# CUDA-C
## (Compute Unified Device Architecture-C): Un linguaggio di programmazione per schede GPGPU
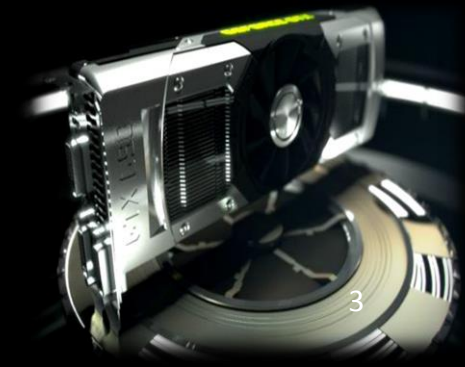
William Spataro

spataro@unical.it

# History

- **2001/2002 – researchers see GPU as data-parallel coprocessor**
  - The *GPGPU* field is born
- **2007 – NVIDIA releases CUDA 1.0**
  - *CUDA* **– Compute Uniform Device Architecture**
  - GPGPU shifts to *GPU Computing*
- **2008 – Khronos releases *OpenCL* specification**
- *2014 – Cuda 6.0*

# History

- **Nvidia creates CUDA to facilitate the development of parallel programs on GPUs (2007)**
- **The CUDA language is ANSI C extended with very few keywords for labeling data-parallel functions (kernels) and their associated data**
- **Nvidia technology benefits from massive economies of scale in the gaming market, CUDA-enabled cards are very inexpensive for the performance they provide**
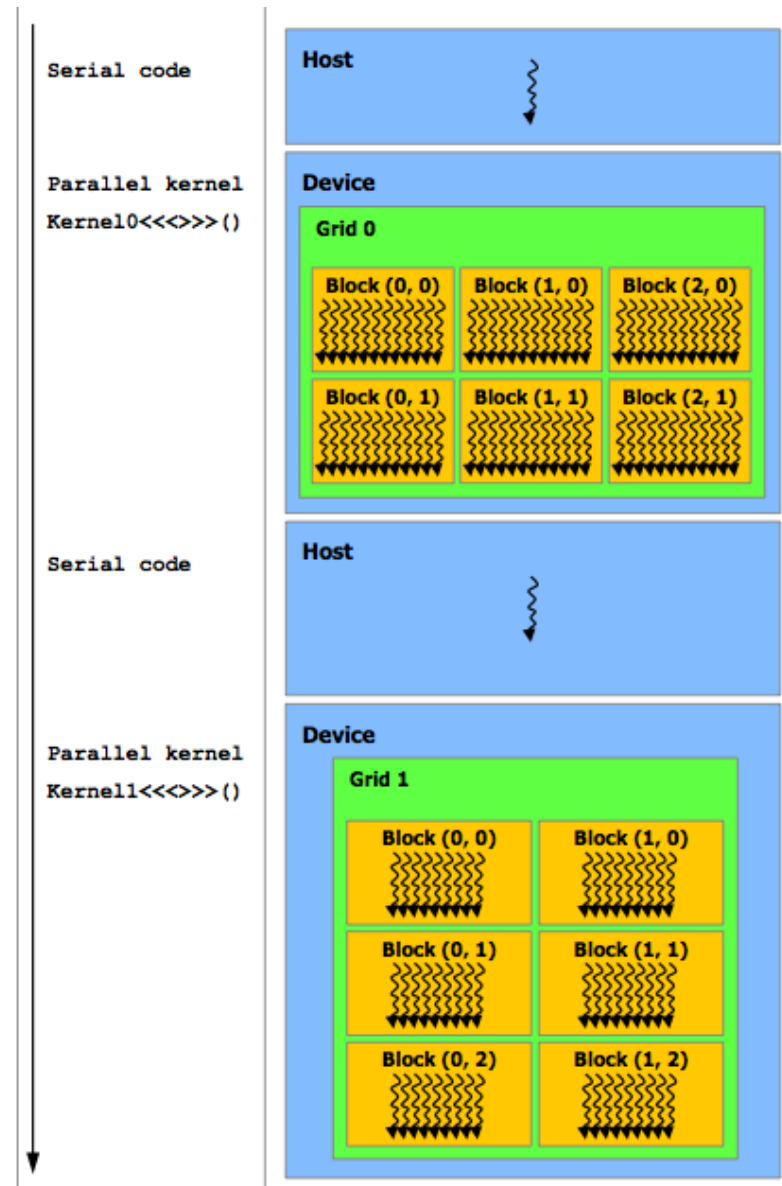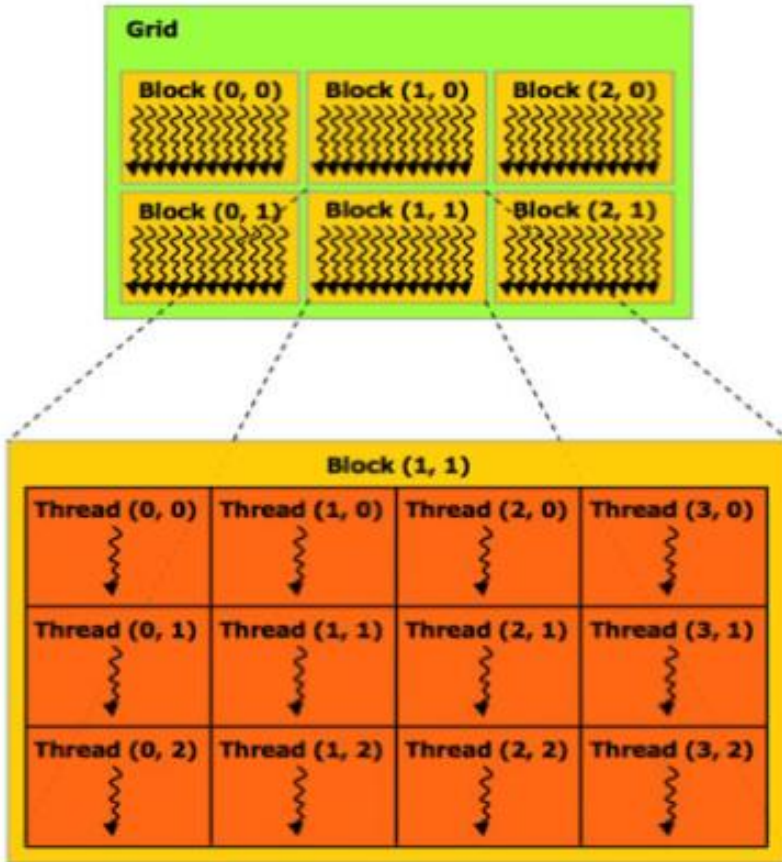
# CUDA

- **Scalable parallel programming model**
- **Minimal extensions to familiar C/C++ environment**
- **Heterogeneous serial-parallel computing**

# Modello di esecuzione

- Un codice CUDA alterna porzioni di codice seriale, eseguito dalla CPU, e di codice parallelo, eseguito dalla GPU.

- Il codice parallelo viene lanciato, ad opera della CPU, sulla GPU come **kernel.**
  - ➤ La GPU esegue un solo kernel alla volta.

- Un kernel è organizzato in **grids di blocks.**
  - ➤ Ogni block contiene lo stesso numero di threads.

- Ogni block viene eseguito da un solo multiprocessore: non può essere spezzato su più SM, mentre più blocks possono risiedere ed essere eseguiti in parallelo dallo stesso multiprocessore.

# Gerarchia dei thread



**Host** = CPU
**Device** = GPU

**Thread**: codice concorrente, eseguibile in parallelo ad altri threads su un device CUDA.

**Warp**: un gruppo di threads che possono essere eseguiti **fisicamente** in parallelo.
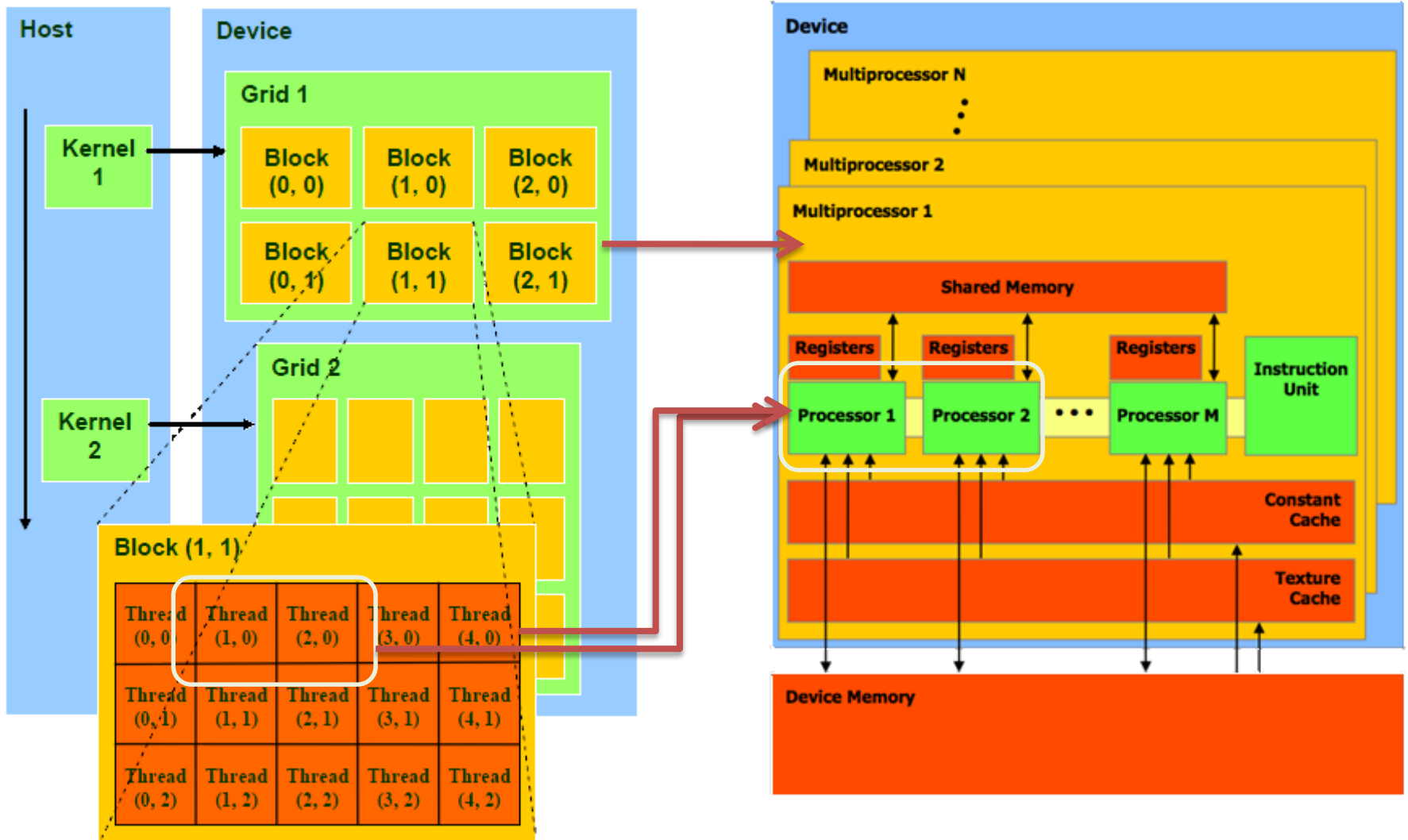
**Half-warp**: una delle 2 metà di un warp spesso eseguiti sullo stesso multiprocessore.

**Block**: un insieme di threads eseguiti sullo stesso Multiprocessore, e che quindi possono condividere memoria (stessa shared memory).

**Grid**: un insieme di thread blocks che eseguono un singolo kernel CUDA, in parallelismo logico, su una singola GPU.

**Kernel**: il codice CUDA che viene lanciato dalla CPU su una o più GPU.

# Esecuzione del codice



In ogni "multiprocessore" Threads è caricato un più processore threads per volta
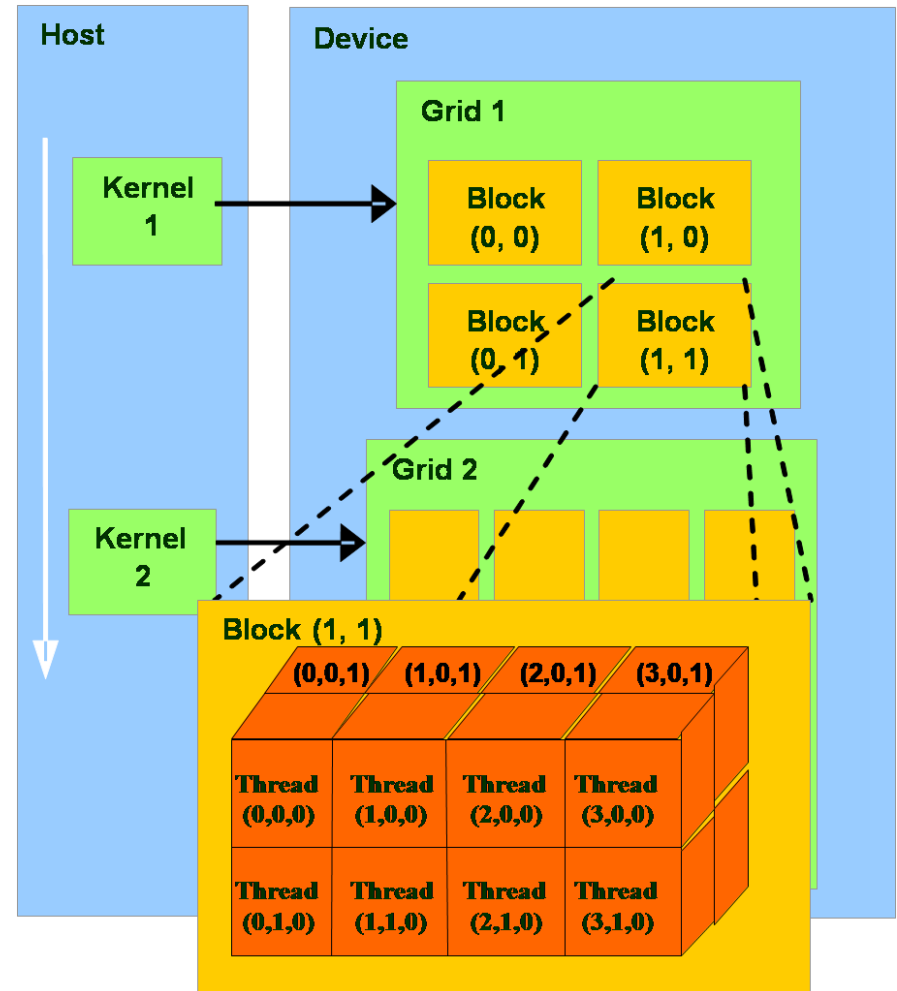
# Multidimensionalità degli IDs

Il codice parallelo viene lanciato, dalla CPU, sulla GPU , questa esegue un solo kernel alla volta.

La dimensione della griglia si misura in blocchi questi possono essere:

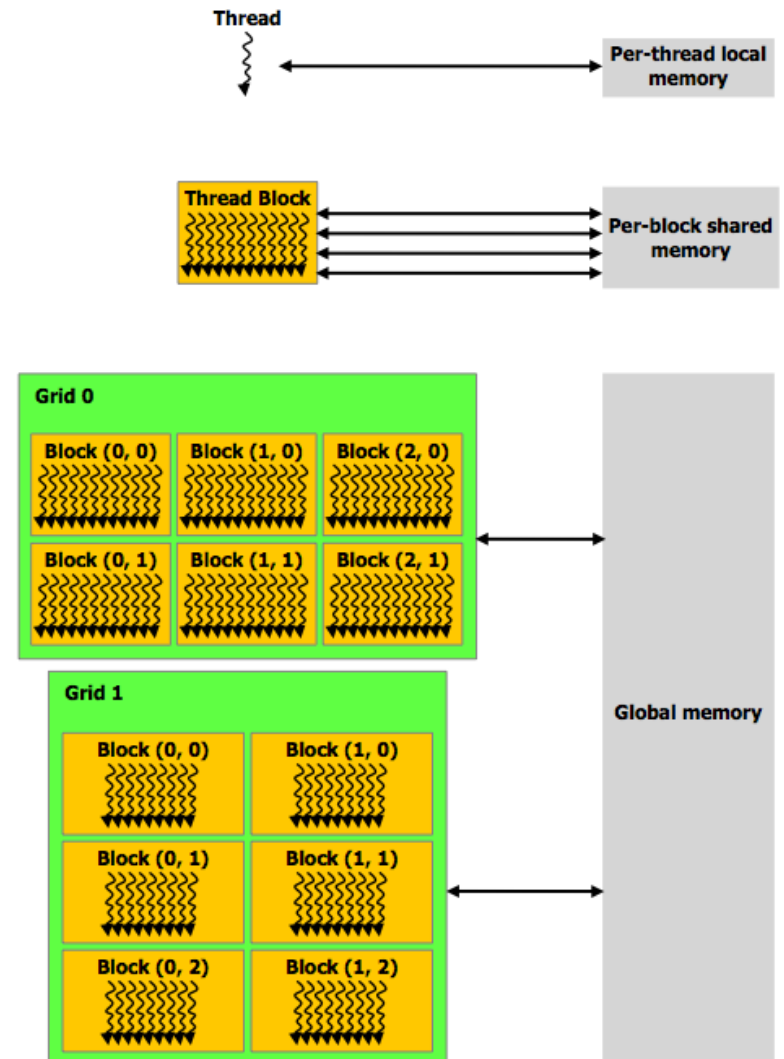Block: 1-D o 2-D (3D da comp. capability 2.0 in poi)

La dimensione dei blocchi si misura in thread
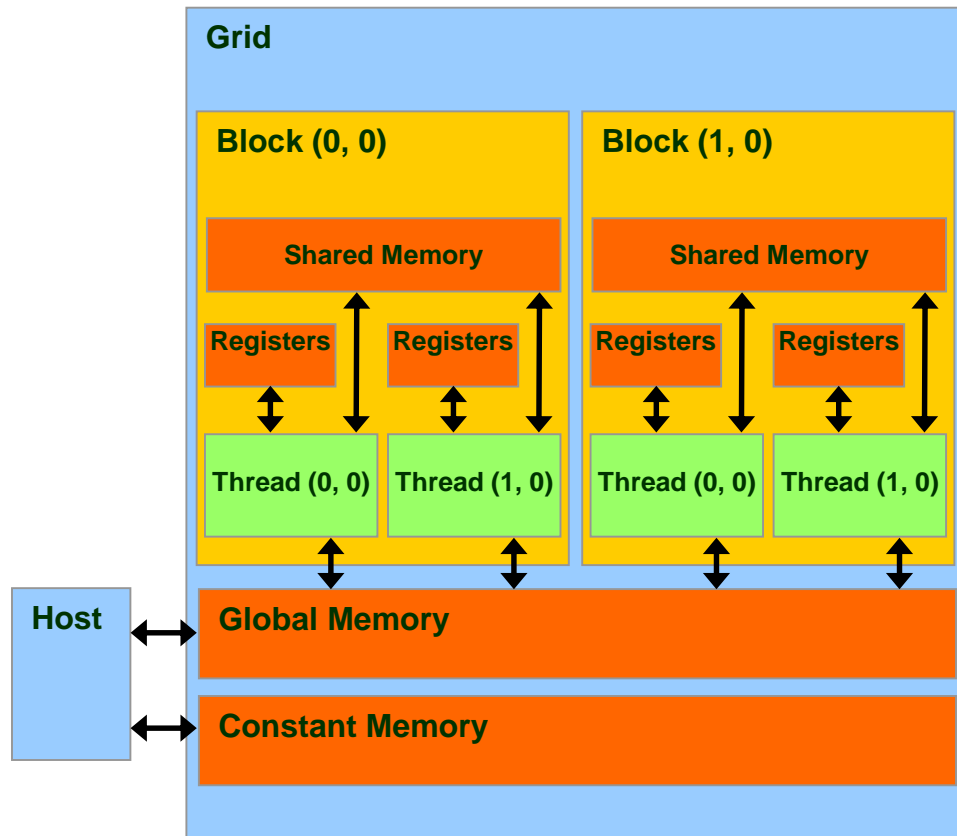
Thread 1-D,2-D,3-D

# Organizzazione gerarchica della memoria

- **Register file**: area di memoria privata di ciascun thread (var. locali).

- **Shared memory**: accessibile a tutti i threads dello stesso block. Può essere usata sia come spazio privato che come spazio condiviso.

- Tutti i threads accedono alla medesima **global memory** (off-chip DRAM).

- Memorie read-only accessibili da tutti i threads: **constant** e **texture memory**.
  ➢dotate di cache locale in ogni SM

- *Global, constant e texture memory sono memorie persistenti tra differenti lanci di kernel della stessa applicazione.*

- *Global memory bandwidth: 2 ordini di grandezza superiori della 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**

# Extended C

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

- **Keywords**
  - **threadIdx, blockIdx**
- **Intrinsics**
  - **__syncthreads**

- **Runtime API**
  - **Memory, symbol, execution management**

- **Function launch**

```
__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);
```

# 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

- **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

# Device 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`

# 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`
- `__device__` and `__host__` can be used together

# 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 from CUDA 1.0 on, explicit synch needed for blocking

# Device Runtime Component: Synchronization Function

- **`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

# Confronto codice seriale e parallelo

**Programma CPU**

```
void add_matrix
( float* a, float* b, float* c, int N ) {
int index;
for ( int i = 0; i < N; ++i )
  for ( int j = 0; j < N; ++j ) {
    index = i + j*N;
    c[index] = a[index] + b[index];
  }

}
int main() {
 add_matrix( a, b, c, N );
}
```
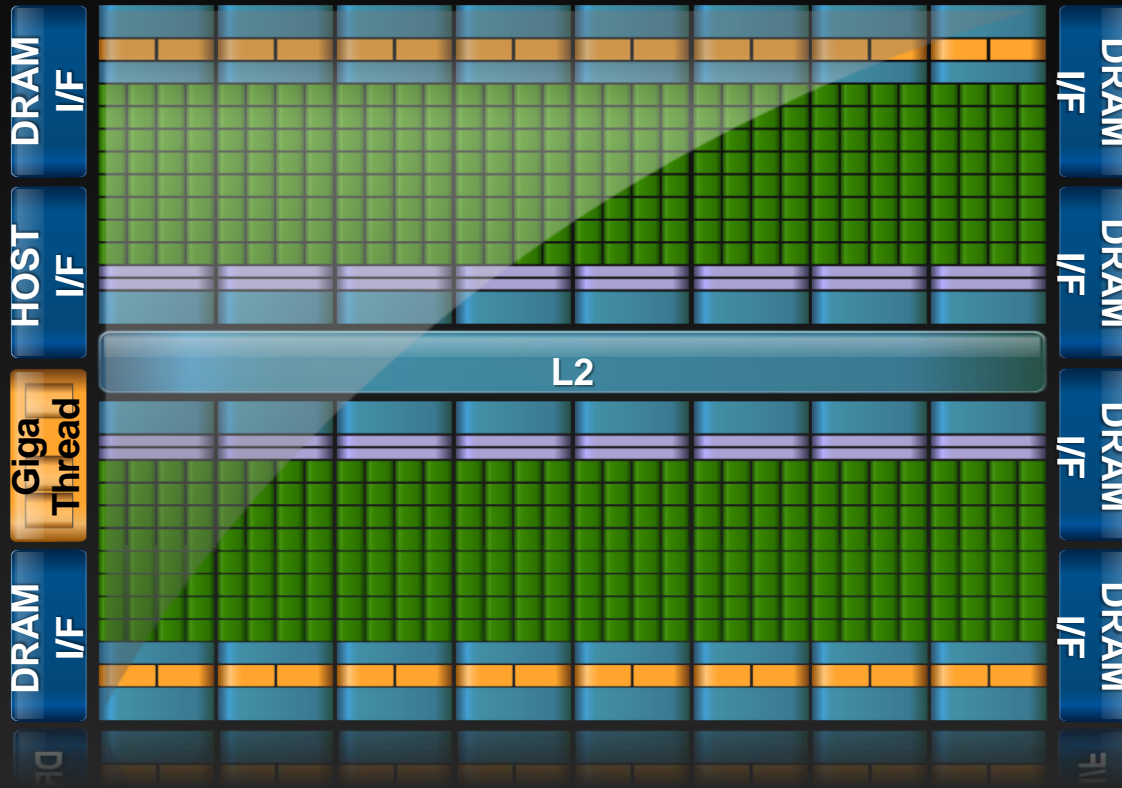
**Programma CUDA**

```
__global__ add_matrix
( float* a, float* b, float* c, int N ) {
int i = blockIdx.x * blockDim.x + threadIdx.x;
int j = blockIdx.y * blockDim.y + threadIdx.y;
int index = i + j*N;
if ( i < N && j < N )
  c[index] = a[index] + b[index];
}

int main() {
dim3 dimBlock( blocksize, blocksize );
dim3 dimGrid( N/dimBlock.x, N/dimBlock.y );
 add_matrix<<<dimGrid, dimBlock>>>( a, b, c, N );
}
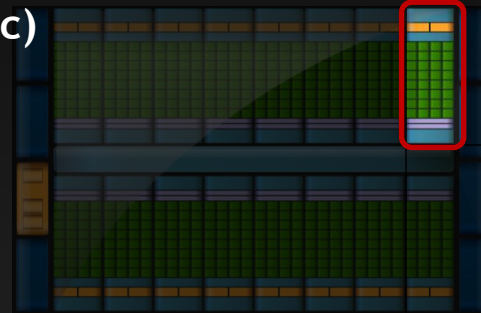```

Il ciclo for è sostituito da una griglia implicita

# NVIDIA GPU Architecture

## Fermi GF100

# SM Multiprocessor

- **32 CUDA Cores per SM (512 total)**

- **Direct load/store to memory**
  - **Usual linear sequence of bytes**
  - **High bandwidth (Hundreds GB/sec)**

- **64KB of fast, on-chip RAM**
  - **Software or hardware-managed**
  - **Shared amongst CUDA cores**
  - **Enables thread communication**

  **SIMT (**Single Instruction Multiple Thread**) execution**



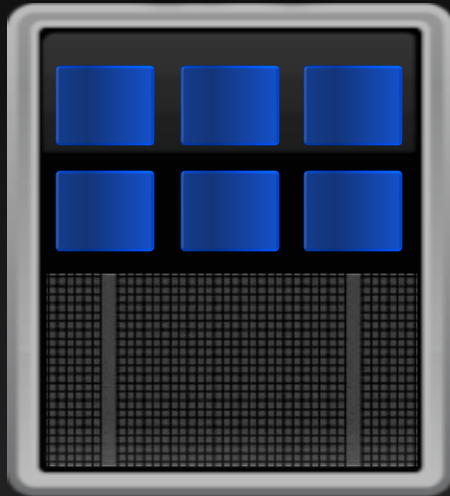| Instruction Cache | |
|---|---|
| Scheduler | Scheduler |
| Dispatch | Dispatch |
| Register File | |
| Core Core Core Core | |
| Core Core Core Core | |
| Core Core Core Core | |
| Core Core Core Core | |
| Core Core Core Core | |
| Core Core Core Core | |
| Core Core Core Core | |
| Core Core Core Core | |
| Load/Store Units x 16 | |
| Special Func Units x 4 | |
| Interconnect Network | |
| 64K Configurable Cache/Shared Mem | |
| Uniform Cache | |

# Compute Capability

- The **compute capability** of a device describes its architecture, e.g.
  - Number of SMs and registers
  - Sizes of memories
  - Features & capabilities

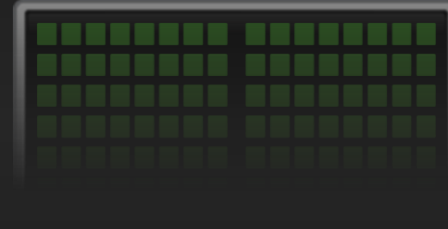| Compute Capability | Selected Features (see CUDA C Programming Guide for complete list) | Tesla models |
|---|---|---|
| 1.0 | Fundamental CUDA support | 870 |
| 1.3 | Double precision, improved memory accesses, atomics | 10-series |
| 2.0 | Caches, 3D grids, surfaces, ECC, P2P, concurrent kernels/copies, function pointers, recursion | 20-series |
| 3.0 | Dynamic parallelism | K-series |

- ## We will concentrate on Fermi devices
  - **Compute Capability >= 2.0**

# Heterogeneous computing

**CPU** + **GPU**

# Heterogeneous computing

**Device**

**Application Code**

**Host**

**GPU**

Compute-Intensive Functions

Use GPU to Parallelize

Rest of Sequential CPU Code

**CPU**

+

# Simple Processing Flow



1. Copy input data from CPU memory to GPU memory

# Simple Processing Flow



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance

# Simple Processing Flow



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance
3. Copy results from GPU memory to CPU memory

# Hello World!

```c
int main(void) {
        printf("Hello World!\n");
        return 0;

}
```

Output:

- **Standard C that runs on the host**

- **NVIDIA compiler (nvcc) can be used to compile programs with no *device* code**

```
$ nvcc
hello_world.
cu
$ a.out
Hello World!
$
```

# Hello World! with Device Code

```c
__global__ void mykernel(void) {
}

int main(void) {
    mykernel<<<1,1>>>();
    printf("Hello World!\n");
    return 0;
}
```

- Two new syntactic elements…

# Hello World! with Device Code

```
__global__ void mykernel(void) {
}
```

- **CUDA C/C++ keyword `__global__` indicates a function that:**
  - Runs on the device
  - Is called from host code

- **`nvcc` separates source code into host and device components**
  - Device functions (e.g. `mykernel()`) processed by NVIDIA compiler
  - Host functions (e.g. `main()`) processed by standard host compiler
    - `gcc, cl.exe`

# Hello World! with Device Code

```
mykernel<<<1,1>>>();
```

- **Triple angle brackets mark a call from *host* code to *device* code**
  - **Also called a "kernel launch"**

- **That's all that is required to execute a function on the GPU!**

# Hello World! with Device Code

```
__global__ void mykernel(void) {
}

int main(void) {
    mykernel<<<1,1>>>();
    printf("Hello World!\n");
    return 0;
}
```

Output:

```
$ nvcc
hello.cu
$ a.out
Hello World!
$
```

- **mykernel()** does nothing....

# Parallel Programming in CUDA C/C++

- GPU computing is about massive parallelism

- We need a more interesting example...

- We'll start by adding two integers and build up to vector addition

a       b       c

# Addition on the Device

- A simple kernel to add two integers

```
__global__ void add(int *a, int *b, int *c) {
    *c = *a + *b;
}
```

- As before __global__ is a CUDA C/C++ keyword meaning
  - `add()` will execute on the device
  - `add()` will be called from the host

# Addition on the Device

- Note that we use pointers for the variables

```
__global__ void add(int *a, int *b, int *c) {
        *c = *a + *b;
    }
```

- `add()` runs on the device, so `a`, `b` and `c` must point to device memory

- We need to allocate memory on the GPU

# Memory Management

- **Host** and **device** memory are separate entities
  - *Device* pointers point to GPU memory
    - May be passed to/from host code
    - May *not* be dereferenced in host code
  - *Host* pointers point to CPU memory
    - May be passed to/from device code
    - May *not* be dereferenced in device code

- Simple CUDA API for handling device memory
  - `cudaMalloc(), cudaFree(), cudaMemcpy()`
  - Similar to the C equivalents `malloc(), free(), memcpy()`

# Addition on the Device: `add()`

- **Returning to our `add()` kernel**

```
__global__ void add(int *a, int *b, int *c) {
        *c = *a + *b;
}
```

- **Let's take a look at main()...**

# Addition on the Device: `main()`

```c
int main(void) {
    int a, b, c;              // host copies of a, b, c
    int *d_a, *d_b, *d_c;     // device copies of a, b, c
    int size = sizeof(int);

    // Allocate space for device copies of a, b, c
    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    // Setup input values
    a = 2;
    b = 7;
```

# Addition on the Device: `main()`

```c
    // Copy inputs to device
    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    // Launch add() kernel on GPU
    add<<<1,1>>>(d_a, d_b, d_c);

    // Copy result back to host
    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    // Cleanup
    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

# Moving to Parallel

- **GPU computing is about massive parallelism**
    - So how do we run code in parallel on the device?

```
add<<< 1, 1 >>>();

add<<< N, 1 >>>();
```

- **Instead of executing `add()` once, execute N times in parallel**

# CUDA logical architecture

- A kernel is launched as a **grid** of **blocks** of **threads**
  - `blockIdx` and `threadIdx` can represent up to 3 dimensions

- Built-in variables:
  - `threadIdx`
  - `blockIdx`
  - `blockDim`
  - `gridDim`

# Vector Addition on the Device

- With `add()` running in parallel we can do vector addition

- Terminology: each parallel invocation of `add()` is referred to as a **block**
  - The set of blocks is referred to as a **grid**
  - Each invocation can refer to its block index using `blockIdx.x`

```
__global__ void add(int *a, int *b, int *c) {
        c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
}
```

- By using `blockIdx.x` to index into the array, each block handles a different index

# Vector Addition on the Device

```
__global__ void add(int *a, int *b, int *c) {
        c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
}
```

- On the device, each block can execute in parallel:

| Block 0 | Block 1 | Block 2 | Block 3 |
|---------|---------|---------|---------|
| `c[0] = a[0] + b[0];` | `c[1] = a[1] + b[1];` | `c[2] = a[2] + b[2];` | `c[3] = a[3] + b[3];` |

# Vector Addition on the Device: `add()`

- Returning to our parallelized `add()` kernel

```
__global__ void add(int *a, int *b, int *c) {
        c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
    }
```

- Let's take a look at main()...

# Vector Addition on the Device: `main()`

```c
#define N 512
int main(void) {
    int *a, *b, *c;              // host copies of a, b, c
    int *d_a, *d_b, *d_c;        // device copies of a, b, c
    int size = N * sizeof(int);

    // Alloc space for device copies of a, b, c
    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    // Alloc space for host copies of a, b, c and setup input values
    a = (int *)malloc(size); random_ints(a, N);
    b = (int *)malloc(size); random_ints(b, N);
    c = (int *)malloc(size);
```

# Vector Addition on the Device: `main()`

```
    // Copy inputs to device
    cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, b, size, cudaMemcpyHostToDevice);

    // Launch add() kernel on GPU with N blocks
    add<<<N,1>>>(d_a, d_b, d_c);

    // Copy result back to host
    cudaMemcpy(c, d_c, size, cudaMemcpyDeviceToHost);

    // Cleanup
    free(a); free(b); free(c);
    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

# CUDA logical architecture

- **A kernel is launched as a grid of blocks of threads**
  - `blockIdx` and `threadIdx` can represent up to 3 dimensions

- Built-in variables:
  - `threadIdx`
  - `blockIdx`
  - `blockDim`
  - `gridDim`

# Threads Hierarchy

- ● A block can be split into parallel threads
    - all threads execute the same sequential program

Thread *t*

- ● Thread block is a group of threads that can:
    - ● Synchronize their execution
    - ● Communicate via shared memory

Block *b*

t0 t1 … tB

# CUDA Threads

- Terminology: a block can be split into parallel **threads**

- Let's change `add()` to use parallel *threads* instead of parallel *blocks*

```
__global__ void add(int *a, int *b, int *c) {
    c[threadIdx.x] = a[threadIdx.x] + b[threadIdx.x];
}
```

- We use `threadIdx.x` instead of `blockIdx.x`

- Need to make one change in `main()`...

# Vector Addition Using Threads: `main()`

```c
#define N 512
int main(void) {
    int *a, *b, *c;                    // host copies of a, b, c
    int *d_a, *d_b, *d_c;              // device copies of a, b, c
    int size = N * sizeof(int);

    // Alloc space for device copies of a, b, c
    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    // Alloc space for host copies of a, b, c and setup input values
    a = (int *)malloc(size); random_ints(a, N);
    b = (int *)malloc(size); random_ints(b, N);
    c = (int *)malloc(size);
```

# Vector Addition Using Threads: `main()`

```
    // Copy inputs to device
    cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, b, size, cudaMemcpyHostToDevice);

    // Launch add() kernel on GPU with N threads
    add<<<1,N>>>(d_a, d_b, d_c);

    // Copy result back to host
    cudaMemcpy(c, d_c, size, cudaMemcpyDeviceToHost);

    // Cleanup
    free(a); free(b); free(c);
    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

# Combining Blocks _and_ Threads

- We've seen parallel vector addition using:
  - Many blocks with one thread each
  - One block with many threads

- Adapting vector addition to use both _blocks_ and _threads_

- Data indexing…

# Indexing Arrays with Blocks and Threads

- No longer as simple as using `blockIdx.x` and `threadIdx.x`
  - Consider indexing an array with one element per thread (8 threads/block)

| threadIdx.x | threadIdx.x | threadIdx.x | threadIdx.x |
|---|---|---|---|
| 0 1 2 3 4 5 6 7 | 0 1 2 3 4 5 6 7 | 0 1 2 3 4 5 6 7 | 0 1 2 3 4 5 6 7 |
| blockIdx.x = 0 | blockIdx.x = 1 | blockIdx.x = 2 | blockIdx.x = 3 |

- With M threads/block a unique index for each thread is given by:

```
int index = blockIdx.x * M + threadIdx.x;
```

# Indexing Arrays: Example

- Which thread will operate on the red element?

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |

M = 8          threadIdx.x = 5

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

blockIdx.x = 2

```
int index = threadIdx.x + blockIdx.x * M;
          =       5      +      2       * 8;
          = 21;
```

# Vector Addition with Blocks and Threads

- Use the built-in variable `blockDim.x` for threads per block

```
int index = threadIdx.x + blockIdx.x * blockDim.x;
```

- Combined version of `add()` to use parallel threads *and* parallel blocks

```
__global__ void add(int *a, int *b, int *c) {
    int index = threadIdx.x + blockIdx.x * blockDim.x;
    c[index] = a[index] + b[index];
}
```

- What changes need to be made in `main()`?

# Addition with Blocks and Threads: `main()`

```c
#define N (2048*2048)
#define THREADS_PER_BLOCK 512
int main(void) {
    int *a, *b, *c;              // host copies of a, b, c
    int *d_a, *d_b, *d_c;        // device copies of a, b, c
    int size = N * sizeof(int);

    // Alloc space for device copies of a, b, c
    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    // Alloc space for host copies of a, b, c and setup input values
    a = (int *)malloc(size); random_ints(a, N);
    b = (int *)malloc(size); random_ints(b, N);
    c = (int *)malloc(size);
```

# Addition with Blocks and Threads: `main()`

```c
    // Copy inputs to device
    cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, b, size, cudaMemcpyHostToDevice);

    // Launch add() kernel on GPU
    add<<<N/THREADS_PER_BLOCK,THREADS_PER_BLOCK>>>(d_a, d_b, d_c);

    // Copy result back to host
    cudaMemcpy(c, d_c, size, cudaMemcpyDeviceToHost);

    // Cleanup
    free(a); free(b); free(c);
    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

# Handling Arbitrary Vector Sizes

- Typical problems are not friendly multiples of `blockDim.x`

- Avoid accessing beyond the end of the arrays:

```
__global__ void add(int *a, int *b, int *c, int n) {
    int index = threadIdx.x + blockIdx.x * blockDim.x;
    if (index < n)
        c[index] = a[index] + b[index];
}
```
- Update the kernel launch:
```
add<<<(N + M-1) / M, M >>>(d_a, d_b, d_c, N);
```

# Kernel with 2D Indexing

```
__global__ void kernel( int *a, int dimx, int dimy )
{
    int ix  = blockIdx.x*blockDim.x + threadIdx.x;
    int iy  = blockIdx.y*blockDim.y + threadIdx.y;
    int idx = iy*dimx + ix;

    a[idx]  = a[idx]+1;
}
```

# Kernel with 2D Indexing

```
__global__ void kernel( int *a, int dimx, int dimy )
{
    int ix  = blockIdx.x*blockDim.x + threadIdx.x;
    int iy  = blockIdx.y*blockDim.y + threadIdx.y;
    int idx = iy*dimx + ix;

    a[idx]  = a[idx]+1;
}
```

```
int main()
{
    int dimx = 16;
    int dimy = 16;
    int num_bytes = dimx*dimy*sizeof(int);

    int *d_a=0, *h_a=0; // device and host pointers

    h_a = (int*)malloc(num_bytes);
    cudaMalloc( (void**)&d_a, num_bytes );

    dim3 grid, block;
    block.x = 4;
    block.y = 4;
    grid.x  = dimx / block.x;
    grid.y  = dimy / block.y;

    kernel<<<grid, block>>>( d_a, dimx, dimy );

    cudaMemcpy( h_a, d_a, num_bytes, cudaMemcpyDeviceToHost );

    ......
}
```

CONCEPTS

| Heterogeneous Computing |
| Blocks |
| Threads |
| Indexing |
| Shared memory |
| __syncthreads() |
| Asynchronous operation |
| Handling errors |
| Managing devices |

# MANAGING THE DEVICE

# Coordinating Host & Device

- **Kernel launches are asynchronous**
  - **Control returns to the CPU immediately**

- **CPU needs to synchronize before consuming the results**

| | |
|---|---|
| `cudaMemcpy()` | Blocks the CPU until the copy is complete<br>Copy begins when all preceding CUDA calls have completed |
| `cudaMemcpyAsync()` | Asynchronous, does not block the CPU |
| `cudaDeviceSynchronize()` | Blocks the CPU until all preceding CUDA calls have completed |

# Reporting Errors

- **All CUDA API calls return an error code (`cudaError_t`)**
  - **Error in the API call itself**
    - **OR**
  - **Error in an earlier asynchronous operation (e.g. kernel)**

- **Get the error code for the last error:**
  ```
  cudaError_t cudaGetLastError(void)
  ```
- **Get a string to describe the error:**
  ```
  char *cudaGetErrorString(cudaError_t)

  printf("%s\n", cudaGetErrorString(cudaGetLastError()));
  ```

# Device Management

- **Application can query and select GPUs**
  - `cudaGetDeviceCount(int *count)`
  - `cudaSetDevice(int device)`
  - `cudaGetDevice(int *device)`
  - `cudaGetDeviceProperties(cudaDeviceProp *prop, int device)`

- **Multiple threads can share a device**

- **A single thread can manage multiple devices**
  - `cudaSetDevice(i)` to select current device
  - `cudaMemcpy(…)` for peer-to-peer copies[†]

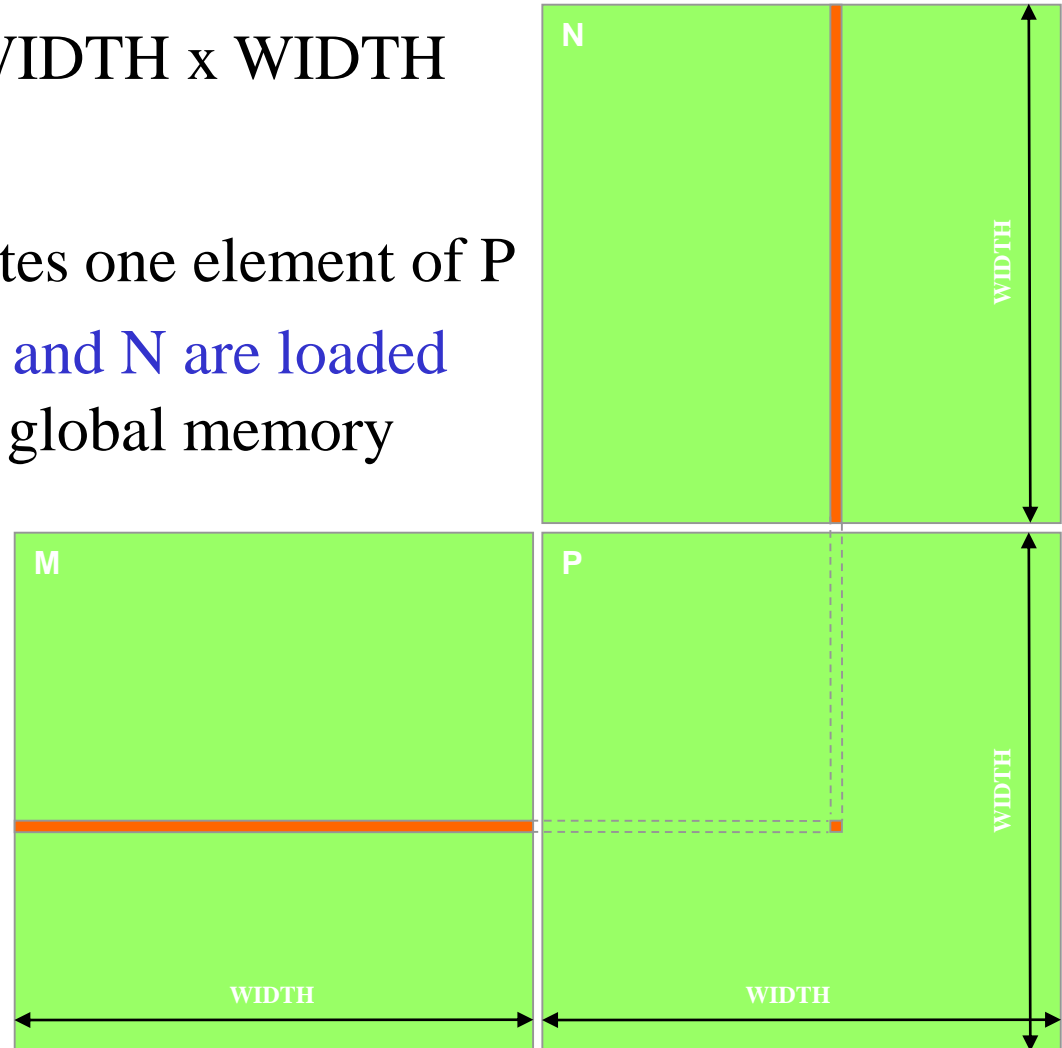[†] requires OS and device support

# Part II

# A Simple Running Example
# Matrix Multiplication

- Let's now see 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

# Programming Model:
# Square Matrix Multiplication Example

- P = M * N of size WIDTH x WIDTH

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

# Memory Layout of a Matrix in C

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

M

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

# Step 1: 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;
        }
}
```

# 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;
```
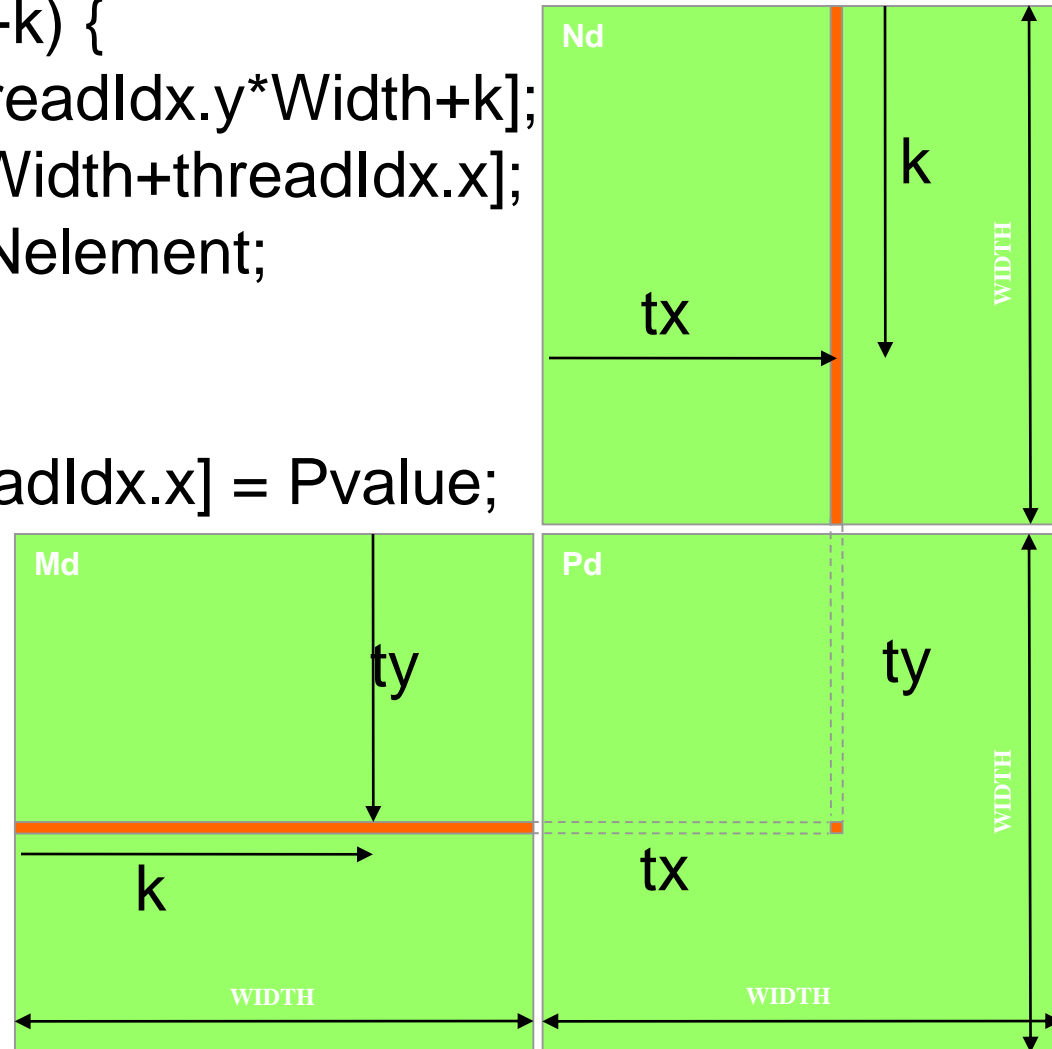
# Step 4: Kernel Function  (cont.)

```
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

tx

WIDTH

Md

ty
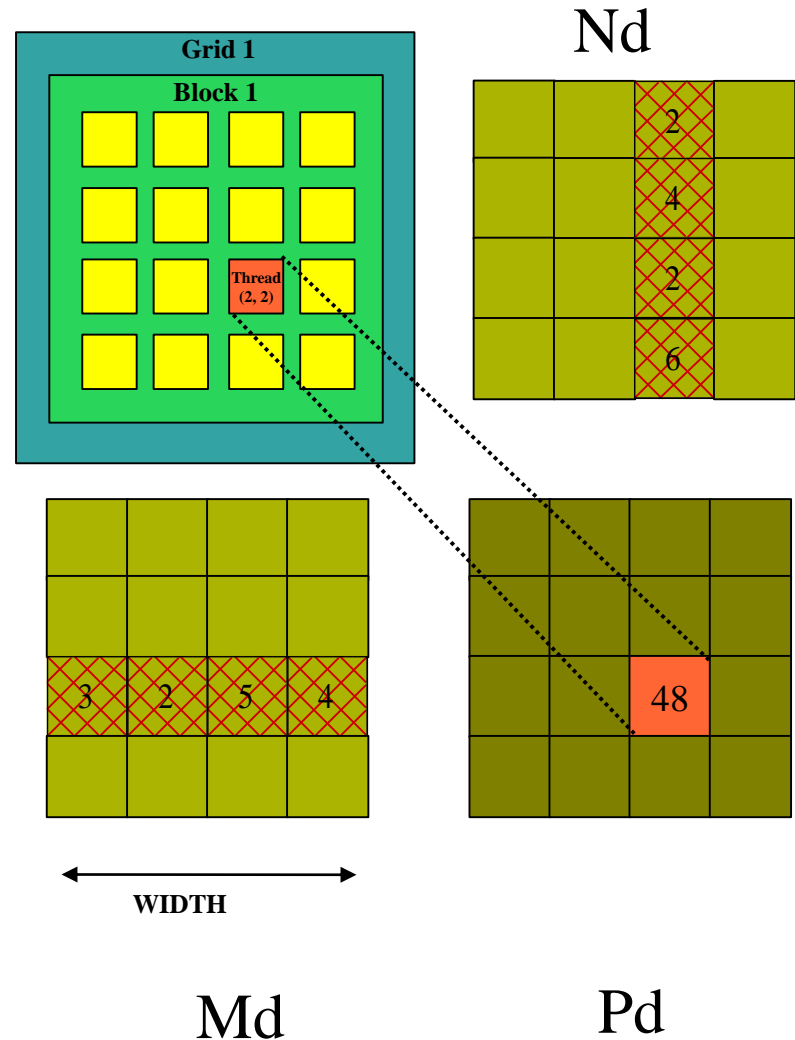
k

WIDTH

Pd

ty

tx

WIDTH

WIDTH

# Step 5: Kernel Invocation
# (Host-side Code)

```
// Setup the execution configuration
  dim3 dimGrid(1, 1);   // one block only!
   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
  - Perform one multiply and addition for each pair of Md and Nd elements
  - Compute to global memory access ratio* close to 1:1 (not very high)
- Size of matrix limited by the number of threads allowed in a thread block ! (e.g. 1024 only!)



Nd

Grid 1

Block 1

Thread (2, 2)

2
4
2
6

WIDTH

Md          Pd

3  2  5  4

48

*CGMA (Computation to Global Memory Access) index: around 20/30:1 to be REALLY good!

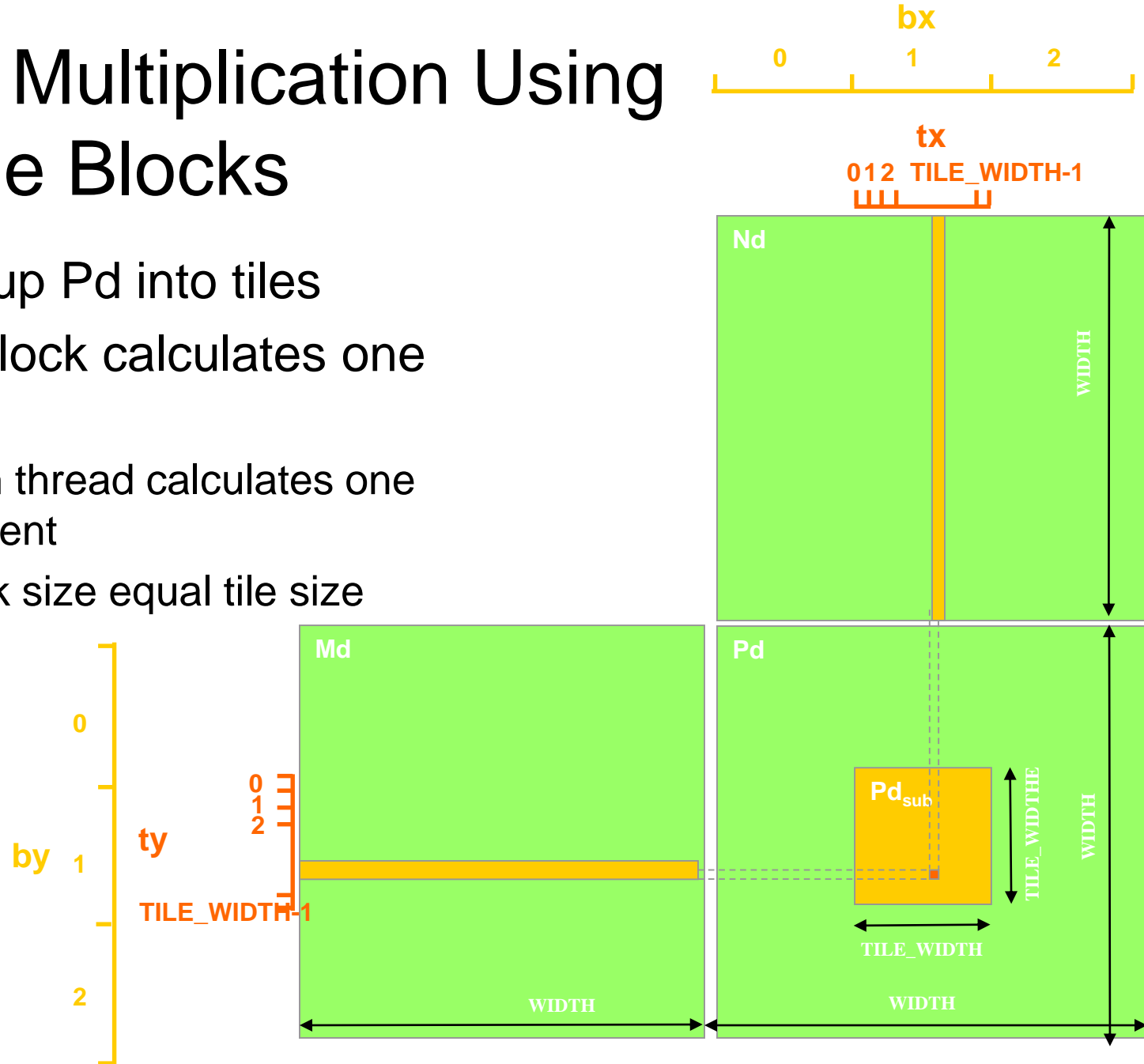# Step 7: Handling Arbitrary Sized Square Matrices

- 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

# 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

# Revised 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 idenx 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;
}
```

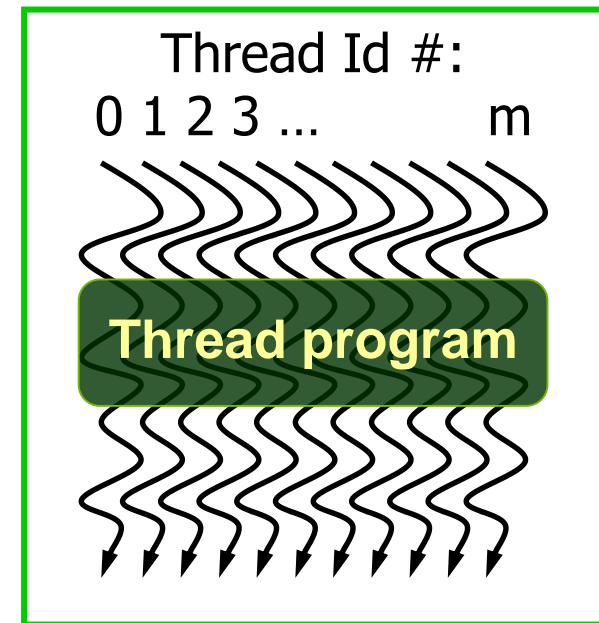# Revised Step 5: Kernel Invocation (Host-side Code)

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


// Launch the device computation threads!
MatrixMulKernel<<<dimGrid, dimBlock>>>(Md, Nd, Pd, Width);
```

# Review: CUDA Thread Block

- All threads in a block execute the same kernel program (SPMD)
- Programmer declares block:
  - Block size 1 to **1024** 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!

**CUDA Thread Block**

Thread Id #:
0 1 2 3 ...          m

**Thread program**

# Transparent Scalability

- Since threads in different blocks **cannot** perform barrier synchronization with each other, the runtime system is free to assigns blocks to any processor at any time, depending on hardware

  - A kernel scales across any number of parallel processors



High-end device

Device

Kernel grid

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

Device

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

Low-cost device

time

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

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

# G80 CUDA mode – A Review

- Processors execute computing threads
- New operating mode/HW interface for computing

# G80 Example: Executing Thread Blocks

**SM 0  SM 1**



**Blocks**

**Blocks**

- Threads are assigned to Streaming Multiprocessors in block granularity
  - Up to **8** (physical) blocks to each SM as resource allows
  - SM in G80 can take up to **768** (physical) 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

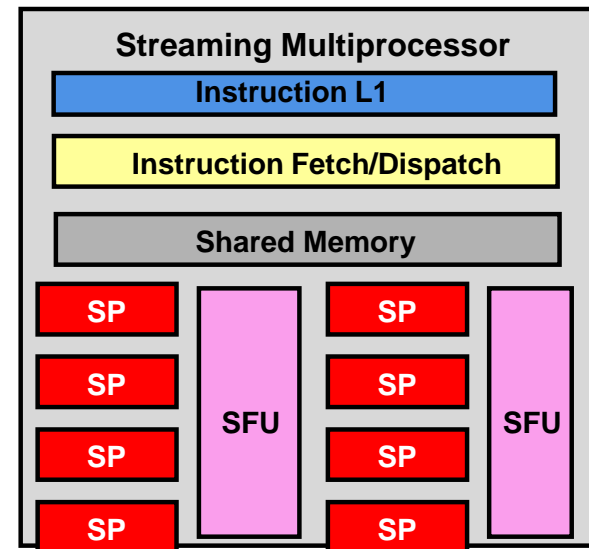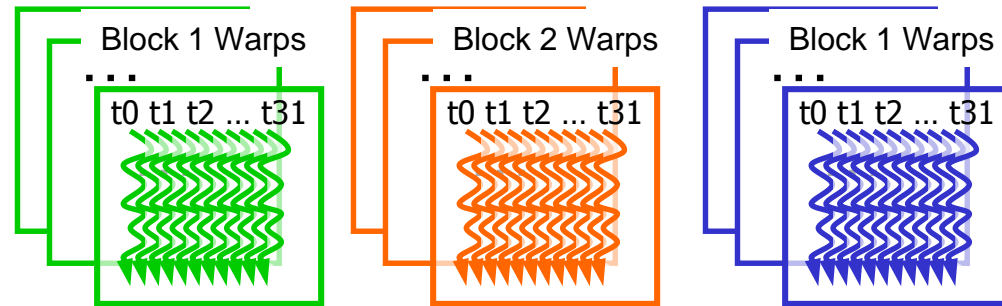# G80 Example: Warps and Thread Scheduling

- Each Block is executed as 32-thread Warps
  - An implementation decision, not part of the CUDA programming model
  - Warps are scheduling units in SM
  - **SIMD** !

- 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



Block 1 Warps
. . .
t0 t1 t2 ... t31

Block 2 Warps
. . .
t0 t1 t2 ... t31

Block 1 Warps
. . .
t0 t1 t2 ... t31

**Streaming Multiprocessor**

**Instruction L1**

**Instruction Fetch/Dispatch**

**Shared Memory**

SP  SP  SFU  SP  SP  SFU

SP  SP

SP  SP

# G80 Example: Thread Scheduling (Cont.)

- Each SM implements **<u>zero-overhead warp scheduling</u>**:

  – Warps whose next instruction has its operands ready for consumption are eligible for execution (i.e. Cuda runtime system maintains a list of warp blocks…)

  – (**Latency tollerance**) Eligible Warps are selected for execution on a prioritized scheduling policy (ex: **Post office queue**)

  – All threads in a warp execute the same instruction when selected (i.e. **<u>SIMD fashion</u>**)

# G80 Block Granularity Considerations (max 8 blocks – 768 threads, which ever comes first!)

- 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 8*64=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. <u>Not even one can fit into an SM!</u> (max 768!)

Best !

* Physical threads for G80 ! 1024 virtual threads per block from Compute capability 2.0 on…

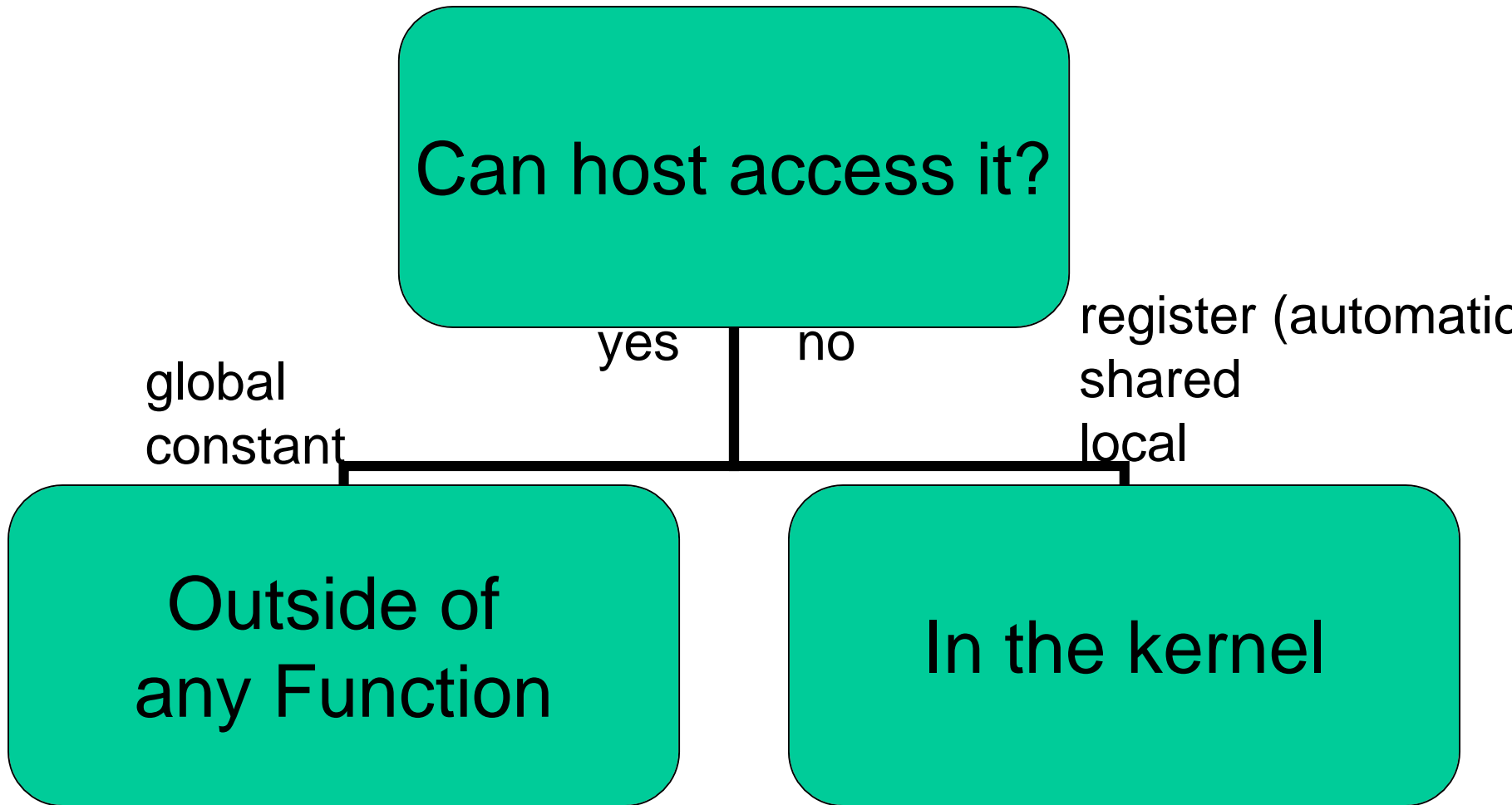# Hints…

1.Cuda occupancy calculator!

2.`cudaGetDeviceProperties()`

# 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

# Global memory access efficiency

- Although having many threads available for execution can theoretically **tolerate long memory access latency**, one can easily run into a situation where traffic congestion (ie. **too much global memory accesses** ) prevents all but few threads from making progress, thus rendering some SM idle!

- A common strategy for reducing global memory traffic (i.e. increasing the number of floating-point operations performed for each access to the global memory) is to **partition the data into subsets called tiles** such that each tile fits into the **shared memory** and the kernel computations on these tiles can be done independently of each other.

- In the simplest form, **the tile dimensions equal those of the block.**

# A Common Programming Strategy

- **Global memory** resides in device memory (DRAM) - much **slower** access than **shared memory** (up to 2 order magnitude!)

- 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 (i.e. **cooperation**)

    - 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)

# Matrix Multiplication using Shared Memory

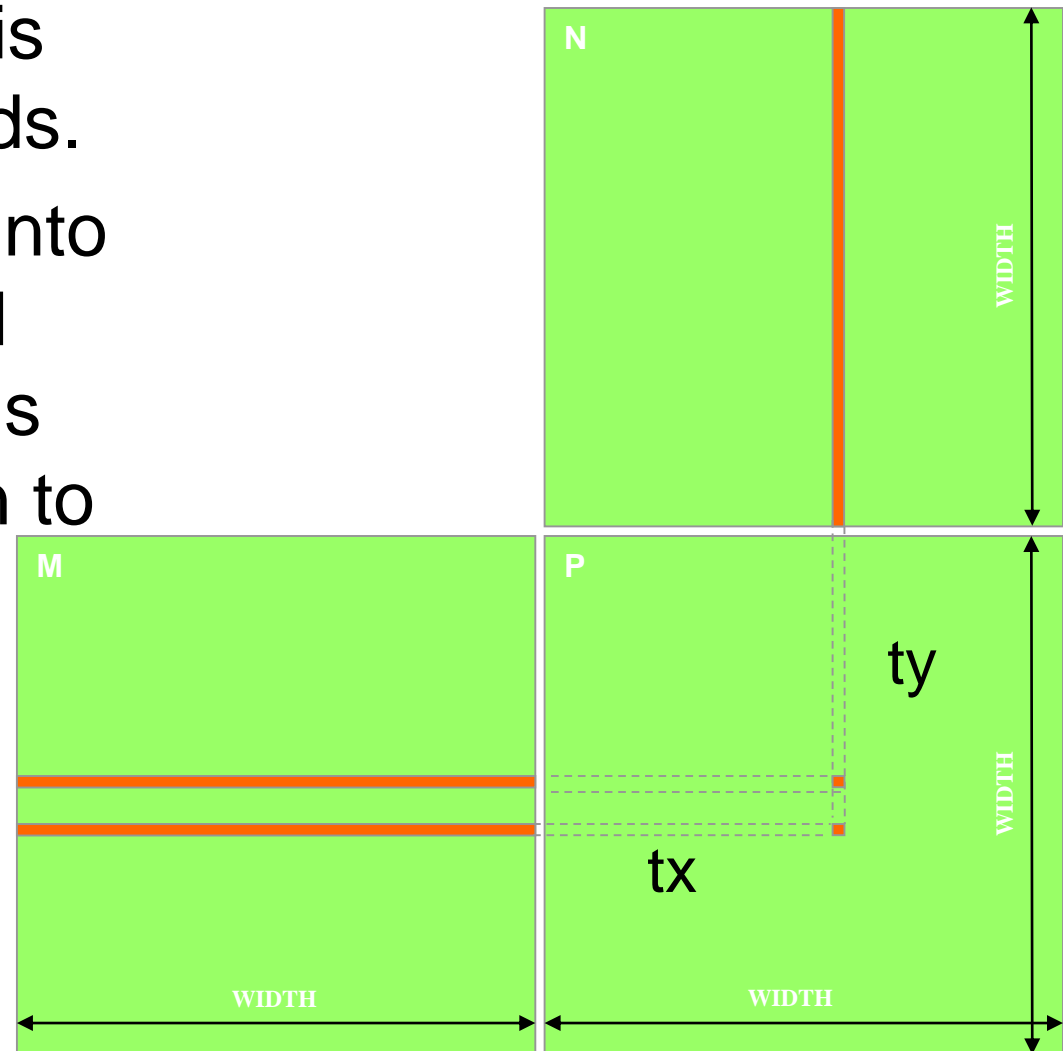# Review: Matrix Multiplication Kernel using Multiple Blocks

```
__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 idenx 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;
}
```

# 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
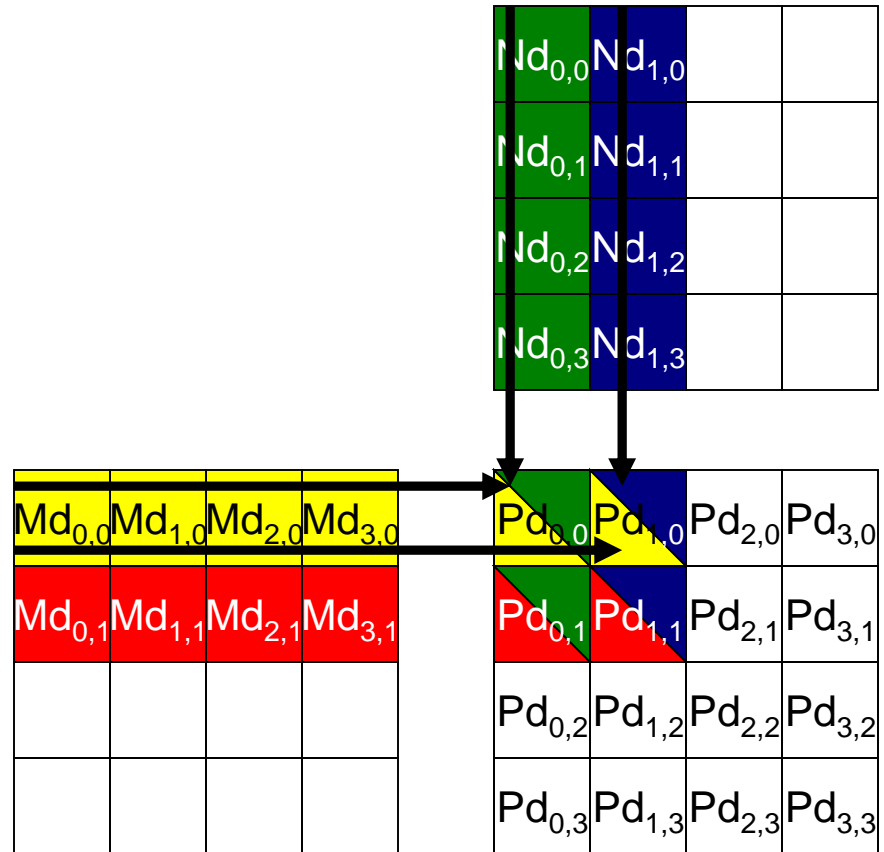  - Tiled algorithms

N

WIDTH

M

P

ty

tx

WIDTH

WIDTH

WIDTH

# 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

# A Small Example

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

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $Md_{0,0}$ | $Md_{1,0}$ | $Md_{2,0}$ | $Md_{3,0}$ | $Pd_{0,0}$ | $Pd_{1,0}$ | $Pd_{2,0}$ | $Pd_{3,0}$ |
| $Md_{0,1}$ | $Md_{1,1}$ | $Md_{2,1}$ | $Md_{3,1}$ | $Pd_{0,1}$ | $Pd_{1,1}$ | $Pd_{2,1}$ | $Pd_{3,1}$ |
| | | | | $Pd_{0,2}$ | $Pd_{1,2}$ | $Pd_{2,2}$ | $Pd_{3,2}$ |
| | | | | $Pd_{0,3}$ | $Pd_{1,3}$ | $Pd_{2,3}$ | $Pd_{3,3}$ |

# Every Md and Nd Element is used **<u>exactly twice</u>** 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

# Breaking Md and Nd into Tiles

- Break up the inner product loop of each thread into **phases**

- At the beginning of each phase, load the Md and Nd elements that **everyone needs during the phase** into shared memory

- Everyone access the Md and Nd elements from the shared memory during the phase

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

| | Phase 1 | | | Phase 2 | | |
|---|---|---|---|---|---|---|
| $T_{0,0}$ | $Md_{0,0}$ $\downarrow$ $Mds_{0,0}$ | $Nd_{0,0}$ $\downarrow$ $Nds_{0,0}$ | $PValue_{0,0}$ += $Mds_{0,0}*Nds_{0,0}$ + $Mds_{1,0}*Nds_{0,1}$ | $Md_{2,0}$ $\downarrow$ $Mds_{0,0}$ | $Nd_{0,2}$ $\downarrow$ $Nds_{0,0}$ | $PValue_{0,0}$ += $Mds_{0,0}*Nds_{0,0}$ + $Mds_{1,0}*Nds_{0,1}$ |
| $T_{1,0}$ | $Md_{1,0}$ $\downarrow$ $Mds_{1,0}$ | $Nd_{1,0}$ $\downarrow$ $Nds_{1,0}$ | $PValue_{1,0}$ += $Mds_{0,0}*Nds_{1,0}$ + $Mds_{1,0}*Nds_{1,1}$ | $Md_{3,0}$ $\downarrow$ $Mds_{1,0}$ | $Nd_{1,2}$ $\downarrow$ $Nds_{1,0}$ | $PValue_{1,0}$ += $Mds_{0,0}*Nds_{1,0}$ + $Mds_{1,0}*Nds_{1,1}$ |
| $T_{0,1}$ | $Md_{0,1}$ $\downarrow$ $Mds_{0,1}$ | $Nd_{0,1}$ $\downarrow$ $Nds_{0,1}$ | $PdValue_{0,1}$ += $Mds_{0,1}*Nds_{0,0}$ + $Mds_{1,1}*Nds_{0,1}$ | $Md_{2,1}$ $\downarrow$ $Mds_{0,1}$ | $Nd_{0,3}$ $\downarrow$ $Nds_{0,1}$ | $PdValue_{0,1}$ += $Mds_{0,1}*Nds_{0,0}$ + $Mds_{1,1}*Nds_{0,1}$ |
| $T_{1,1}$ | $Md_{1,1}$ $\downarrow$ $Mds_{1,1}$ | $Nd_{1,1}$ $\downarrow$ $Nds_{1,1}$ | $PdValue_{1,1}$ += $Mds_{0,1}*Nds_{1,0}$ + $Mds_{1,1}*Nds_{1,1}$ | $Md_{3,1}$ $\downarrow$ $Mds_{1,1}$ | $Nd_{1,3}$ $\downarrow$ $Nds_{1,1}$ | $PdValue_{1,1}$ += $Mds_{0,1}*Nds_{1,0}$ + $Mds_{1,1}*Nds_{1,1}$ |

time ⟶

Mds, Nds= shared

# 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 (#phases)
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[(m*TILE_WIDTH + ty)*Width + Col];
11.      __syncthreads();

12.    for (int k = 0; k < TILE_WIDTH; ++k)
13.        Pvalue += Mds[ty][k] * Nds[k][tx];
14.    __syncthreads();
     }
15. Pd[Row*Width + Col] = Pvalue;
}
```
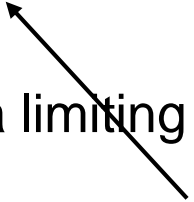
# CUDA Code – Kernel Execution Configuration
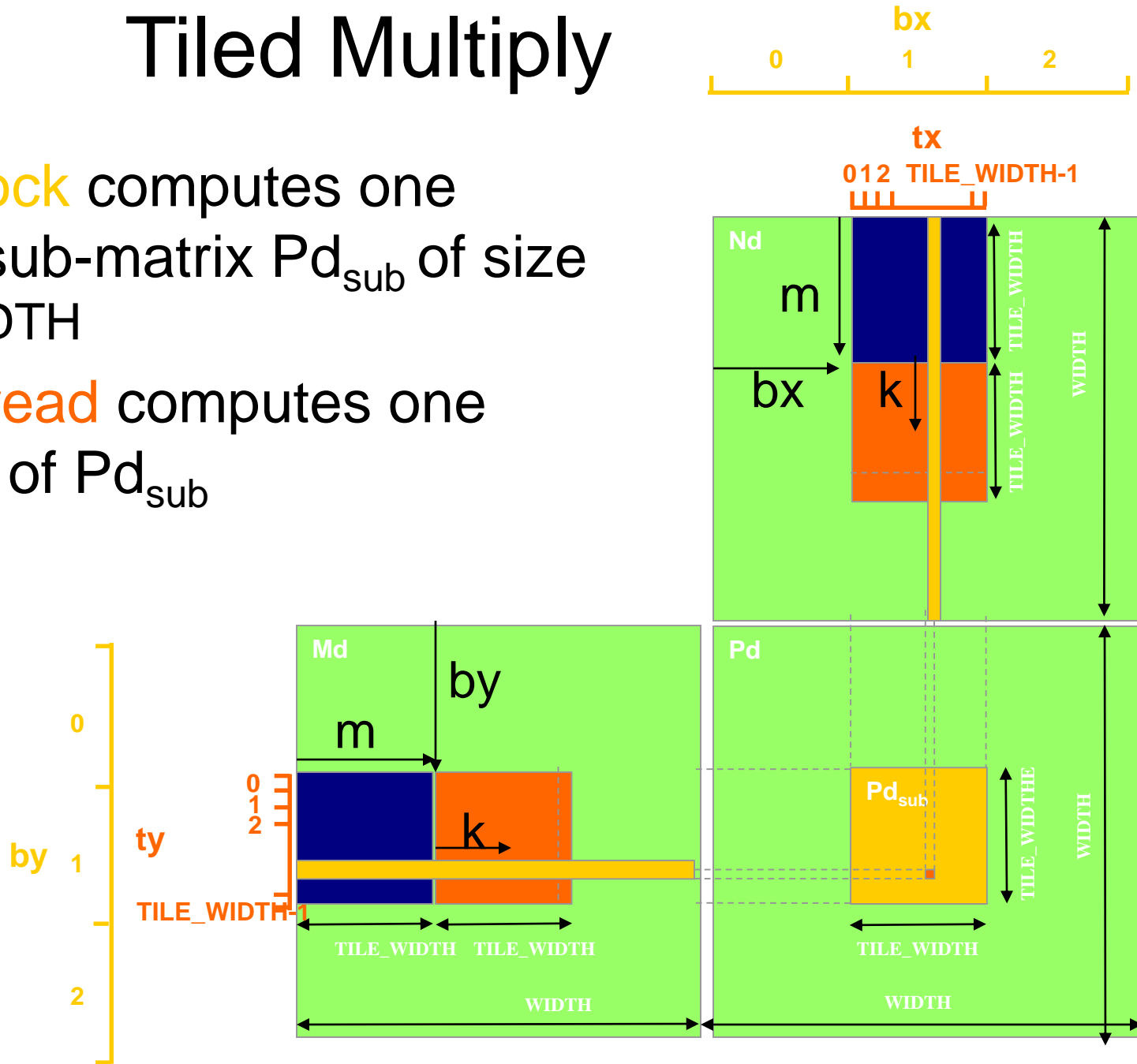
```
// Setup the execution configuration
dim3 dimBlock(TILE_WIDTH, TILE_WIDTH);
dim3 dimGrid(Width  / TILE_WIDTH,
             Width /  TILE_WIDTH);
```

# First-order Size Considerations in G80

- Each thread block should have many threads
  - Global memory accesses reduced by N = #tile width !
  - 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
  - TILE_WIDTH of 16 gives each SM 3 blocks, 768 threads (**full capacity**)

- 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

inner product

# Tiled Multiply

- Each block computes one square sub-matrix $Pd_{sub}$ of size TILE_WIDTH

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

# Memory as a limiting factor to parallelism

- The limited amount of CUDA shared memory (e.g. 16KB) limits the number of threads that can simultaneously reside in the SM!

- For the matrix multiplication example, shared memory **can become** a limiting factor:

- TILE_WIDTH = 16,  so each block requires 16x16x4 = 1kB of storage for Mds + 1kB for Nds
  - 2kB of shared memory per block

- The 16-kB shared memory allows 8 blocks to simultaneously reside in an SM. Ok!

- But the maximum number of threads per SM is 1024 (for Tesla T10)
  - For 1024*1024 matrix only 1024/256 = 4 blocks are allowed in each SM
  - only 4 x 2kB = 8kB of the shared memory will be used.

**Hint: Use occupancy calculator**

# 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[(m*TILE_WIDTH + ty)*Width + Col];
11.      __syncthreads();

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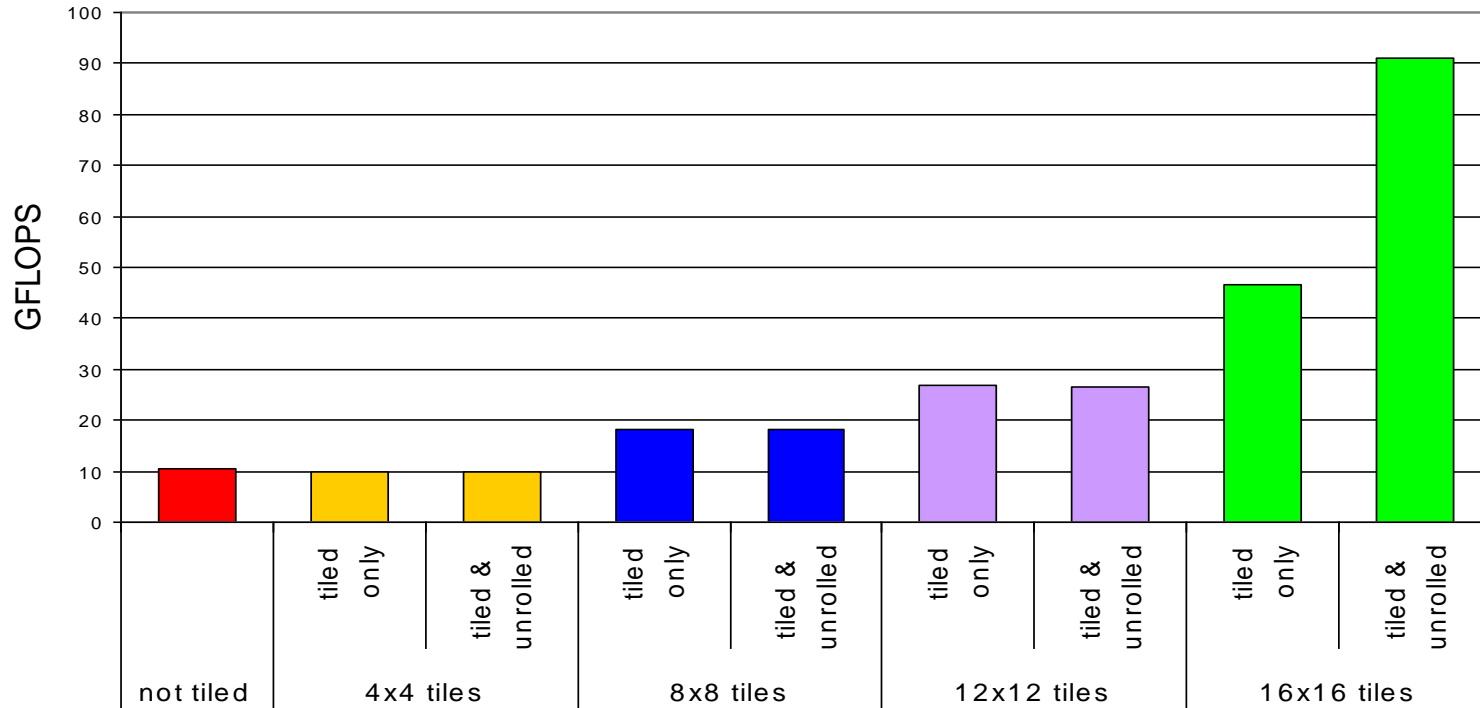
# Tiling Size Effects

# Summary- Typical Structure of a CUDA Program

- Global variables declaration
  - __host__
  - __device__... __global__, __constant__, __texture__
- Function prototypes
  - __global__ void kernelOne(…)
  - float handyFunction(…)
- Main ()
  - allocate memory space on the device – cudaMalloc(&d_GlblVarPtr, bytes )
  - transfer data from host to device – cudaMemCpy(d_GlblVarPtr, h_Gl…)
  - execution configuration setup
  - kernel call – kernelOne<<<execution configuration>>>( args… );
  - transfer results from device to host – cudaMemCpy(h_GlblVarPtr,…)
  - optional: compare against <u>golden</u> (host computed) solution
- Kernel – void kernelOne(type args,…)
  - variables declaration -  __local__, __shared__
    - automatic variables transparently assigned to registers or local memory
  - syncthreads()…
- Other functions
  - float handyFunction(int inVar…);

repeat
as
needed