# COL380

# Introduction to
# Parallel & Distributed Programming

# GPU and manycore

CPU+GPU

PCH: PCI, SCSI/ SATA, USB, Audio..

DMI

DMA Engine

16 GBps

4x 16-lane PCIe

4x 16-lane PCIe
16 GB/s
Each direction

| Memory Controller | C | C | C | •• | C | Network Controller |
| | L1 | L1 | L1 | | L1 | |
| | L2 | L2 | L2 | | L2 | |
| | L3 | | | | | |
| | C | C | C | •• | C | |
| | L1 | L1 | L1 | | L1 | |
| | L2 | L2 | L2 | | L2 | |

UPI links
= 33.6 G
Each dir

| Memory Controller | C | C | C | •• | C | Network Controller |
| | L1 | L1 | L1 | | L1 | |
| | L2 | L2 | L2 | | L2 | |
| | L3 | | | | | |
| | C | C | C | •• | C | |
| | L1 | L1 | L1 | | L1 | |
| | L2 | L2 | L2 | | L2 | |

1A   1B        2H
2A   2B        2H

Memory

Ch A   Ch B        Ch H

1A   1B        2H
2A   2B        2H

Memory

Ch A   Ch B        Ch H

Subodh Kumar

PCI Express 3.0 Host Interface

GigaThread Engine

GPC

TPC TPC TPC TPC TPC TPC TPC

SM

L2 Cache

GPC GPC GPC

High-Speed Hub

NVLink   NVLink   NVLink   NVLink   NVLink   NVLink

Memory Controller   Memory Controller   Memory Controller   Memory Controller

HBM2   HBM2   HBM2   HBM2

# GPU



Source: Nvidia

```
#pragma omp target teams num_teams(n) distribute
#pragma omp parallel for
   for(int i=0; i<N; i++)
      vec2[i] += vec1[i];
```

```
accelerate(sum, size, vec, vec2);
```

```
sum (const int size, __global float * vec1, __global float * vec2)
{
      int ii = get_global_id(0);
      if (ii < size) vec2[ii] += vec1[ii];
}
```

// Matrix multiplication kernel – per thread code

```
__global__ void
MatrixMulKernel(float* Md, float* Nd, float* Pd, int Width
{
    // Pvalue stores the matrix element 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;
}
```
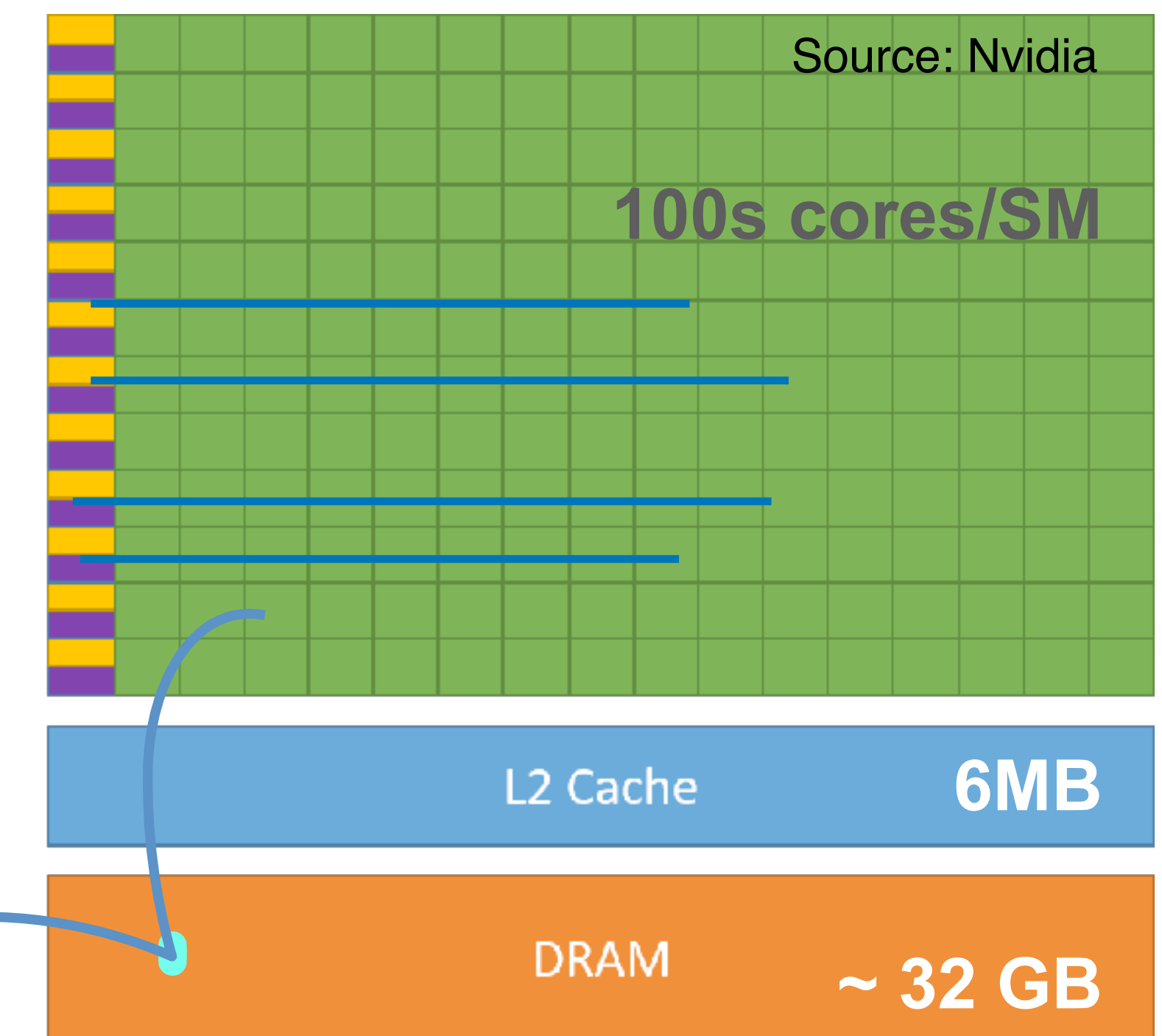
Subodh Kumar

- # Run many threads

  ‣ Memory latency needs to be hidden

  ‣ SIMD => Explicit parallel
  read/write

~ 128KB L1/SM

Source: Nvidia

100s cores/SM

L1: 32 KB/core
L2: 512 KB/core

Core | Control
L1 Cache
Core | Control
L1 Cache
Core | Control
L1 Cache
Core | Control
L1 Cache
L2 Cache
L2 Cache
L3 Cache — 2 MB/core
DRAM — Many GB
CPU

L2 Cache — 6MB
DRAM — ~ 32 GB
GPU

Subodh Kumar

Source: Nvidia
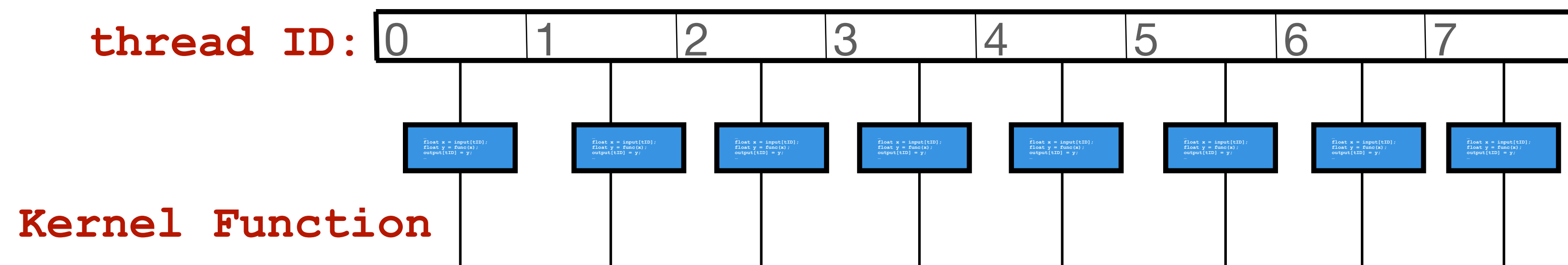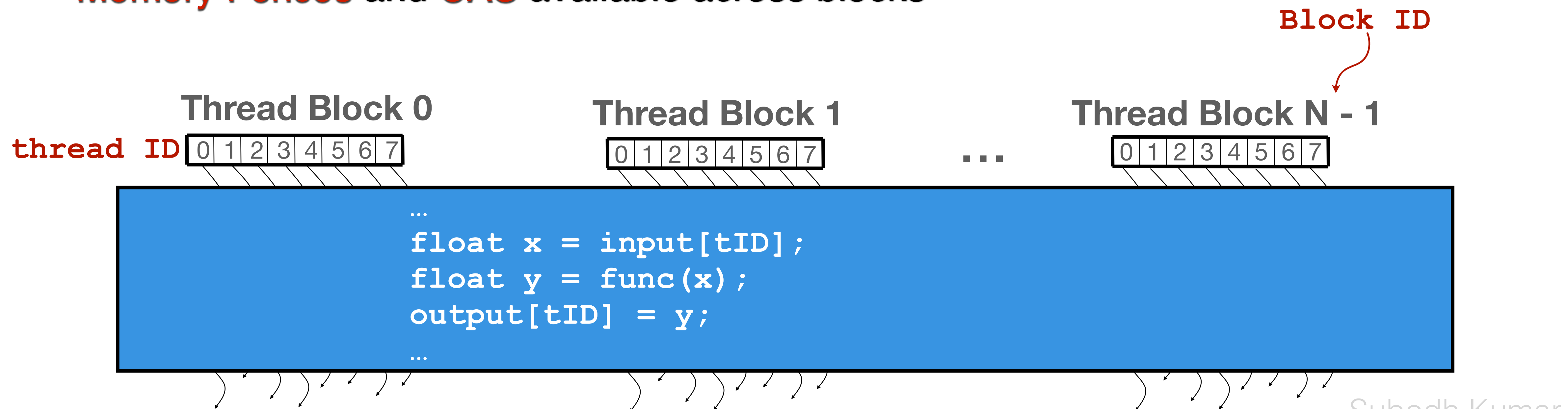
- Co-processor

  ‣ Many cores

  ‣ CPU code offloads multi-threaded tasks

    ➨ Message passing and shared memory models

- GPU Threads are organized hierarchically

  ‣ Grids, Blocks, Warps

  ‣ Data-parallelism focus

- Shared memory

  ‣ Memory hierarchy exposed

  ‣ Parallel read/write

- Each GPU is also called a device

  ‣ CPU is host

- GPU memory is device memory

- Device executes GPU code (kernel)

  ‣ Executed by many threads in parallel

- GPU is massively parallel

  ‣ Thousands of cores

- GPU threads are lightweight

  ➡ Very little creation and scheduling overhead

  ➡ Context remains live

  ➡ Many more threads then the number if processors

- # CPU thread 'forks' a grid of threads, each executes the kernel

  - ▸ Grid executes asynchronously, CPU thread continues

  - ▸ Grid has a number of blocks, each block has a number of threads

    - ➥ Threads execute in groups of 32 called warps

    - ➥ Execute as SIMT: Single instruction multiple threads
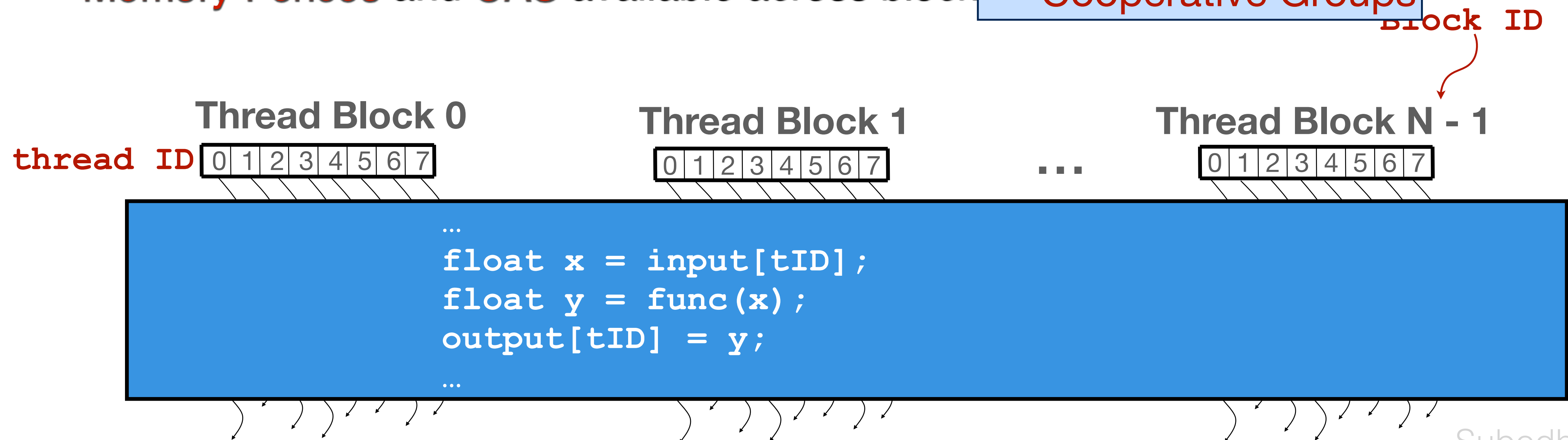
      - But a warp can diverge

thread ID: | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

Kernel Function

- ## Threads within a block share <span style="color:#b22">shared memory</span> and <span style="color:#b22">global memory</span>

  - ▸ Threads in different blocks only access a common global memory

  - ▸ Intra-warp register-based consensus

- ## <span style="color:#b22">Barrier</span> and finer-grained synchronization only within warps and blocks

  - ▸ <span style="color:#b22">Memory Fences</span> and <span style="color:#b22">CAS</span> available across blocks

**Block ID**

**Thread Block 0**     **Thread Block 1**     **Thread Block N - 1**

thread ID | 0 1 2 3 4 5 6 7 |     | 0 1 2 3 4 5 6 7 |     ...     | 0 1 2 3 4 5 6 7 |

```
…
float x = input[tID];
float y = func(x);
output[tID] = y;
…
```
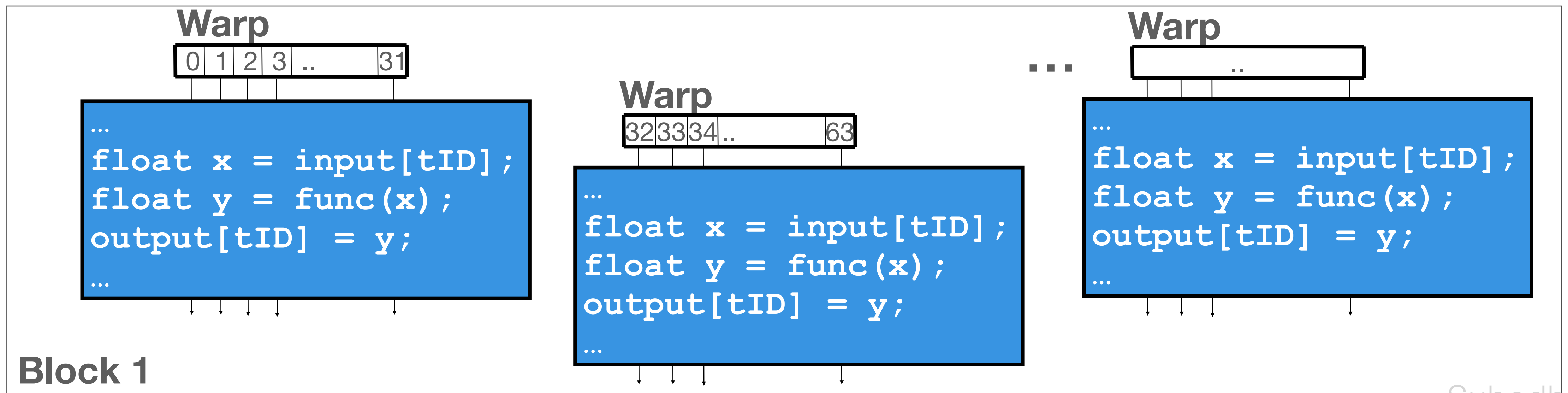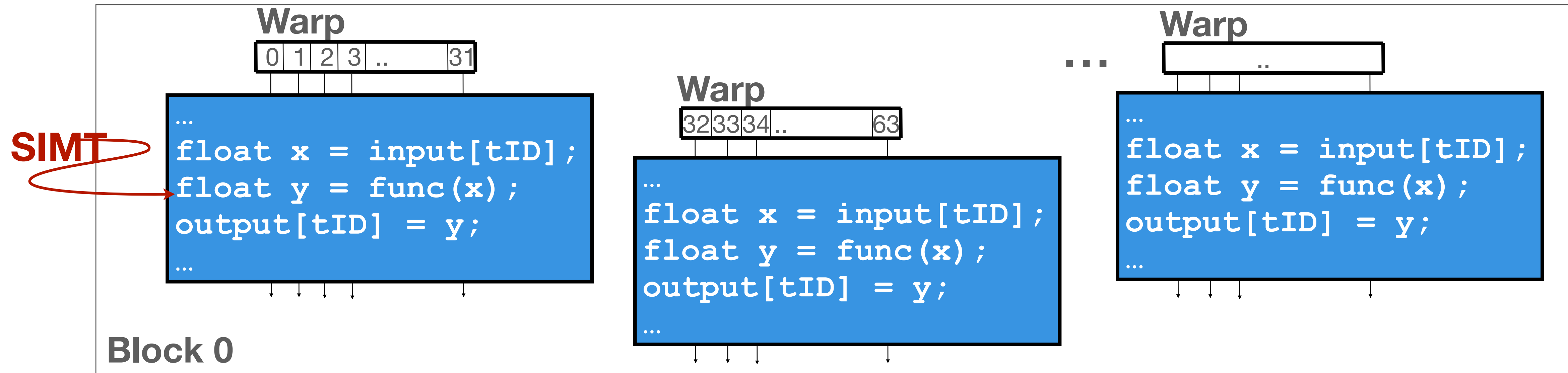
Subodh Kumar

- Threads within a block share shared memory and global memory
  - ▸ Threads in different blocks only access a common global memory
  - ▸ Intra-warp register-based consensus

- Barrier and finer-grained synchronization only within warps and blocks
  - ▸ Memory Fences and CAS available across block
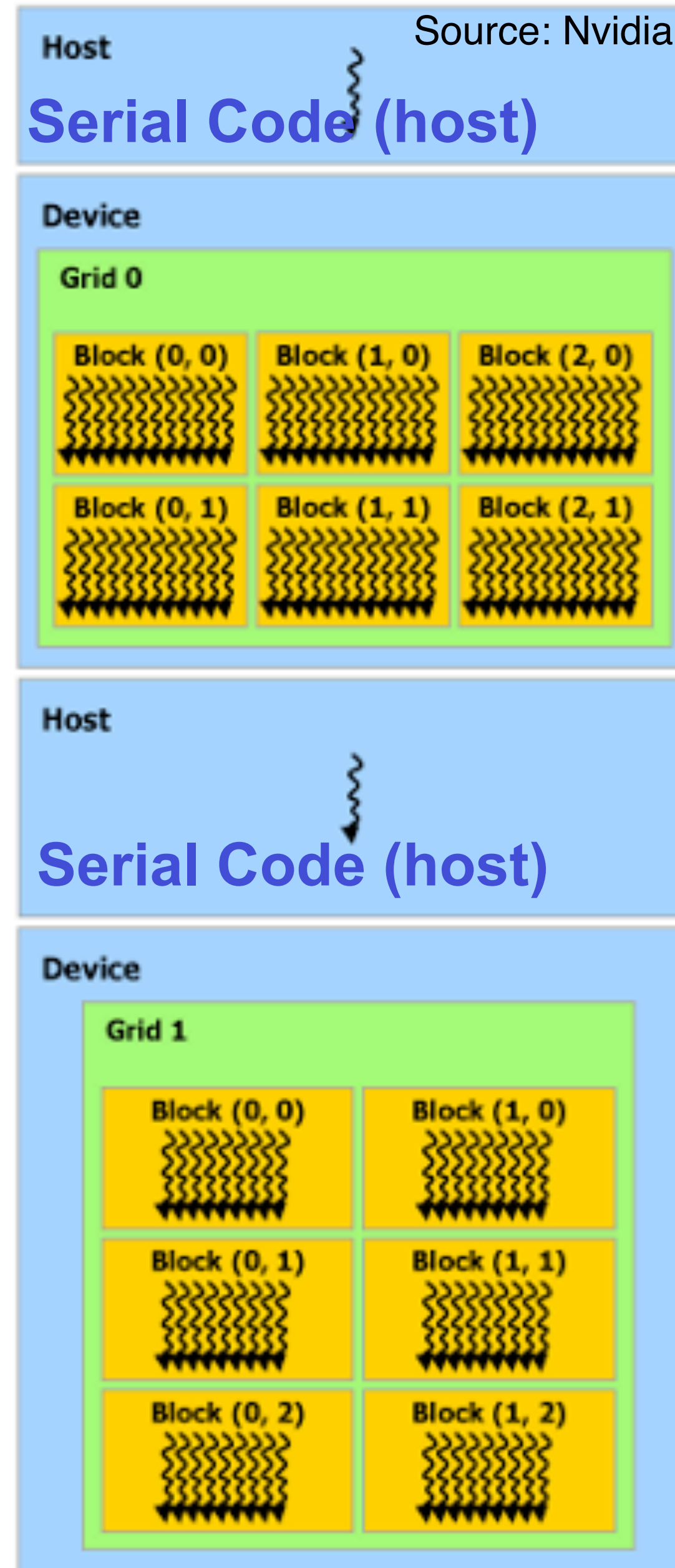
Also see:
    Cooperative Groups

Block ID

**Thread Block 0**

thread ID `0 1 2 3 4 5 6 7`

**Thread Block 1**

`0 1 2 3 4 5 6 7`

**...**

**Thread Block N - 1**

`0 1 2 3 4 5 6 7`

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

Subodh Kumar

# Warps

**Warp**

| 0 | 1 | 2 | 3 | .. | | 31 |

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

**SIMT**

**Warp**

| 32 | 33 | 34 | .. | | 63 |

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

**Warp**

| | .. | |

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

...

**Block 0**

**Warp**

| 0 | 1 | 2 | 3 | .. | | 31 |

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

**Warp**

| 32 | 33 | 34 | .. | | 63 |

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

**Warp**

| | .. | |

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

...

**Block 1**

# CUDA is C-like



Source: Nvidia

**Execute on GPU:**

**KernelA<<< nBlk, nTid >>>(args);**

**Execute on GPU**

**KernelB<<< nBlk, nTid >>>(args);**

- Integrated host+device app Cuda program
  - ‣ Serial or modestly parallel parts in host C code
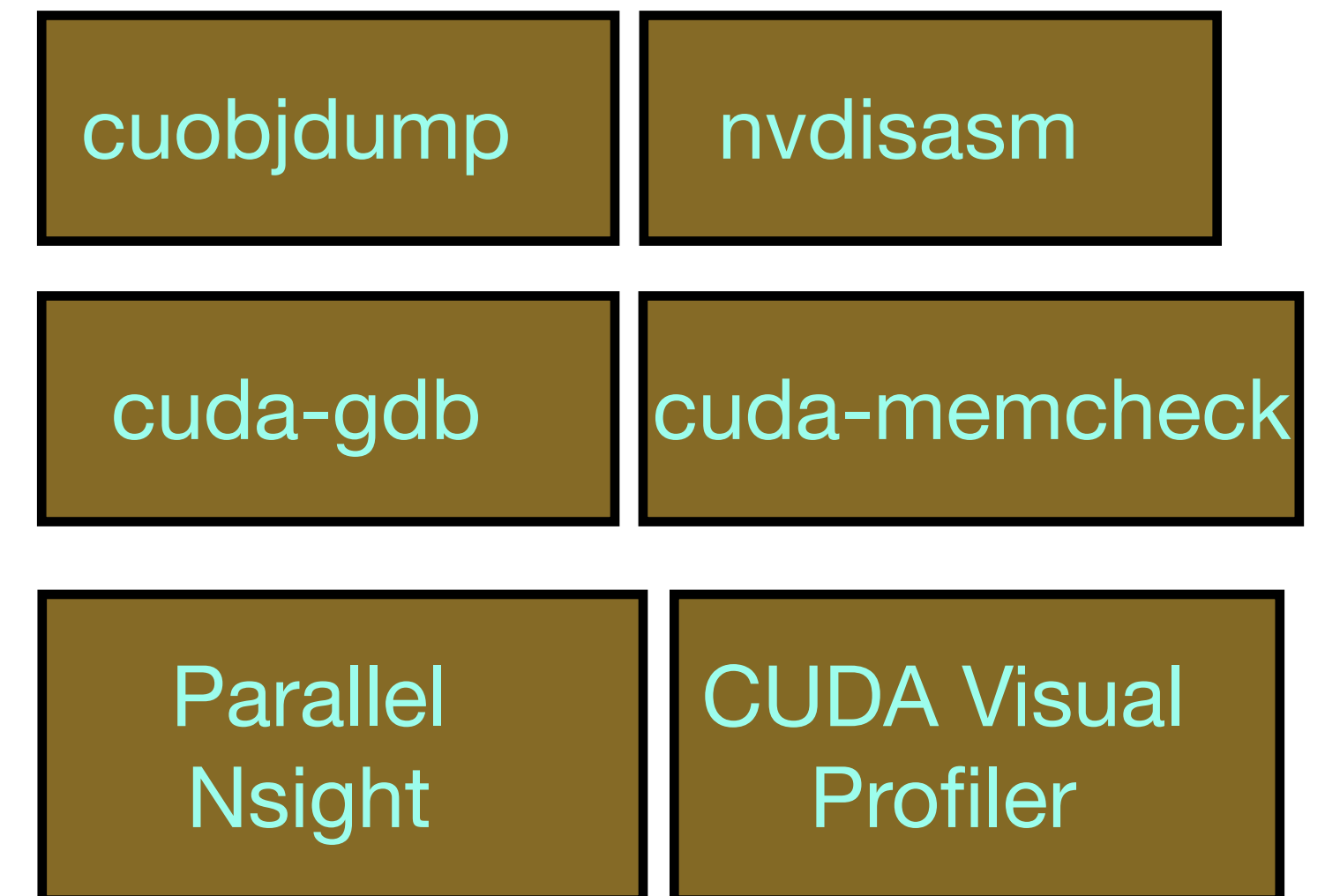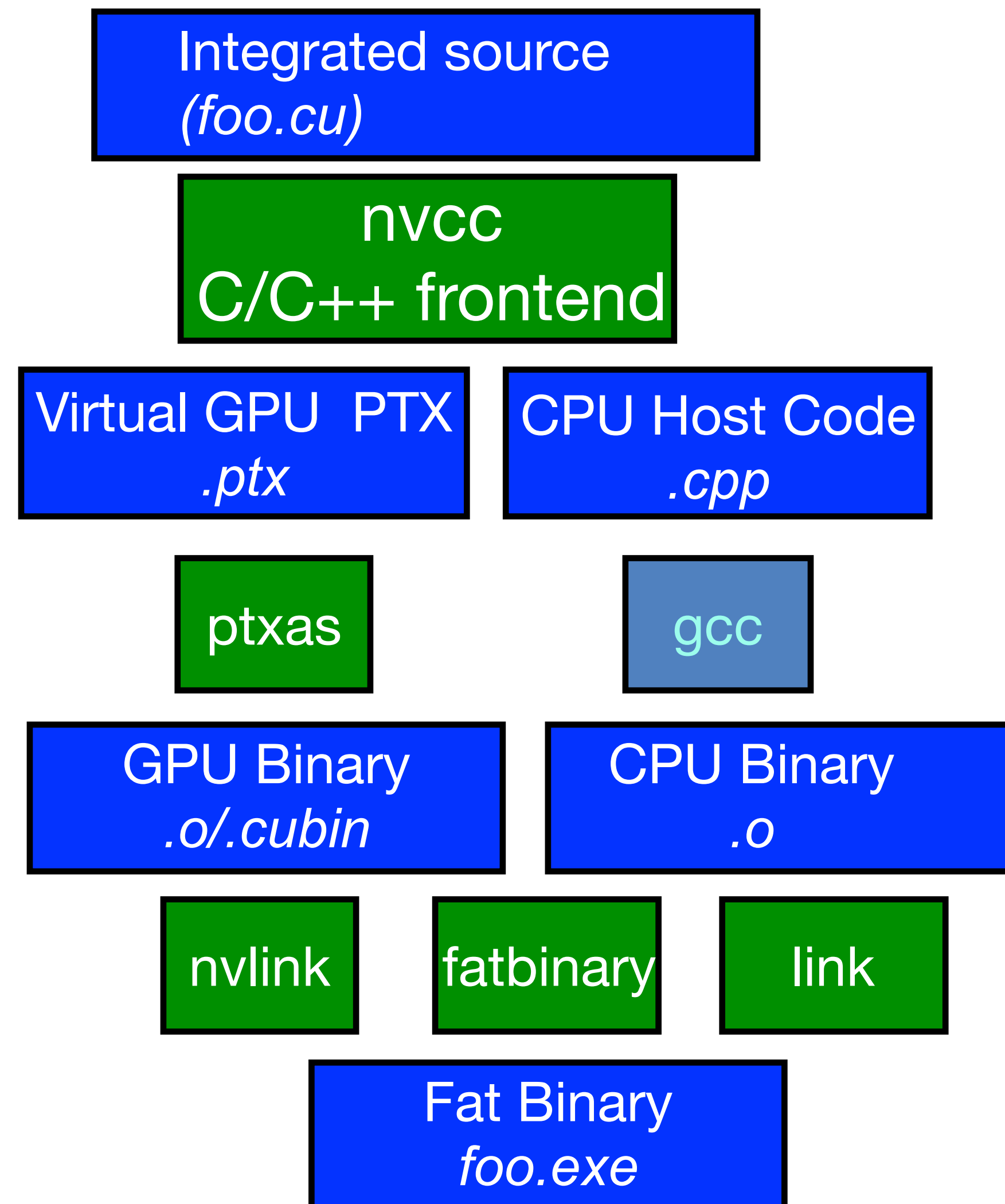  - ‣ Highly parallel parts in device Cuda code

Subodh Kumar

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

- Built-in variables
  - threadIdx, blockIdx, blockDim

- Intrinsics
  - __syncthreads

- Runtime API
  - Memory, symbol, execution management

- Kernel launch

Built-in Types:
char2, int4,
__half2, dim3

```
__device__ float filter[N];

__global__ void convolve (float *image){

  __shared__ float region[M];
  ...
  region[threadIdx.x] = image[i];
  __syncthreads()
  ...
  image[j] = result;
}
...
// Allocate GPU memory
void *myimage = cudaMalloc(bytes)

// 100 blocks, 10 threads per block
convolve<<<100, 10>>> (myimage);
```

- Source files have a mix of host and device code

- nvcc separates device code from host code
  ‣ compiles device code into PTX/cubin
  ‣ host code is output as C source (and invoke compiler)

- Applications link to the generated host code
  ‣ host code includes PTX/cubin code as a global initialized data array
  ‣ and cudart (CUDA C runtime) function calls to load and launch kernels

- Possible to load and execute the PTX/cubin using the CUDA driver API

# Cuda Dev Tools

Integrated source
*(foo.cu)*

nvcc
C/C++ frontend

Virtual GPU  PTX
*.ptx*

CPU Host Code
*.cpp*

ptxas

gcc

GPU Binary
*.o/.cubin*

CPU Binary
*.o*

nvlink

fatbinary

link

Fat Binary
*foo.exe*

cuobjdump

nvdisasm

cuda-gdb

cuda-memcheck

Parallel
Nsight

CUDA Visual
Profiler

Subodh Kumar

```
cudaEvent_t start, stop;
cudaEventCreate(&start);
cudaEventCreate(&stop);
cudaEventRecord(start, 0);
GPUThings();
cudaEventRecord(stop, 0);
cudaEventSynchronize(stop);
float ms;
cudaEventElapsedTime(&ms, start, stop);
cudaEventDestroy(start);
cudaEventDestroy(stop);
```

| | |
|---|---|
| cuobjdump | nvdisasm |
| cuda-gdb | cuda-memcheck |
| Parallel Nsight | CUDA Visual Profiler |

ptxas

gcc

GPU Binary
*.o/.cubin*

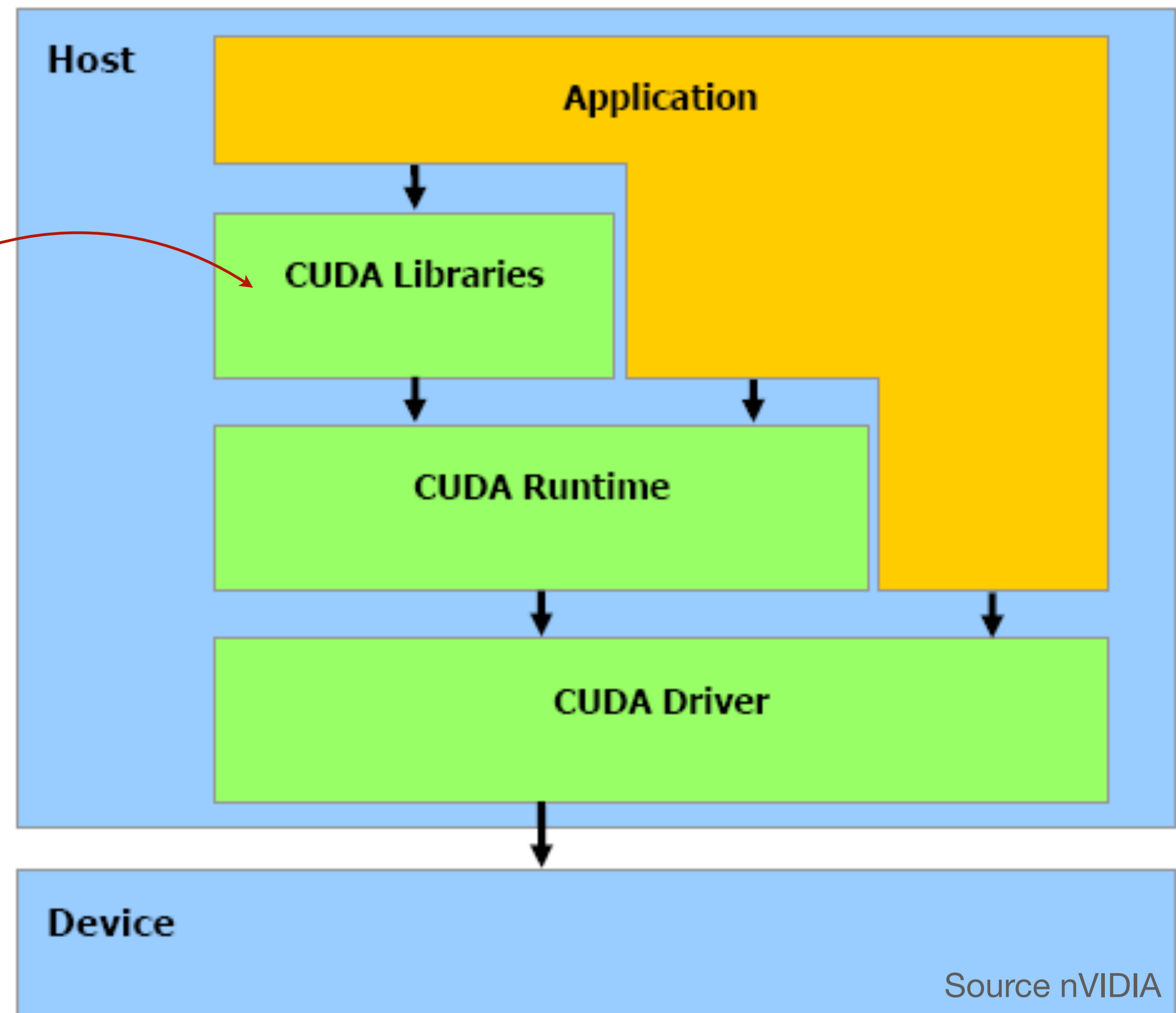GPU Binary

nvlink

foo.exe

On HPC: Need GPU nodes
    Login: gpu.hpc.iitd.ac.in (K40, CC3.5)
    module add compiler/cuda/11.0/compilervars
    Available: V100 (CC7.0)

Subodh Kumar

Provides library functions for host as well as device
Implement subset of stdlib



Host

Application

CUDA Libraries

CUDA Runtime

CUDA Driver

Device

Source nVIDIA

Subodh Kumar

- Run *k* instances of kernel function *f*

  ‣ Each instance in a thread

  ‣ Kernel 'pushed' to GPU

  ➡ Declared with **__global__** specifier

- *k* threads are organized into blocks

```
int main()
{
    cpuCode();

    f<<<1, N>>>(F); // GPU Kernel
}
```

Subodh Kumar

- Run *k* instances of kernel function *f*
  - ‣ Each instance in a thread
  - ‣ Kernel 'pushed' to GPU
    - ➡ Declared with **__global__** specifier

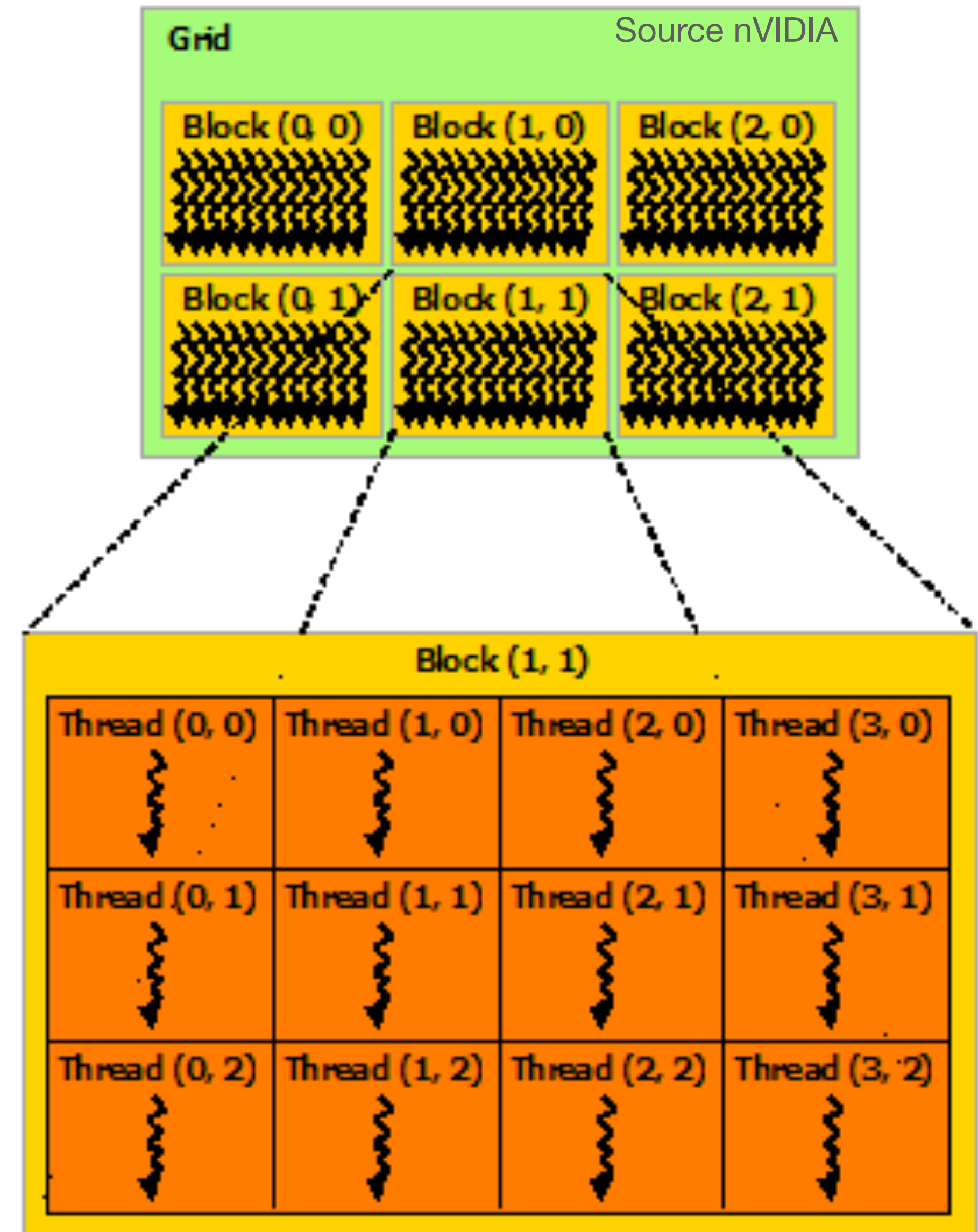- *k* threads are organized into blocks

```
int main()
{
    cpuCode();

    f<<<1, N>>>(F); // GPU Kernel

}
```

```
// Kernel definition
__global__ void f(float* F)
{
    int id = threadIdx.x;
    …
}
```

Subodh Kumar

- Run *k* instances of kernel function *f*
  - ‣ Each instance in a thread
  - ‣ Kernel 'pushed' to GPU
    - ➡ Declared with **__global__** specifier
- *k* threads are organized into blocks

Also see: cudaGraphcudaGraphCreate(..)
Specify task-graph of
Kernels, CPU functions,
memory operations,
Synchronization

```
int main()
{
    cpuCode();

    f<<<1, N>>>(F); // GPU Kernel
}
```

```
// Kernel definition
__global__ void f(float* F)
{
    int id = threadIdx.x;
    …
}
```

- **Invocation: <<<A,B>>>**

- **A** and **B** can be (up to) 3-D vectors

  - ‣ **dim3** B(a, b, c); // a,b,c are ints

    - ➡ a x b x c <= 1024 for blocks

    - c is the most significant dimension, a the least

    - Dereference: **B**.x, **B**.y and **B**.z

    - Thread ID = (**B**.x + **B**.y * a + **B**.z * a*b)

    - Inbuilt: **dim3** threadIdx, blockIdx, blockDim;



Source nVIDIA

- All threads of a block need not execute in SIMD fashion

- Threads executed in groups of upto 32 parallel threads (There need not be 32 physical cores)

  ‣ called *warps*

- All threads of a warp start together

  ‣ But may diverge due to conditional branching

  ‣ Different sub-warps are serialized until they converge back   See __syncwarp()

  ‣ Important efficiency consideration