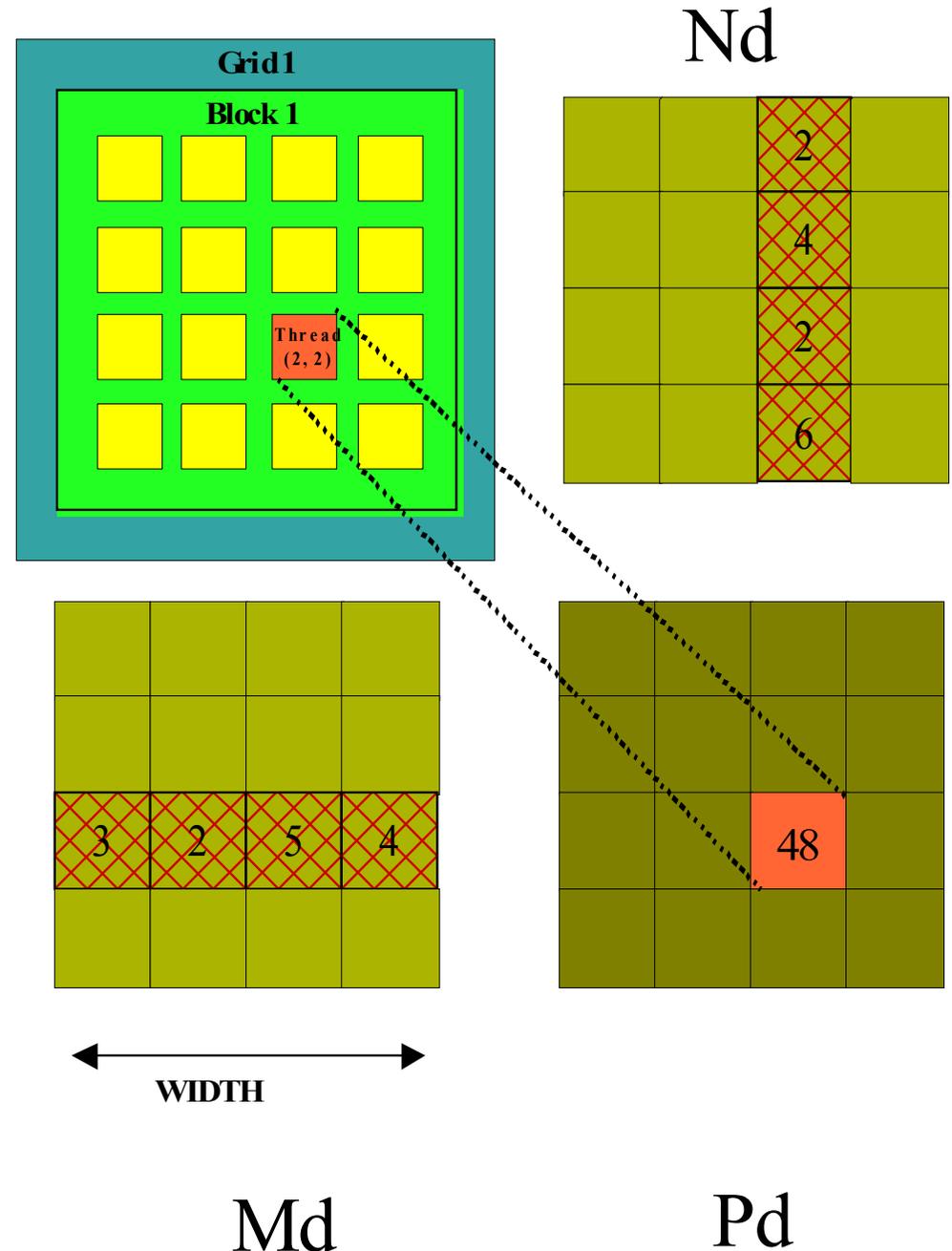Nd

k

WIDTH

tx

Md

ty

k

WIDTH

Pd

ty

tx

WIDTH

WIDTH

```
// Setup the execution configuration
dim3 dimGrid(1, 1);
dim3 dimBlock(Width, Width);


// Launch the device computation threads!
MatrixMulKernel<<<dimGrid, dimBlock>>>(Md, Nd, Pd, Width);
```
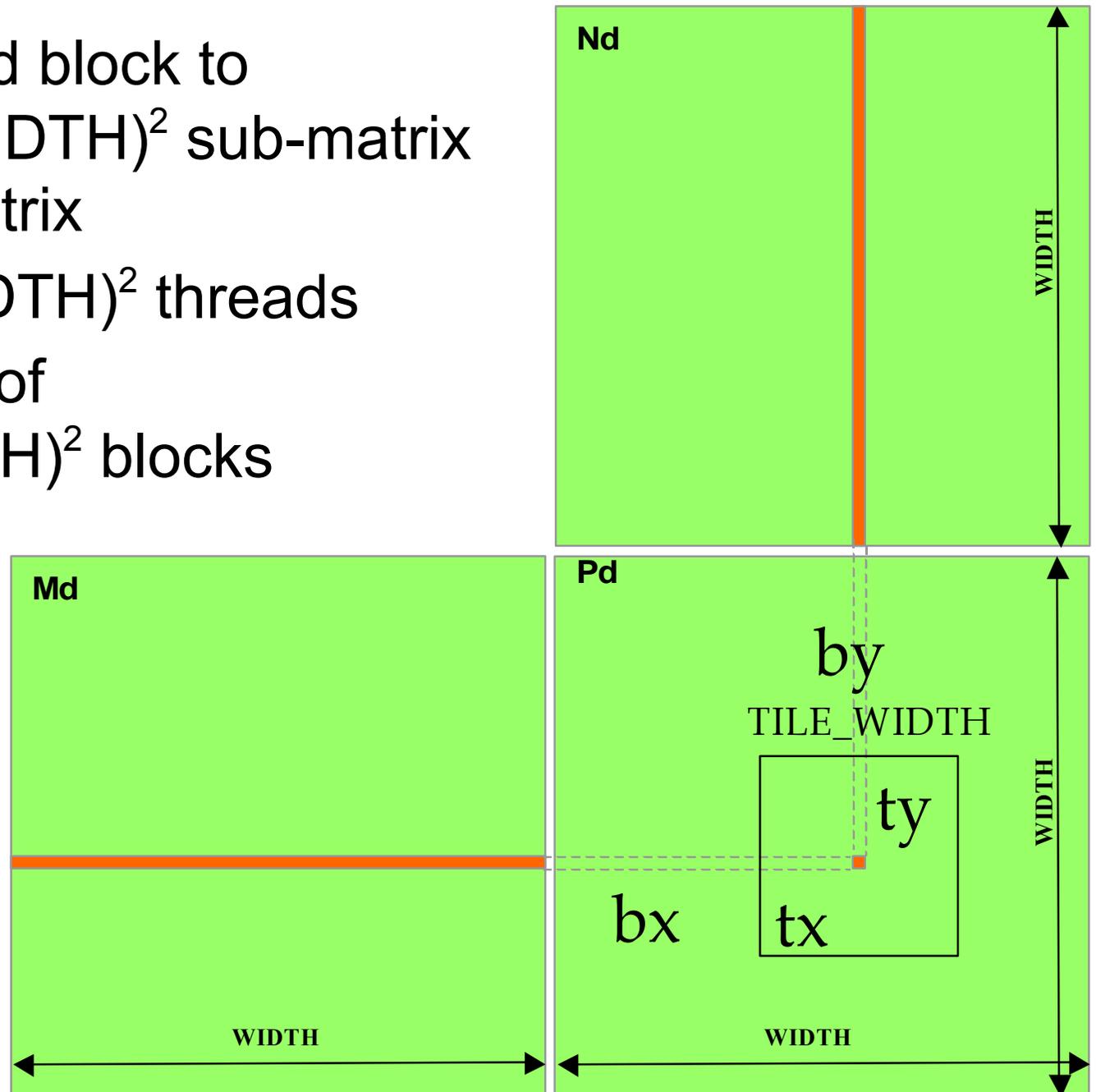
# Only One Thread Block Used

- One Block of threads compute matrix Pd

  - Each thread computes one element of Pd

- Each thread

  - Loads a row of matrix Md

  - Loads a column of matrix Nd

  - Performs one multiply and addition for each pair of Md and Nd 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

Nd

Grid 1

Block 1

Thread (2, 2)

2

4

2

6

3 2 5 4

48

WIDTH

Md

Pd

- Have each 2D thread block to compute a $(TILE\_WIDTH)^2$ sub-matrix (tile) of the result matrix

- Each has $(TILE\_WIDTH)^2$ threads

- Generate a 2D Grid of $(WIDTH/TILE\_WIDTH)^2$ blocks

You still need to put a loop around the kernel call for cases where WIDTH/TILE_WIDTH is greater than max grid size (64K)!

**Nd**

WIDTH

**Md**

**Pd**

by

TILE_WIDTH

ty

bx

tx

WIDTH

WIDTH

WIDTH

# Compiling a CUDA Program



- Parallel Thread eXecution (PTX)
  - Virtual Machine and ISA
  - Programming model
  - Execution resources and state

# Compilation

- Any source file containing CUDA language extensions must be compiled with NVCC
- NVCC is a compiler driver
  - Works by invoking all the necessary tools and compilers like cudacc, g++, cl, ...
- NVCC outputs:
  - C code (host CPU Code)
    - Must then be compiled with the rest of the application using another tool
  - PTX
    - Object code directly
    - Or, PTX source, interpreted at runtime

# Linking

- Any executable with CUDA code requires two dynamic libraries:
  - The CUDA runtime library (`cudart`)
  - The CUDA core library (`cuda`)

# Debugging Using the Device Emulation Mode

- An executable compiled in device emulation mode
  (`nvcc -deviceemu`) runs completely on the host using the
  CUDA runtime
  - No need of any device and CUDA driver
  - Each device thread is emulated with a host thread
- Running in device emulation mode, one can:
  - Use host native debug support (breakpoints, inspection, etc.)
  - Access any device-specific data from host code and vice-versa
  - Call any host function from device code (e.g. `printf`) and
    vice-versa
  - Detect deadlock situations caused by improper usage of
    `__syncthreads`

# Device Emulation Mode Pitfalls

- Emulated device threads execute sequentially, so simultaneous accesses of the same memory location by multiple threads could produce different results.

- Dereferencing device pointers on the host or host pointers on the device can produce correct results in device emulation mode, but will generate an error in device execution mode

# Floating Point

- Results of floating-point computations will slightly differ because of:
  - Different compiler outputs, instruction sets
  - Use of extended precision for intermediate results
    - There are various options to force strict single precision on the host

# Matrix Multiplication Using Multiple Blocks

- Break-up Pd into tiles
- Each block calculates one tile
- Each thread calculates one element
- Block size equal tile size

# A Small Example

Block(0,0)          Block(1,0)

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

TILE_WIDTH = 2

Block(0,1)          Block(1,1)

# A Small Example: Multiplication

# Matrix Multiplication Kernel using Multiple Blocks
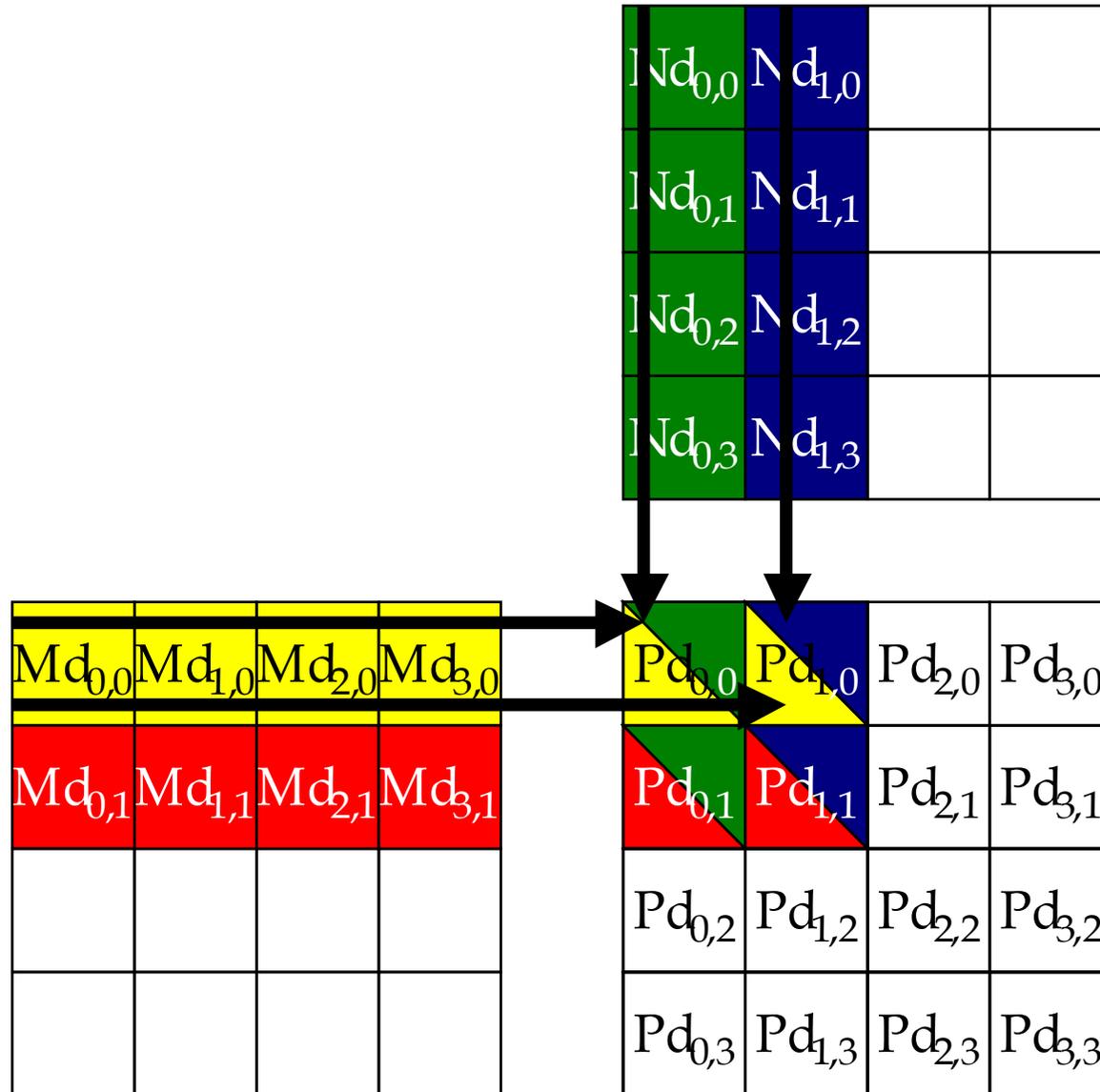
```c
__global__ void MatrixMulKernel(float* Md, float* Nd, float* Pd, int Width)
{

  // Calculate the row index of the Pd element and M
  int Row = blockIdx.y*TILE_WIDTH + threadIdx.y;

  // Calculate the column index of Pd and N
  int Col = blockIdx.x*TILE_WIDTH + threadIdx.x;

  float Pvalue = 0;
  // each thread computes one element of the block sub-matrix
  for (int k = 0; k < Width; ++k)
    Pvalue += Md[Row*Width+k] * Nd[k*Width+Col];

  Pd[Row*Width+Col] = Pvalue;
}
```
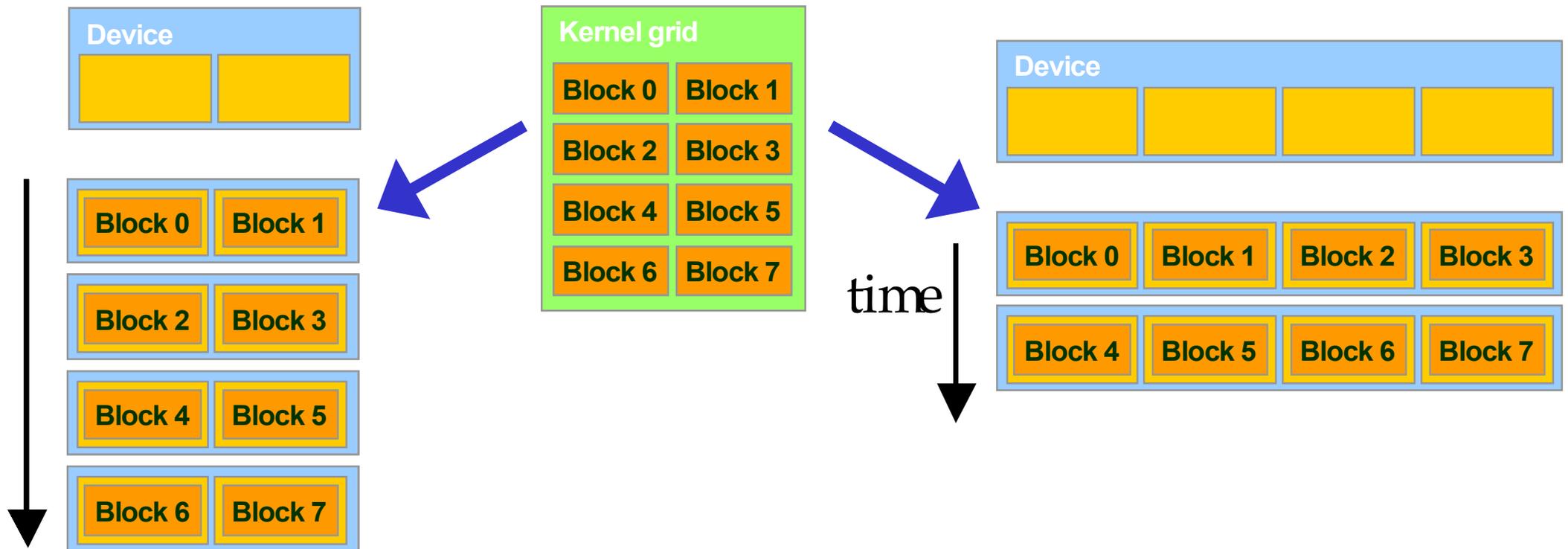
# CUDA Thread Block

- All threads in a block execute the same kernel program (SPMD)
  - Programmer declares block:
  - Block size 1 to 512 concurrent threads
  - Block shape 1D, 2D, or 3D
  - Block dimensions in threads
- Threads have thread id numbers within block
  - Thread program uses thread id to select work and address shared data
- Threads in the same block share data and synchronize while doing their share of the work
- Threads in different blocks cannot cooperate
  - Each block can execute in any order relative to other blocks!

# Transparent Scalability

- Hardware is free to assigns blocks to any processor at any time

- A kernel scales across any number of parallel processors

- Each block can execute in any order relative to other blocks.

# G80 Example: Executing Thread Blocks



**Blocks**

SM 0  SM 1

- Threads are assigned to Streaming Multiprocessors in block granularity
  - Up to 8 blocks to each SM as resource allows
  - SM in G80 can take up to 768 threads
    - Could be 256 (threads/block) * 3 blocks
    - Or 128 (threads/block) * 6 blocks, etc.
- Threads run concurrently
  - SM maintains thread/block id #s
  - SM manages/schedules thread execution

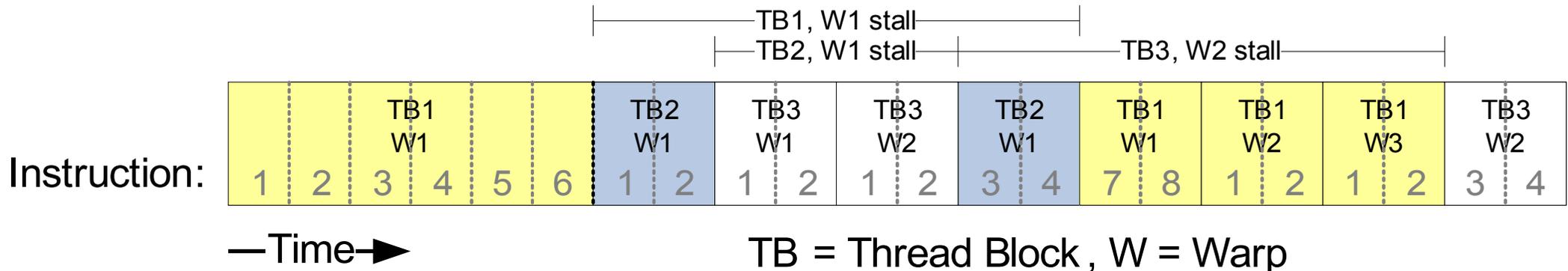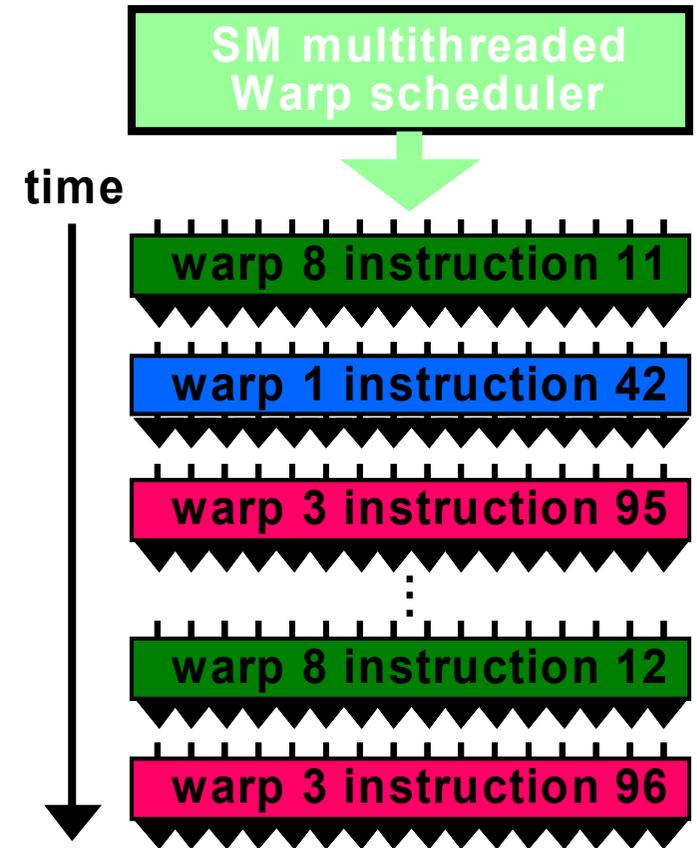# G80 Example: Thread Scheduling

- Each Block is executed as 32-thread Warps
  - An implementation decision, not part of the CUDA programming model
  - Warps are scheduling units in SM
- If 3 blocks are assigned to an SM and each block has 256 threads, how many Warps are there in an SM?
  - Each Block is divided into 256/32 = 8 Warps
  - There are 8 * 3 = 24 Warps

- SM implements zero-overhead warp scheduling
  - At any time, only one of the warps is executed by SM
  - Warps whose next instruction has its operands ready for consumption are eligible for execution
  - Eligible Warps are selected for execution on a prioritized scheduling policy
  - All threads in a warp execute the same instruction when selected
- 4 clock cycles needed to dispatch the same instruction for all threads in a Warp in G80
  - If one global memory access is needed for every 4 instructions, a minimal of 13 Warps are needed to fully tolerate 200-cycle memory latency

**time**

SM multithreaded
Warp scheduler

warp 8 instruction 11

warp 1 instruction 42

warp 3 instruction 95

warp 8 instruction 12

warp 3 instruction 96

TB1, W1 stall

TB2, W1 stall

TB3, W2 stall

| | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TB1 W1 | | | | | | TB2 W1 | | TB3 W1 | | TB3 W2 | | TB2 W1 | | TB1 W1 | | TB1 W2 | | TB1 W3 | | TB3 W2 | |
| 1 | 2 | 3 | 4 | 5 | 6 | 1 | 2 | 1 | 2 | 1 | 2 | 3 | 4 | 7 | 8 | 1 | 2 | 1 | 2 | 3 | 4 |

Instruction:

—Time→

TB = Thread Block , W = Warp

41

# G80 Block Granularity Considerations

- For Matrix Multiplication using multiple blocks, should I use 8X8, 16X16 or 32X32 blocks?

  - For 8X8, we have 64 threads per Block. Since each SM can take up to 768 threads, there are 12 Blocks. However, each SM can only take up to 8 Blocks, only 512 threads will go into each SM!

  - For 16X16, we have 256 threads per Block. Since each SM can take up to 768 threads, it can take up to 3 Blocks and achieve full capacity unless other resource considerations overrule.

  - For 32X32, we have 1024 threads per Block. Not even one can fit into an SM!

# Application Programming Interface

- The API is an extension to the C programming language

- It consists of:

  - Language extensions

    - To target portions of the code for execution on the device

  - A runtime library split into:

    - A common component providing built-in vector types and a subset of the C runtime library in both host and device codes

    - A host component to control and access one or more devices from the host

    - A device component providing device-specific functions

# Language Extensions: Built-in Variables

- **`dim3 gridDim;`**
  - Dimensions of the grid in blocks (**`gridDim.z`** unused)
- **`dim3 blockDim;`**
  - Dimensions of the block in threads
- **`dim3 blockIdx;`**
  - Block index within the grid
- **`dim3 threadIdx;`**
  - Thread index within the block

# Common Runtime Component: Mathematical Functions

- **`pow, sqrt, cbrt, hypot`**
- **`exp, exp2, expm1`**
- **`log, log2, log10, log1p`**
- **`sin, cos, tan, asin, acos, atan, atan2`**
- **`sinh, cosh, tanh, asinh, acosh, atanh`**
- **`ceil, floor, trunc, round`**
- Etc.
  - When executed on the host, a given function uses the C runtime implementation if available
  - These functions are only supported for scalar types, not vector types

# Common Runtime Component: Mathematical Functions

- Some mathematical functions (e.g. `sin(x)`) have a less accurate, but faster device-only version (e.g. `__sin(x)`)
  - `__pow`
  - `__log, __log2, __log10`
  - `__exp`
  - `__sin, __cos, __tan`

# Host Runtime Component

- Provides functions to deal with:
  - Device management (including multi-device systems)
  - Memory management
  - Error handling
- Initializes the first time a runtime function is called
- A host thread can invoke device code on only one device
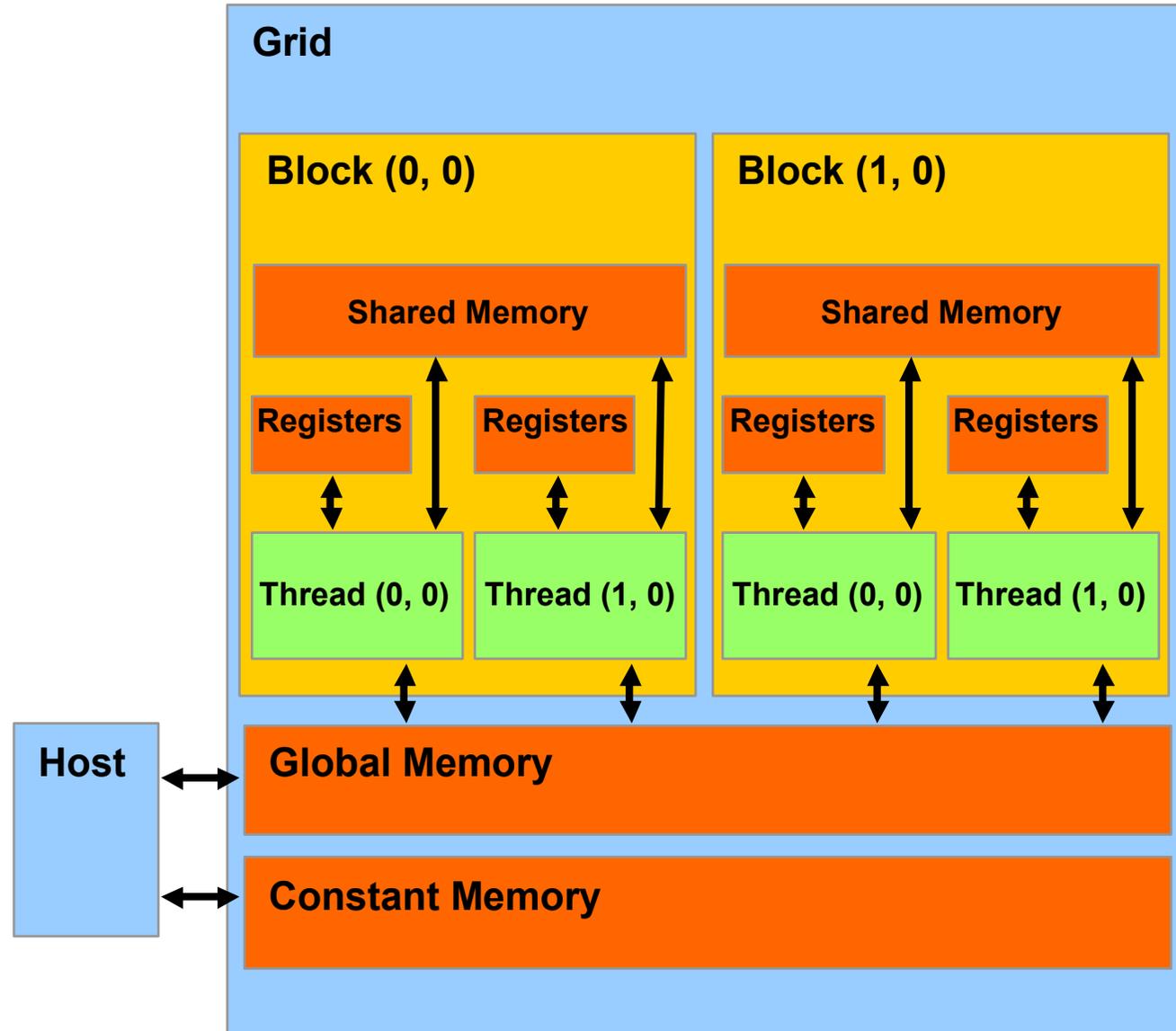  - Multiple host threads required to run on multiple devices

- **`void __syncthreads();`**
- Synchronizes all threads in a block
- Once all threads have reached this point, execution resumes normally
- Used to avoid RAW / WAR / WAW hazards when accessing shared or global memory
- Allowed in conditional constructs only if the conditional is uniform across the entire thread block

# G80 Implementation of CUDA Memories

- 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
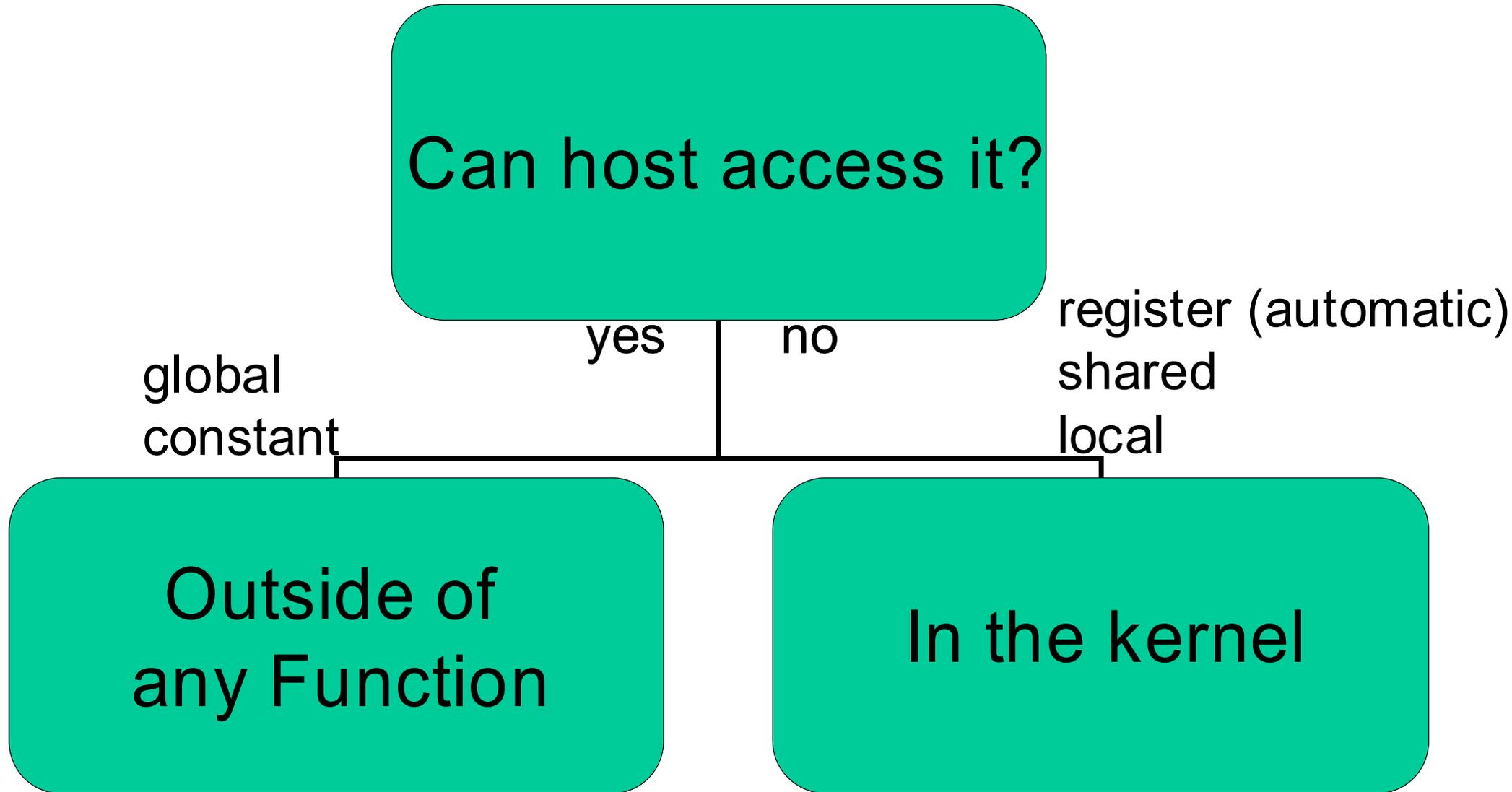  - Read/only per-grid constant memory

# CUDA Variable Type Qualifiers

| Variable declaration | Memory | Scope | Lifetime |
|---|---|---|---|
| `__device__ __local__`     `int LocalVar;` | local | thread | thread |
| `__device__ __shared__`     `int SharedVar;` | shared | block | block |
| `__device__`     `int GlobalVar;` | global | grid | application |
| `__device__ __constant__` `int ConstantVar;` | 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

# Where to Declare Variables?

Can host access it?

yes    no

global
constant

register (automatic)
shared
local

Outside of
any Function

In the kernel

# 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;`

# A Common Programming Strategy

- Global memory resides in device memory (DRAM) - much slower access than shared memory

- So, a profitable way of performing computation on the device is to tile data to take advantage of fast shared memory:

  - Partition data into subsets that fit into shared memory

  - Handle each data subset with one thread block by:

    - Loading the subset from global memory to shared memory, using multiple threads to exploit memory-level parallelism

    - Performing the computation on the subset from shared memory; each thread can efficiently multi-pass over any data element

    - Copying results from shared memory to global memory

# A Common Programming Strategy (Cont.)

- Constant memory also resides in device memory (DRAM) - much slower access than shared memory
  - But… cached!
  - Highly efficient access for read-only data
- Carefully divide data according to access patterns
  - R/Only → constant memory (very fast if in cache)
  - R/W shared within Block → shared memory (very fast)
  - R/W within each thread → registers (very fast)
  - R/W inputs/results → global memory (very slow)

# GPU Atomic Integer Operations

- Atomic operations on integers in global memory:
  - Associative operations on signed/unsigned ints
  - add, sub, min, max, ...
  - and, or, xor
  - Increment, decrement
  - Exchange, compare and swap
- Requires hardware with compute capability 1.1 and above.

# Review: Matrix Multiplication Kernel using Multiple Blocks

```
__global__ void MatrixMulKernel(float* Md, float* Nd, float* Pd, int Width)
{

  // Calculate the row index of the Pd element and M
  int Row = blockIdx.y*TILE_WIDTH + threadIdx.y;

  // Calculate the column index of Pd and N
  int Col = blockIdx.x*TILE_WIDTH + threadIdx.x;

  float Pvalue = 0;
  // each thread computes one element of the block sub-matrix
  for (int k = 0; k < Width; ++k)
    Pvalue += Md[Row*Width+k] * Nd[k*Width+Col];

  Pd[Row*Width+Col] = Pvalue;
}
```
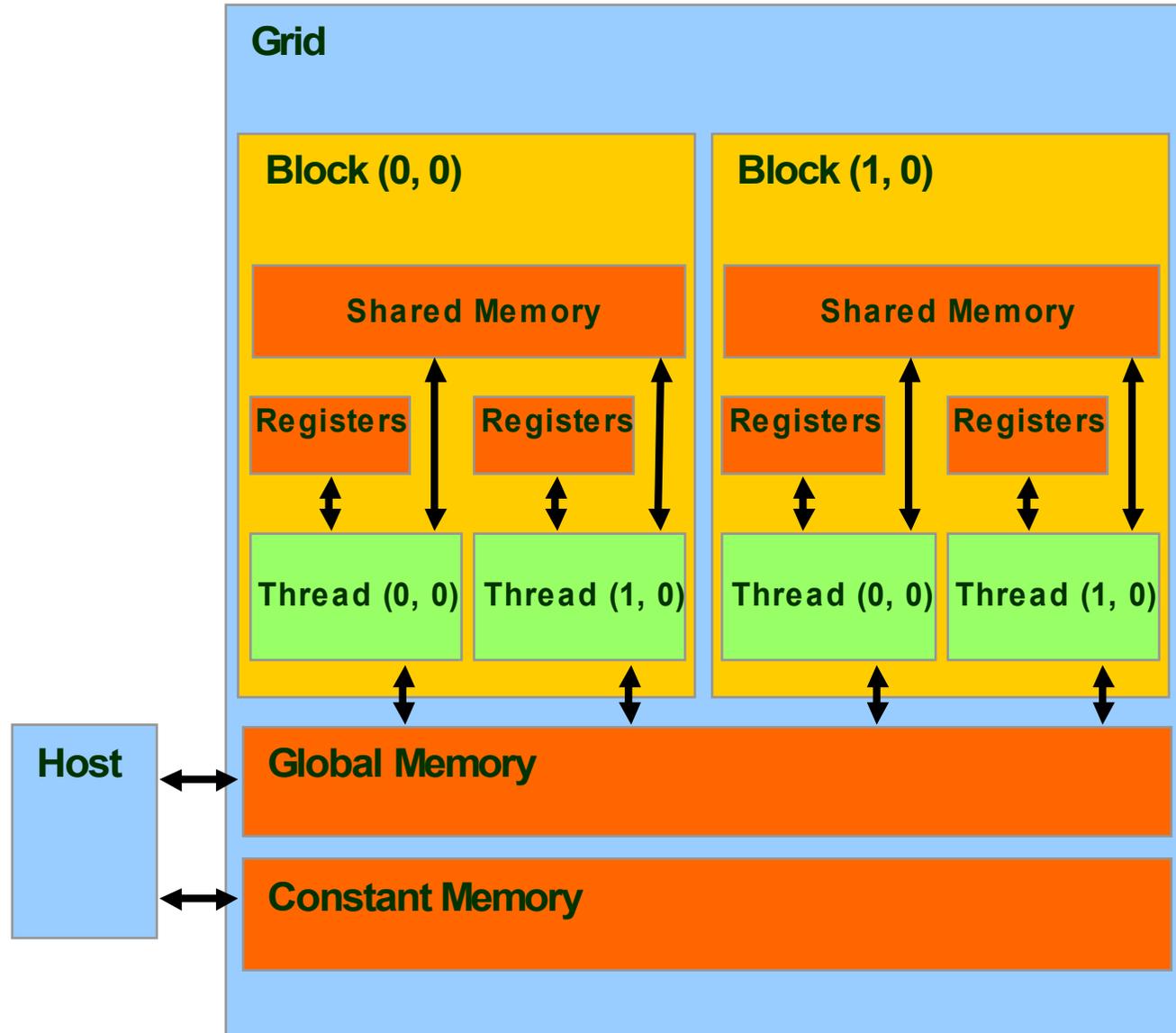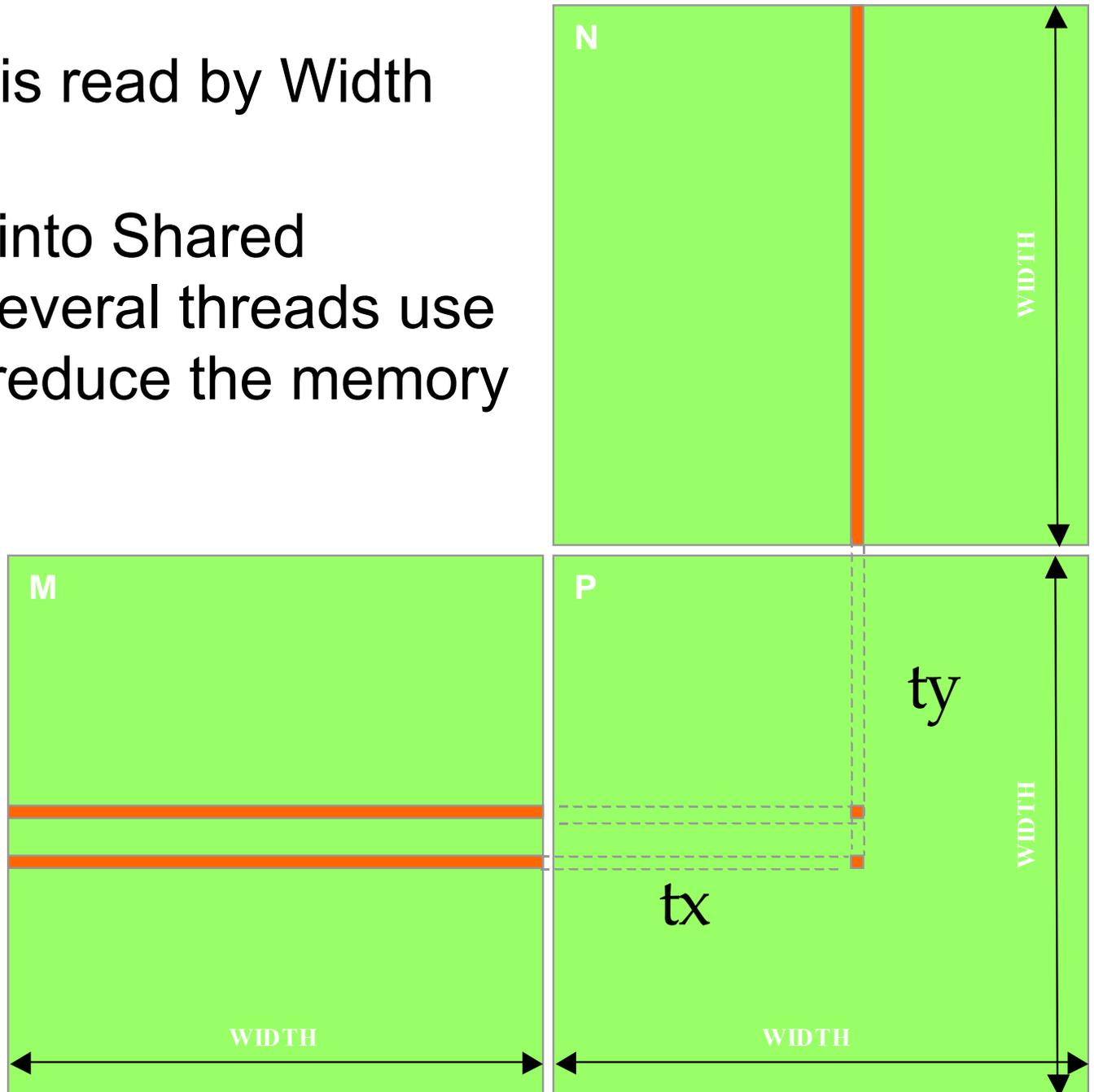
# How about performance on G80?

- All threads access global memory for their input matrix elements
  - Two memory accesses (8 bytes) per floating point multiply-add
  - 4B/s of memory bandwidth/FLOPS
  - 4*346.5 = 1386 GB/s required to achieve peak FLOP rating
  - 86.4 GB/s limits the code at 21.6 GFLOPS
- The actual code runs at about 15 GFLOPS
- Need to drastically cut down memory accesses to get closer to the peak 346.5 GFLOPS

**Grid**

**Block (0, 0)**

Shared Memory

Registers | Registers

Thread (0, 0) | Thread (1, 0)

**Block (1, 0)**

Shared Memory

Registers | Registers

Thread (0, 0) | Thread (1, 0)

**Host**

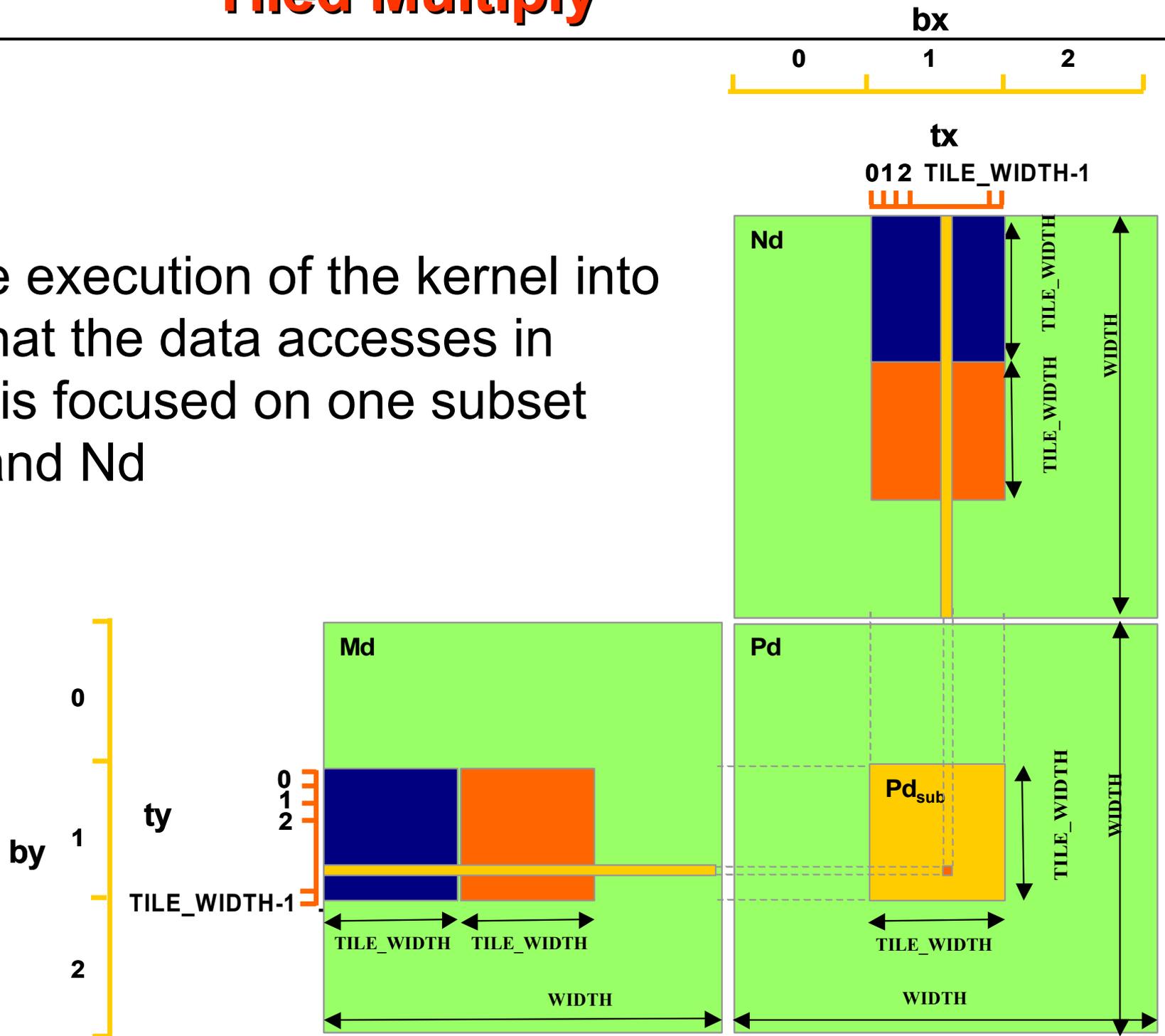**Global Memory**

**Constant Memory**

# Idea: Use Shared Memory to reuse global memory data

- Each input element is read by Width threads.
- Load each element into Shared Memory and have several threads use the local version to reduce the memory bandwidth
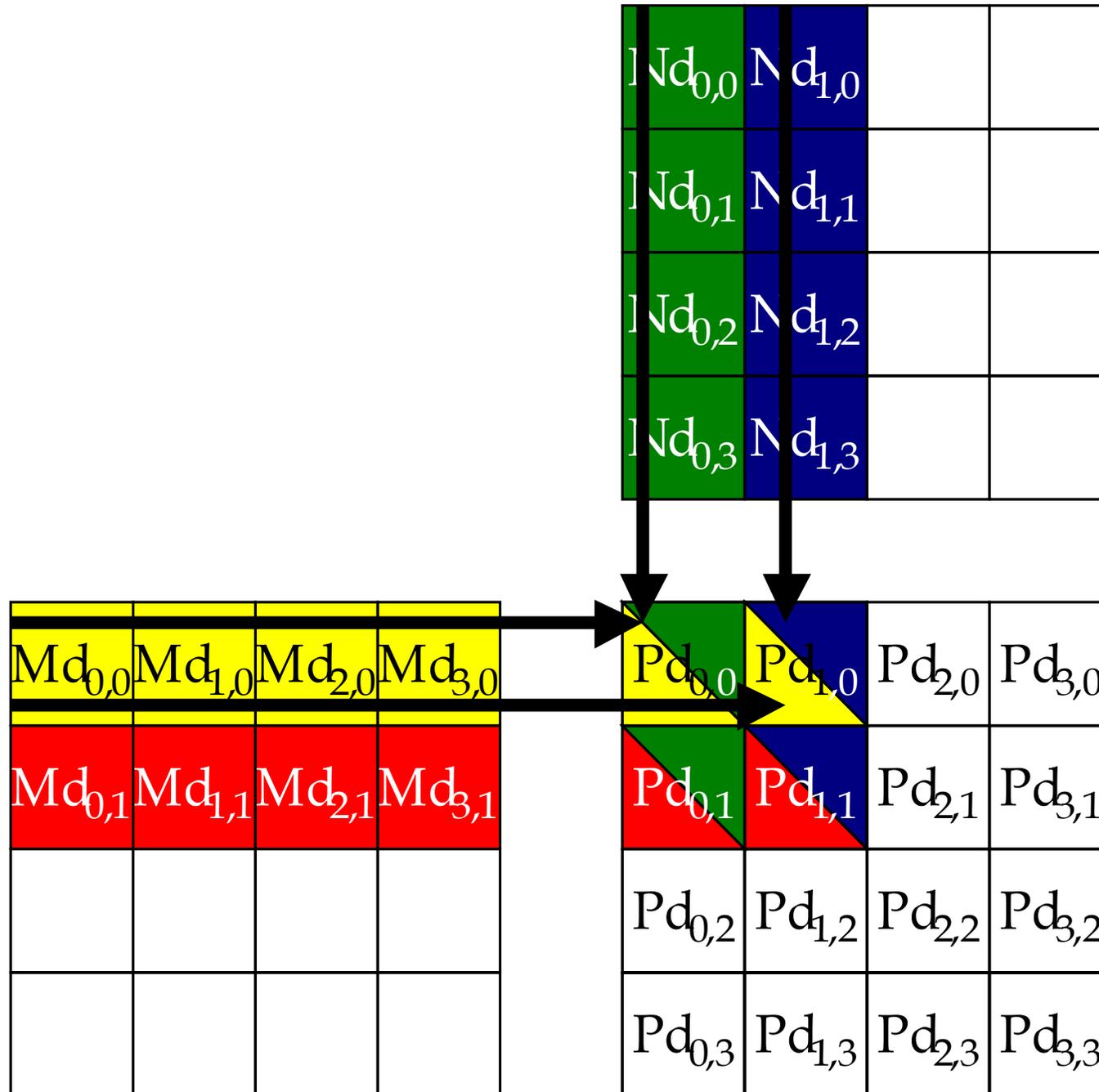- Tiled algorithms

# Tiled Multiply

- Break up the execution of the kernel into phases so that the data accesses in each phase is focused on one subset (tile) of Md and Nd

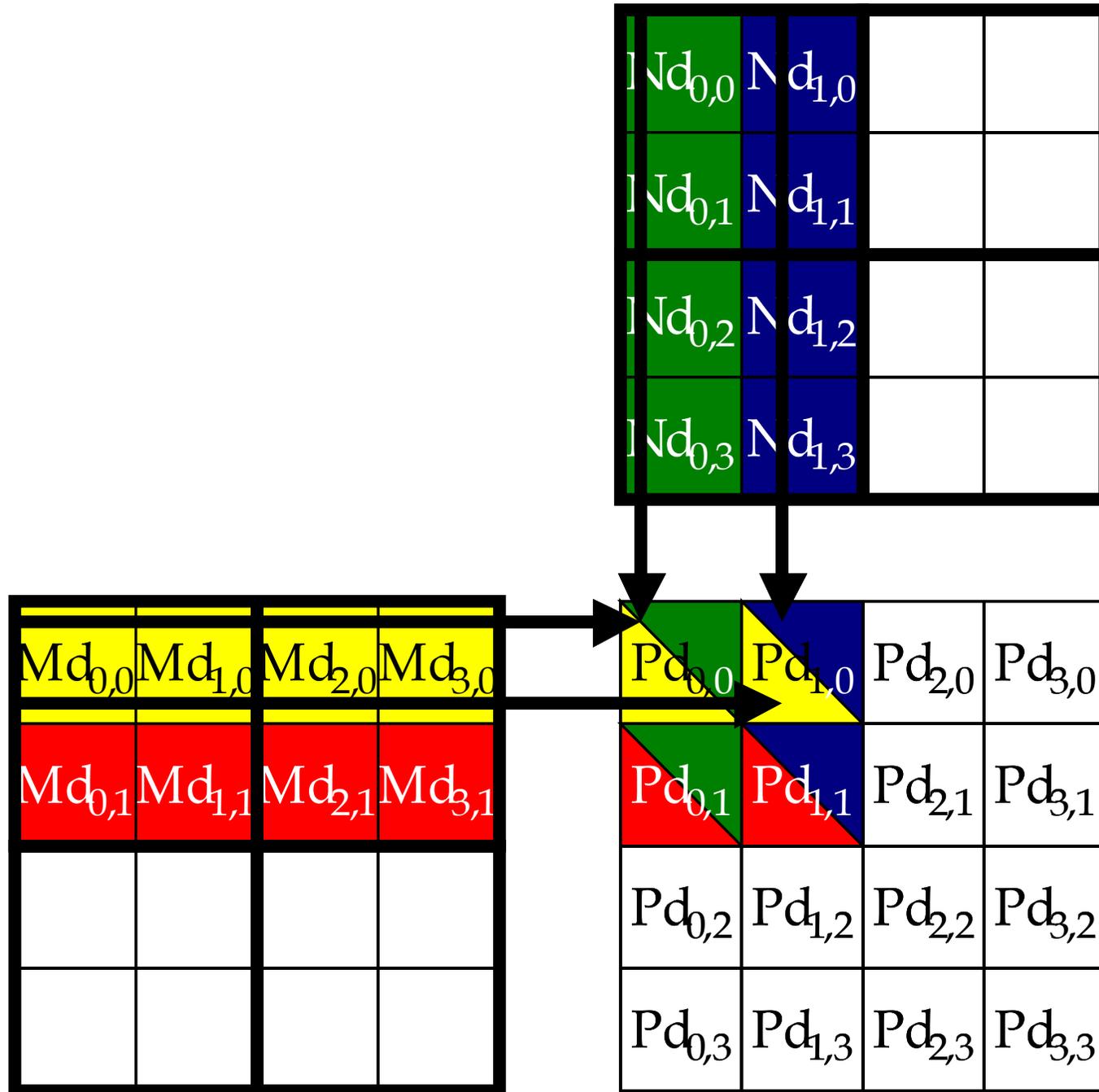# Every Md and Nd Element is used exactly twice in generating a 2X2 tile of P

| $P_{0,0}$ thread$_{0,0}$ | $P_{1,0}$ thread$_{1,0}$ | $P_{0,1}$ thread$_{0,1}$ | $P_{1,1}$ thread$_{1,1}$ |
|---|---|---|---|
| $M_{0,0} * N_{0,0}$ | $M_{0,0} * N_{1,0}$ | $M_{0,1} * N_{0,0}$ | $M_{0,1} * N_{1,0}$ |
| $M_{1,0} * N_{0,1}$ | $M_{1,0} * N_{1,1}$ | $M_{1,1} * N_{0,1}$ | $M_{1,1} * N_{1,1}$ |
| $M_{2,0} * N_{0,2}$ | $M_{2,0} * N_{1,2}$ | $M_{2,1} * N_{0,2}$ | $M_{2,1} * N_{1,2}$ |
| $M_{3,0} * N_{0,3}$ | $M_{3,0} * N_{1,3}$ | $M_{3,1} * N_{0,3}$ | $M_{3,1} * N_{1,3}$ |

Access order

# Each phase of a Thread Block uses one tile from Md and one from Nd

|  | Phase 1 | | | Phase 2 | | |
|---|---|---|---|---|---|---|
| $T_{0,0}$ | $\textbf{Md}_{0,0}$ $\downarrow$ $Mds_{0,0}$ | $\textbf{Nd}_{0,0}$ $\downarrow$ $Nds_{0,0}$ | $PValue_{0,0}$ += $Mds_{0,0}*Nds_{0,0}$ + $Mds_{1,0}*Nds_{0,1}$ | $\textbf{Md}_{2,0}$ $\downarrow$ $Mds_{0,0}$ | $\textbf{Nd}_{0,2}$ $\downarrow$ $Nds_{0,0}$ | $PValue_{0,0}$ += $Mds_{0,0}*Nds_{0,0}$ + $Mds_{1,0}*Nds_{0,1}$ |
| $T_{1,0}$ | $\textbf{Md}_{1,0}$ $\downarrow$ $Mds_{1,0}$ | $\textbf{Nd}_{1,0}$ $\downarrow$ $Nds_{1,0}$ | $PValue_{1,0}$ += $Mds_{0,0}*Nds_{1,0}$ + $Mds_{1,0}*Nds_{1,1}$ | $\textbf{Md}_{3,0}$ $\downarrow$ $Mds_{1,0}$ | $\textbf{Nd}_{1,2}$ $\downarrow$ $Nds_{1,0}$ | $PValue_{1,0}$ += $Mds_{0,0}*Nds_{1,0}$ + $Mds_{1,0}*Nds_{1,1}$ |
| $T_{0,1}$ | $\textbf{Md}_{0,1}$ $\downarrow$ $Mds_{0,1}$ | $\textbf{Nd}_{0,1}$ $\downarrow$ $Nds_{0,1}$ | $PdValue_{0,1}$ += $Mds_{0,1}*Nds_{0,0}$ + $Mds_{1,1}*Nds_{0,1}$ | $\textbf{Md}_{2,1}$ $\downarrow$ $Mds_{0,1}$ | $\textbf{Nd}_{0,3}$ $\downarrow$ $Nds_{0,1}$ | $PdValue_{0,1}$ += $Mds_{0,1}*Nds_{0,0}$ + $Mds_{1,1}*Nds_{0,1}$ |
| $T_{1,1}$ | $\textbf{Md}_{1,1}$ $\downarrow$ $Mds_{1,1}$ | $\textbf{Nd}_{1,1}$ $\downarrow$ $Nds_{1,1}$ | $PdValue_{1,1}$ += $Mds_{0,1}*Nds_{1,0}$ + $Mds_{1,1}*Nds_{1,1}$ | $\textbf{Md}_{3,1}$ $\downarrow$ $Mds_{1,1}$ | $\textbf{Nd}_{1,3}$ $\downarrow$ $Nds_{1,1}$ | $PdValue_{1,1}$ += $Mds_{0,1}*Nds_{1,0}$ + $Mds_{1,1}*Nds_{1,1}$ |

time →

# First-order Size Considerations in G80

- Each thread block should have many threads
  - TILE_WIDTH of 16 gives 16*16 = 256 threads
- There should be many thread blocks
  - A 1024*1024 Pd 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

# CUDA Code – Kernel Execution Configuration

```
// Setup the execution configuration
dim3 dimBlock(TILE_WIDTH, TILE_WIDTH);
dim3 dimGrid(Width  / TILE_WIDTH,
             Width /  TILE_WIDTH);
```

# Tiled Matrix Multiplication Kernel

```c
__global__ void MatrixMulKernel(float* Md, float* Nd, float* Pd, int Width)
{
    __shared__ float Mds[TILE_WIDTH][TILE_WIDTH];
    __shared__ float Nds[TILE_WIDTH][TILE_WIDTH];

    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
    int Row = by * TILE_WIDTH + ty;
    int Col = bx * TILE_WIDTH + tx;

    float Pvalue = 0;
// Loop over the Md and Nd tiles required to compute the Pd element
    for (int m = 0; m < Width/TILE_WIDTH; ++m) {
// Coolaborative loading of Md and Nd tiles into shared memory
        Mds[ty][tx] = Md[Row*Width + (m*TILE_WIDTH + tx)];
        Nds[ty][tx] = Nd[Col + (m*TILE_WIDTH + ty)*Width];
        __syncthreads();

    .   for (int k = 0; k < TILE_WIDTH; ++k)
            Pvalue += Mds[ty][k] * Nds[k][tx];
        __syncthreads();
    }
    Pd[Row*Width+Col] = Pvalue;
}
```
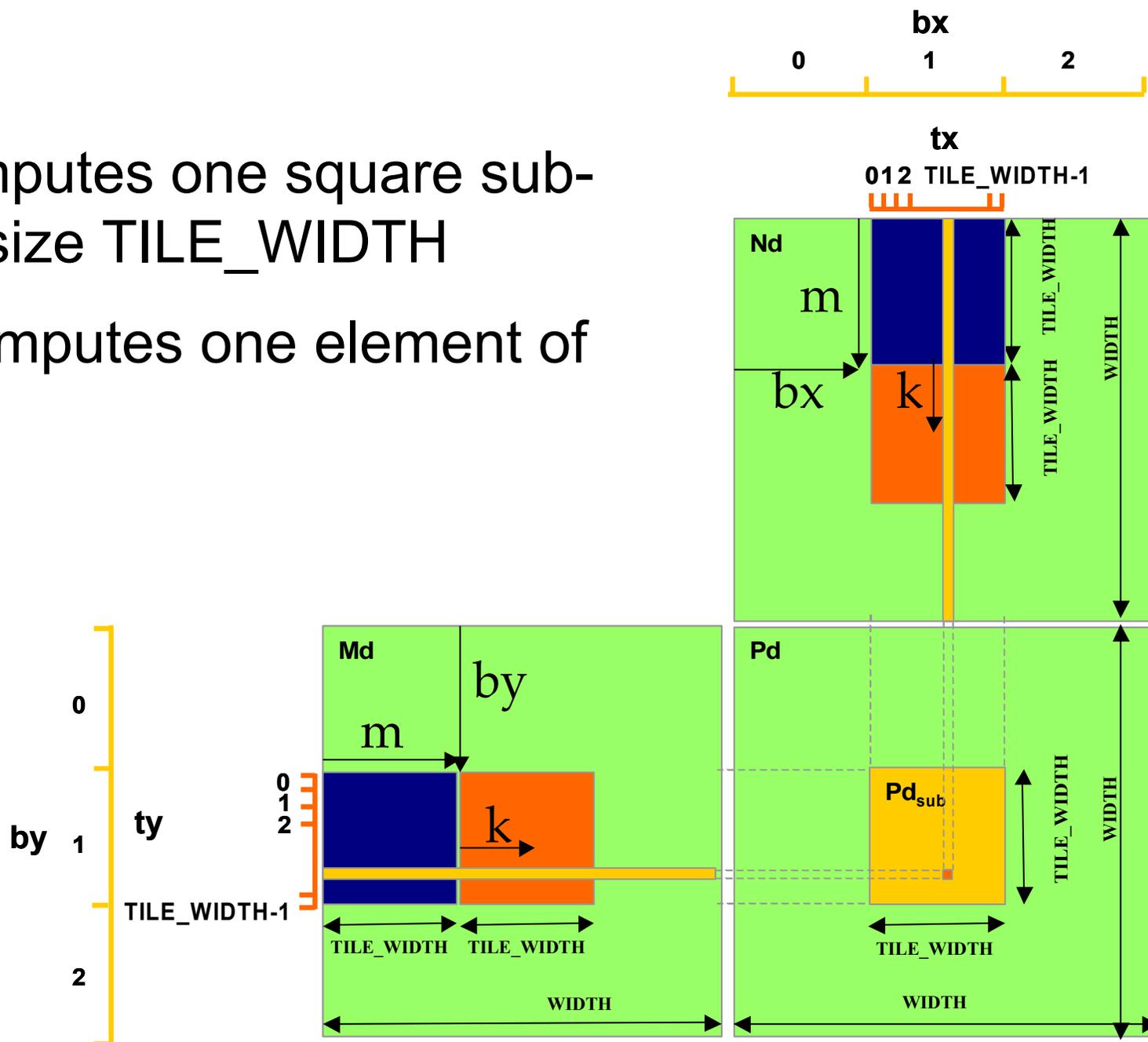
- Each block computes one square sub-matrix $Pd_{sub}$ of size TILE_WIDTH

- Each thread computes one element of $Pd_{sub}$

# G80 Shared Memory and Threading

- Each SM in G80 has 16KB shared memory
    - SM size is implementation dependent!
    - For TILE_WIDTH = 16, each thread block uses 2*256*4B = 2KB of shared memory.
    - Can potentially have up to 8 Thread Blocks actively executing
        - This allows up to 8*512 = 4,096 pending loads. (2 per thread, 256 threads per block)
    - The next TILE_WIDTH 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
- Using 16x16 tiling, we reduce the accesses to the global memory by a factor of 16
    - The 86.4B/s bandwidth can now support (86.4/4)*16 = 347.6 GFLOPS!