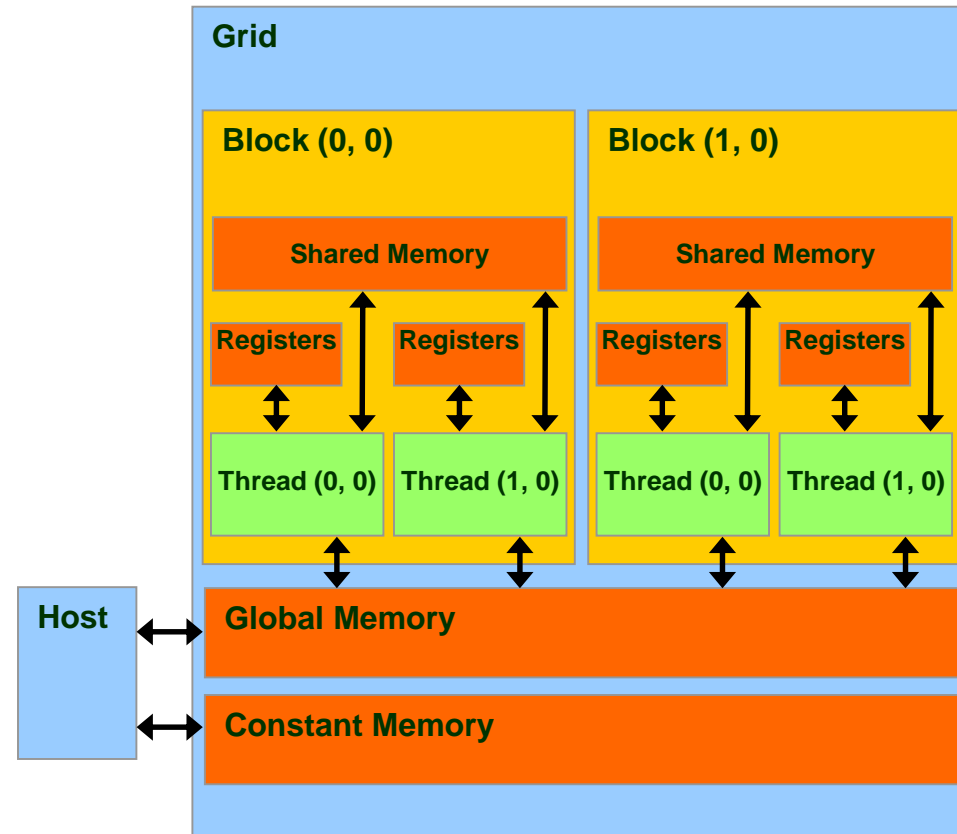


Two vertical lines, one blue and one orange, are positioned on the left side of the slide.

# Matrix Multiplication using Shared Memory

# 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**



# 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

# Review: Matrix Multiplication Kernel using Multiple Blocks

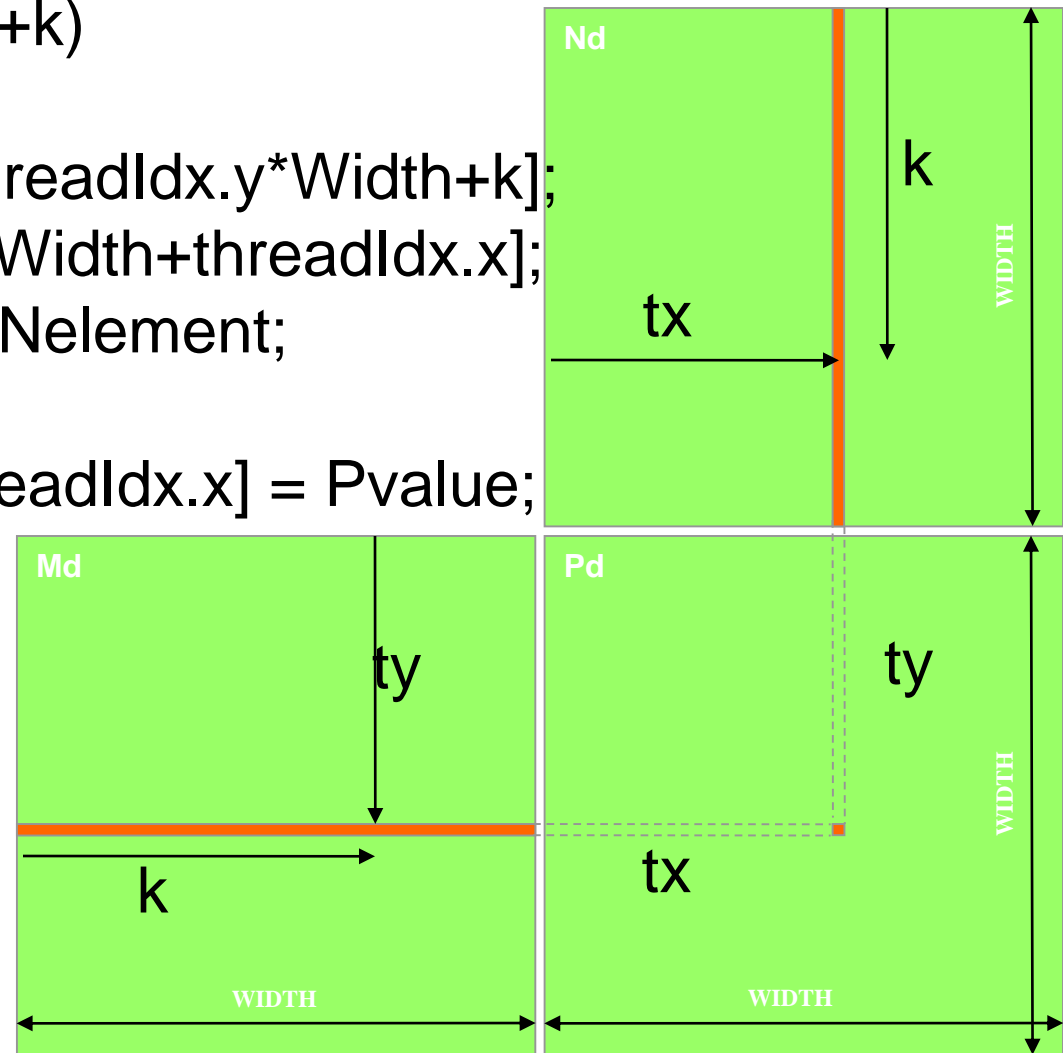
```
__global__ void MatrixMulKernel(float* Md, float* Nd, float* Pd, int Width)
{
    // Calculate the row index of the Pd element and M
    int Row = blockIdx.y*TILE_WIDTH + threadIdx.y;
    // Calculate the column index of Pd and N
    int Col = blockIdx.x*TILE_WIDTH + threadIdx.x;

    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;
}
```

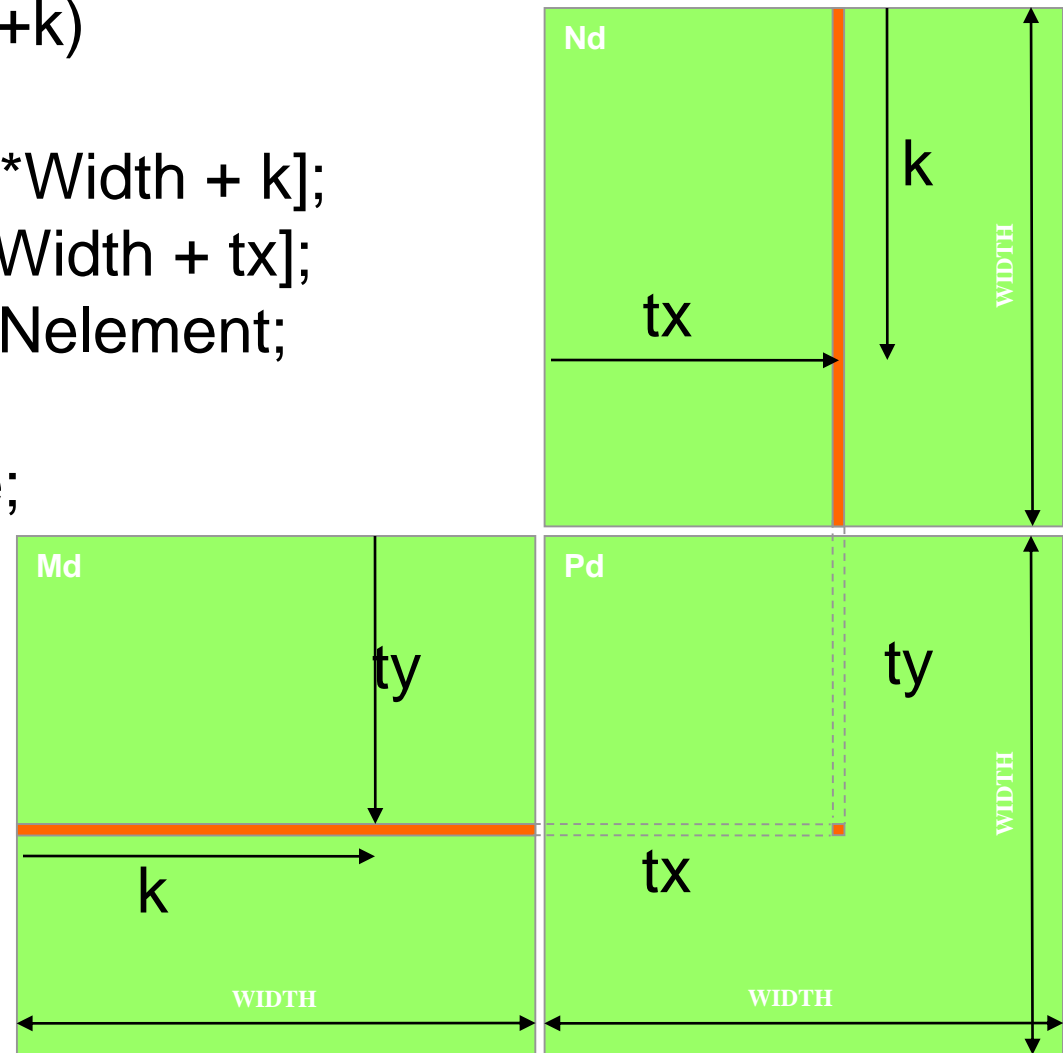
# Recall - Step 4: Kernel Function

```
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;
}
```



# Recall - Step 4: Kernel Function

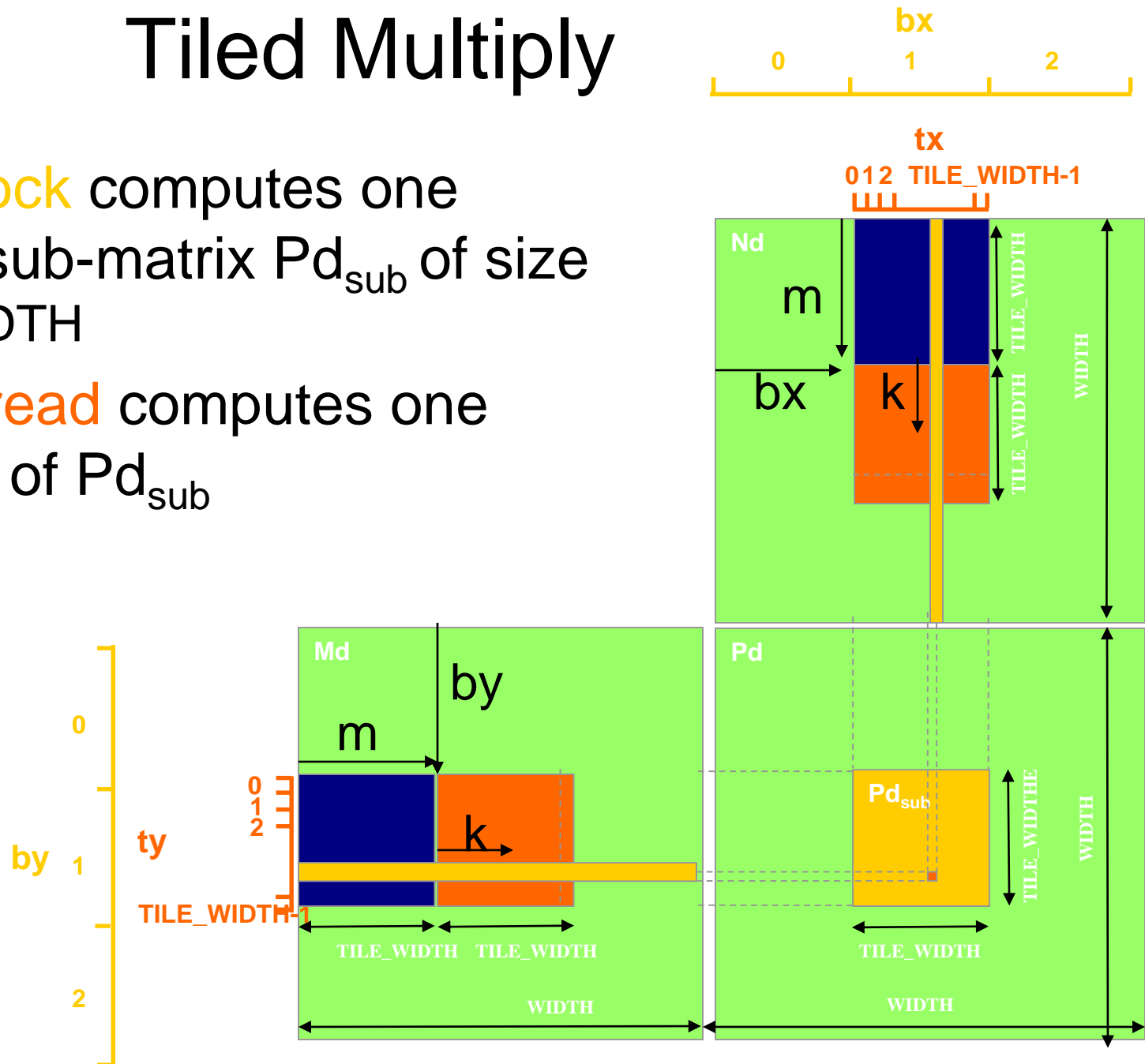
```
for (int k = 0; k < Width; ++k)
{
    float Melement = Md[ty*Width + k];
    float Nelement = Nd[k*Width + tx];
    Pvalue += Melement * Nelement;
}
Pd[ty*Width + tx] = Pvalue;
```





# Tiled Multiply

- Each **block** computes one square sub-matrix  $Pd_{\text{sub}}$  of size  $TILE\_WIDTH$
- Each **thread** computes one element of  $Pd_{\text{sub}}$



# 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 perform  $2 * 256 = 512$  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

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

// Identify the row and column of the Pd element to work on
5.  int Row = by * TILE_WIDTH + ty;
6.  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 loading of Md and Nd tiles into shared memory
9.      Mds[ty][tx] = Md[Row*Width + (m*TILE_WIDTH + tx)];
10.     Nds[ty][tx] = Nd[Col + (m*TILE_WIDTH + ty)*Width];
11.     __syncthreads();

11.     for (int k = 0; k < TILE_WIDTH; ++k)
12.         Pvalue += Mds[ty][k] * Nds[k][tx];
13.     __syncthreads();
14. }
13. Pd[Row*Width+Col] = Pvalue;
}
```