# Programming Languages for Accelerators
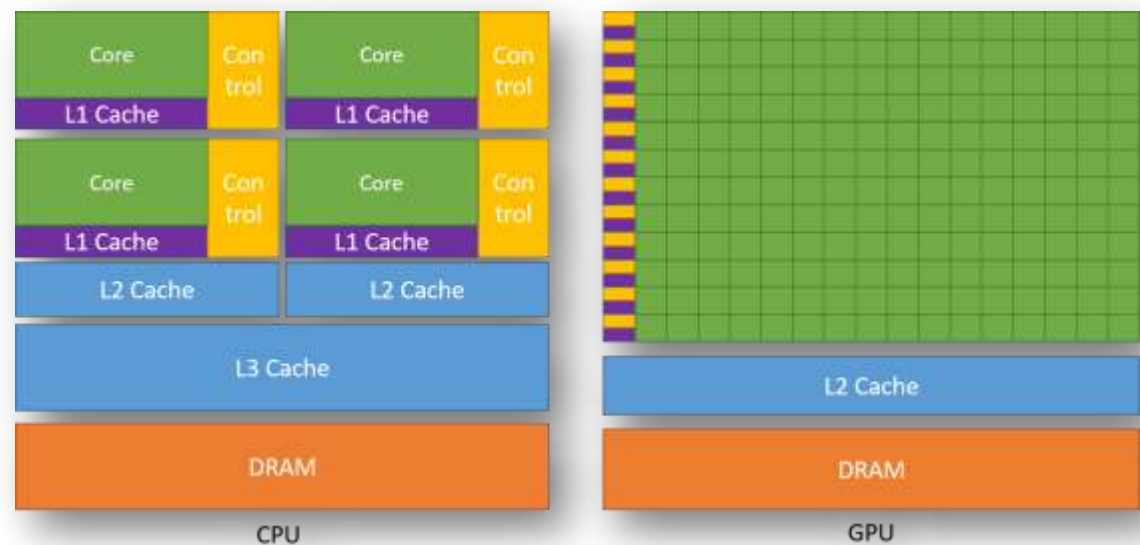
## Manos Pavlidakis[1,2]

manospavl@ics.forth.gr

[1] Institute of Computer Science, Foundation for Research and Technology - Hellas, Greece
[2] Computer Science Department, University of Crete, Greece

# What is an accelerator?

- A device that performs <u>some functions</u> more efficiently than general-purpose CPU
  - Due to massive pallelism
  - GPUs are perfect for Vector Add

- General Purpose Graphic Processing Unit (NVIDIA, AMD)

- Field-Programmable Gate Array (Xilinx, Intel Altera)

- Application-Specific Integrated Circuit
  - TPU: Tensor Processing Unit (Google)



[Programming Guide :: CUDA Toolkit Documentation (nvidia.com)](#)

# How we program accelerators?

# Well known programming languages

- CUDA → NVIDIA GPUs

- HIP:
    - ROcM → AMD GPUs
    - CUDA → NVIDIA GPUs

- oneAPI/SYCL → Heterogeneous accelerators

# CUDA: Compute Unified Device Architecture

# What is CUDA?

- CUDA
  - Based on industry-standard C/C++
  - Small set of extensions to enable heterogeneous programming → Kernel code
  - APIs to manage devices → Transfers, Allocations etc.

# Terminology

- Host: The CPU and its memory
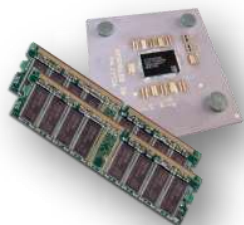- Device: The GPU/Accelerator and its memory
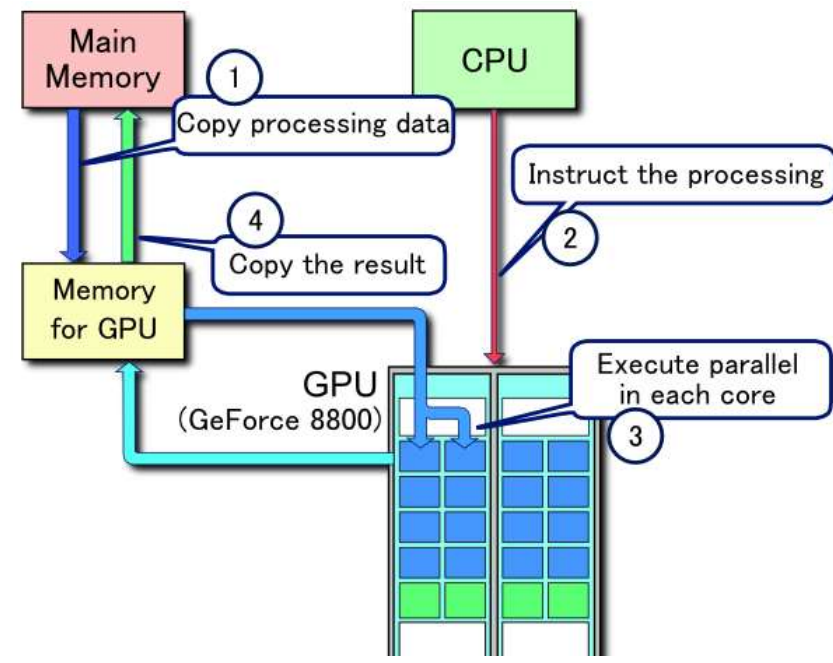
Host

Device

# Host and Device code



```
 1  #define THREADS_PER_BLOCK 256
 2  const int N=2048;
 3
 4  int main(void) {
 5      int *a, *b, *c; // host copies of a, b, c
 6      int *d_a, *d_b, *d_c; // device copies of a, b, c
 7      int size = N * sizeof(int);
 8
 9      // Alloc space for device
10      cudaMalloc((void **)&d_a, size);
11      cudaMalloc((void **)&d_b, size);
12      cudaMalloc((void **)&d_c, size);
13
14      // Alloc space on host
15      a = (int*)malloc(size); random_ints(a, N);
16      b = (int*)malloc(size); random_ints(b, N);
17      c = (int*)malloc(size);
18
19      // Copy inputs to device
20      cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice);      ①
21      cudaMemcpy(d_b, b, size, cudaMemcpyHostToDevice);
22
23      // Launch add() kernel on GPU
24      add<<<N/THREADS_PER_BLOCK,THREADS_PER_BLOCK>>>(d_    ② b, d_c);
25
26      // Copy result back to host
27      cudaMemcpy(c, d_c, size, cudaMemcpyDeviceToHost);      ④
28
29      // Cleanup
30      free(a); free(b); free(c);
31      cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
32      return 0;
33  }
34  //Kernel
35  __global__ void add(int *a, int *b, int *c) {
36      int index = threadIdx.x + blockIdx.x * blockDim.x;    ③
37          c[index] = a[index] + b[index];
38  }
```

Host code

Device code

(CUDA kernel)

Executed by one
**Host** Thread

Executed by multiple
**CUDA** Threads

# CUDA by an example

# Add two integers with CUDA

- A simple kernel to add two integers

```
__global__ void add(int *a, int *b, int *c){

    *c = *a + *b;

}
```

- __global__ is a CUDA C/C++ keyword meaning:
  - add() will execute on the device
  - add() will be called from the host
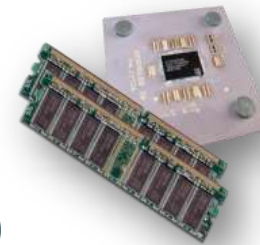
# Add two integers with CUDA

- 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 <u>device</u> memory

- We need to allocate memory on the GPU for a, b, c

# Memory Management

- Host and device memory are separate entities

- **Device** pointers point to GPU memory
  - May be passed to/from **host** code

- **Host** pointers point to CPU memory
  - May be passed to/from **device** code

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

# Add two integers: main()

```
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;
```

# Add two integers: main()

```
// 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 results back to host

cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

// Free device memory

cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);

return 0;

} // End of main()
```

# Use parallelism

- Performance gain of GPUs is based on massive parallelism
  - CPUs have several cores(i.e. hundreds)
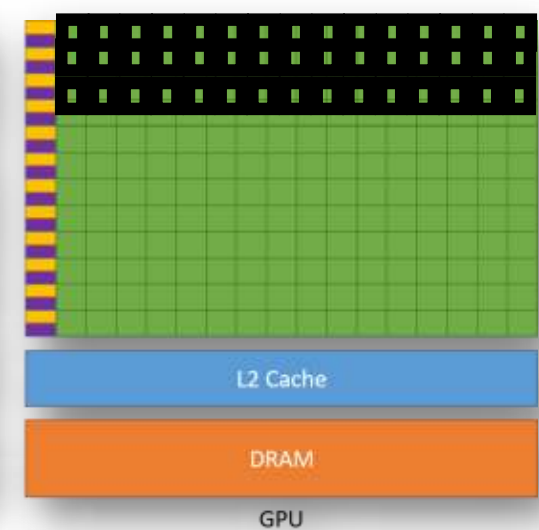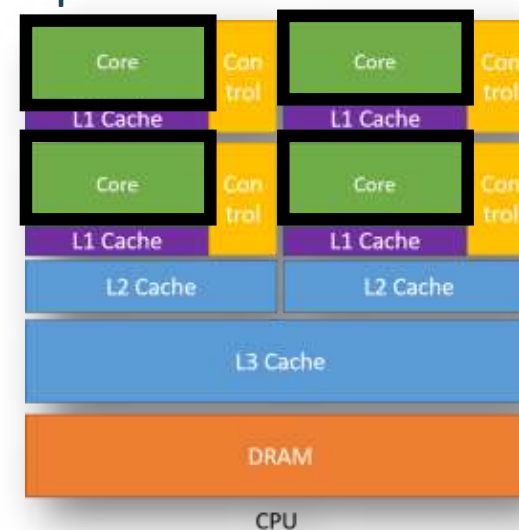  - GPUs have many cores(i.e. thousands)

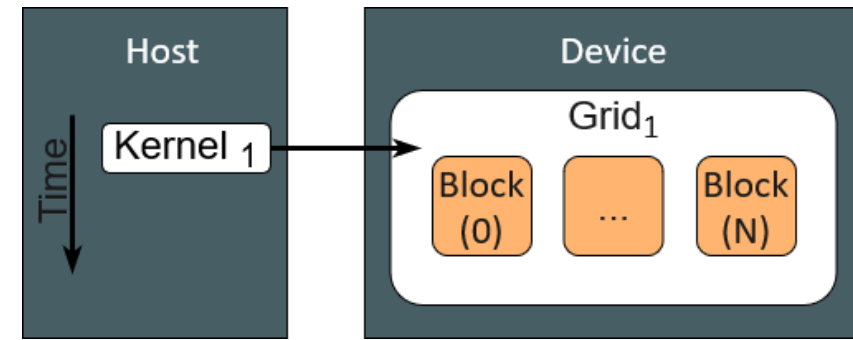- How we run the add() kernel on many cores?

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

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

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

# Vector addition on the Device



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

- 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

# Sequential Vector Addition using C++

```cpp
int main(void) {

        int a, b, c;

        int size = sizeof(int);

        // Allocate and initialize a, b, c

        a = (int*)malloc(size);  random_ints(a, N);

        b = (int*)malloc(size);  random_ints(b, N);

        c = (int*)malloc(size);

        for(int i = 0; i<size; i++)// Add arrays

                c[ i ] = a[ i ] + b[ i ];

        free(a); free(b); free(c);

        return 0;

}
```
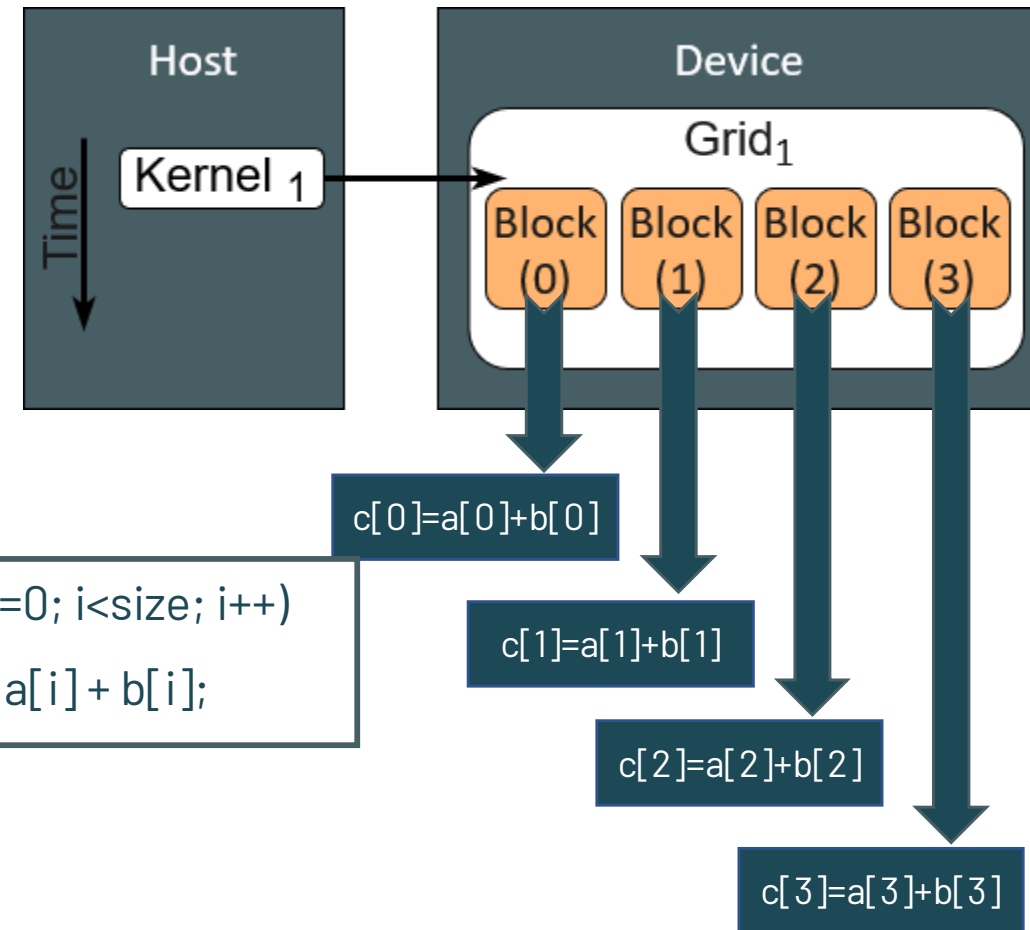
# Vector addition with 4 blocks

- If we want to create 4 parallel blocks
- We will call the kernel from host with N=4

  add<<< **4**, 1 >>>(…);

- On the device the kernel code

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

```
for(int i =0; i<size; i++)
    c[i] = a[i] + b[i];
```



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 main()

```c
#define N 4

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 random input values
    a = (int *)malloc(size); random_ints(a, N);
    b = (int *)malloc(size); random_ints(b, N);
    c = (int *)malloc(size);
```
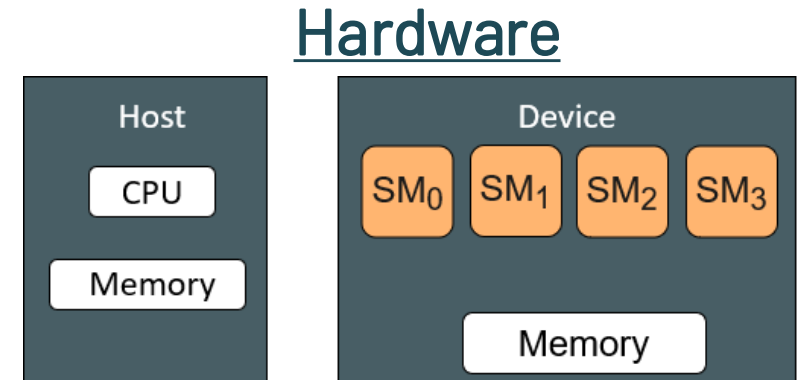
# Vector addition main()

```
// 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);

// Free device and host memory

free(a); free(b); free(c);

cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);

return 0;

}
```
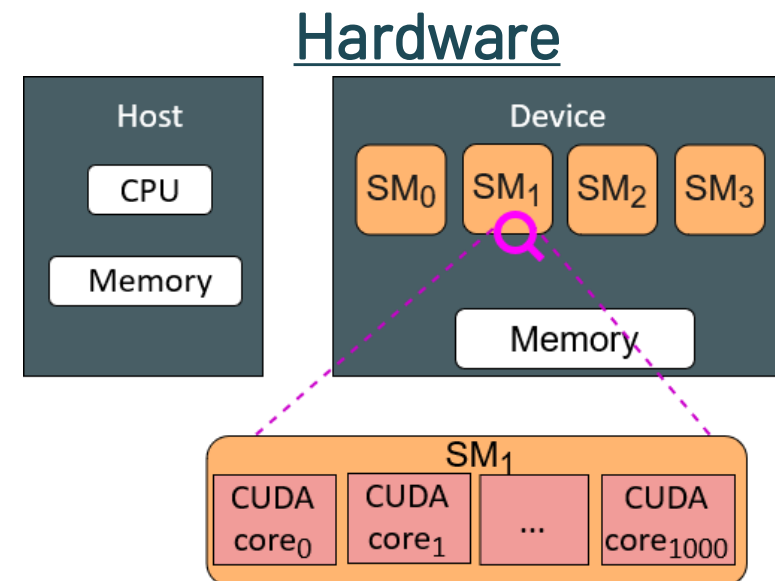
# More parallelism using CUDA cores

**Hardware**

- NVIDIA GPUs consists of:
  - Streaming Multiprocessors(SM)
  - An SM consists of CUDA cores

Host
CPU
Memory

Device
$SM_0$ $SM_1$ $SM_2$ $SM_3$
Memory

# More parallelism using CUDA cores



Hardware

- NVIDIA GPUs consists of:
  - Streaming Multiprocessors(SM)
  - An SM consists of CUDA cores
  - A CUDA core is like a **thin** processor

- A GPU has:
  - Hundreds of SMs (e.g. A100 GPU has **128)**
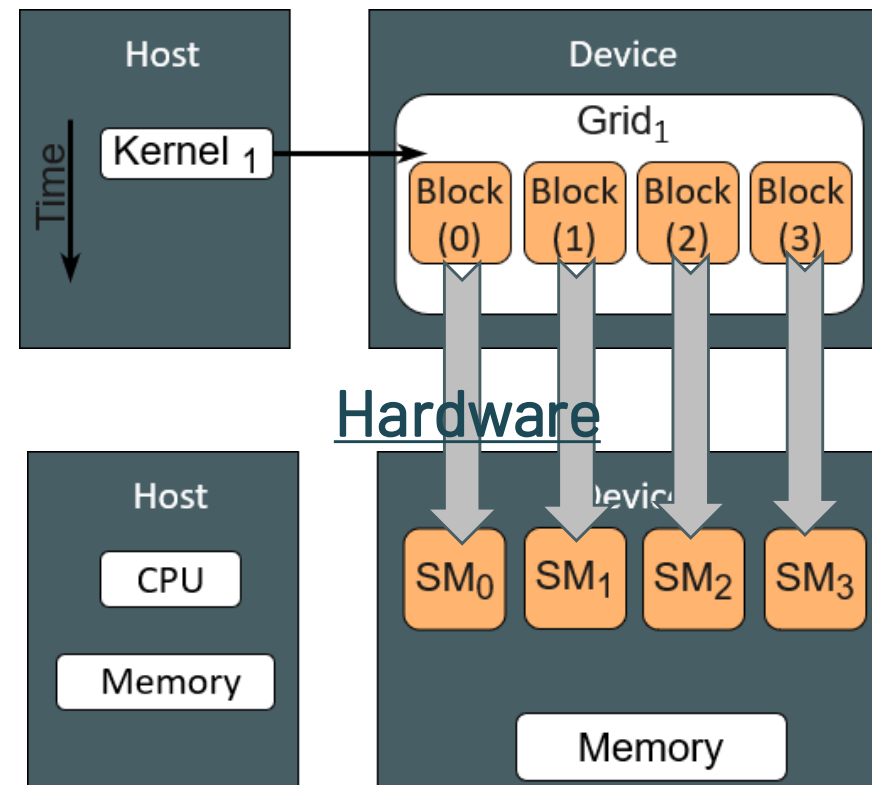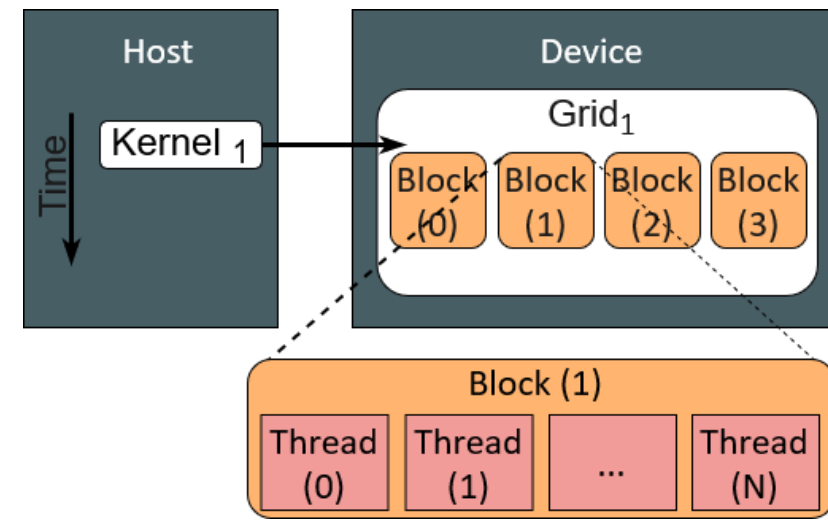  - Thousand CUDA cores (e.g. A100 GPU has **7000)**

# More parallelism using CUDA cores

- NVIDIA GPUs consists of:
  - Streaming Multiprocessors(SM)
  - An SM consists of CUDA cores
  - A CUDA core is like a **thin** processor

- A GPU has:
  - Hundreds of SMs (e.g. A100 GPU has **128)**
  - Thousand CUDA cores (e.g. A100 GPU has **7000)**

- In our example we use only Blocks
  - A Block is assigned to an SM
  - Blocks in the same SM execute <u>concurrently</u>!
    - Not in parallel!!
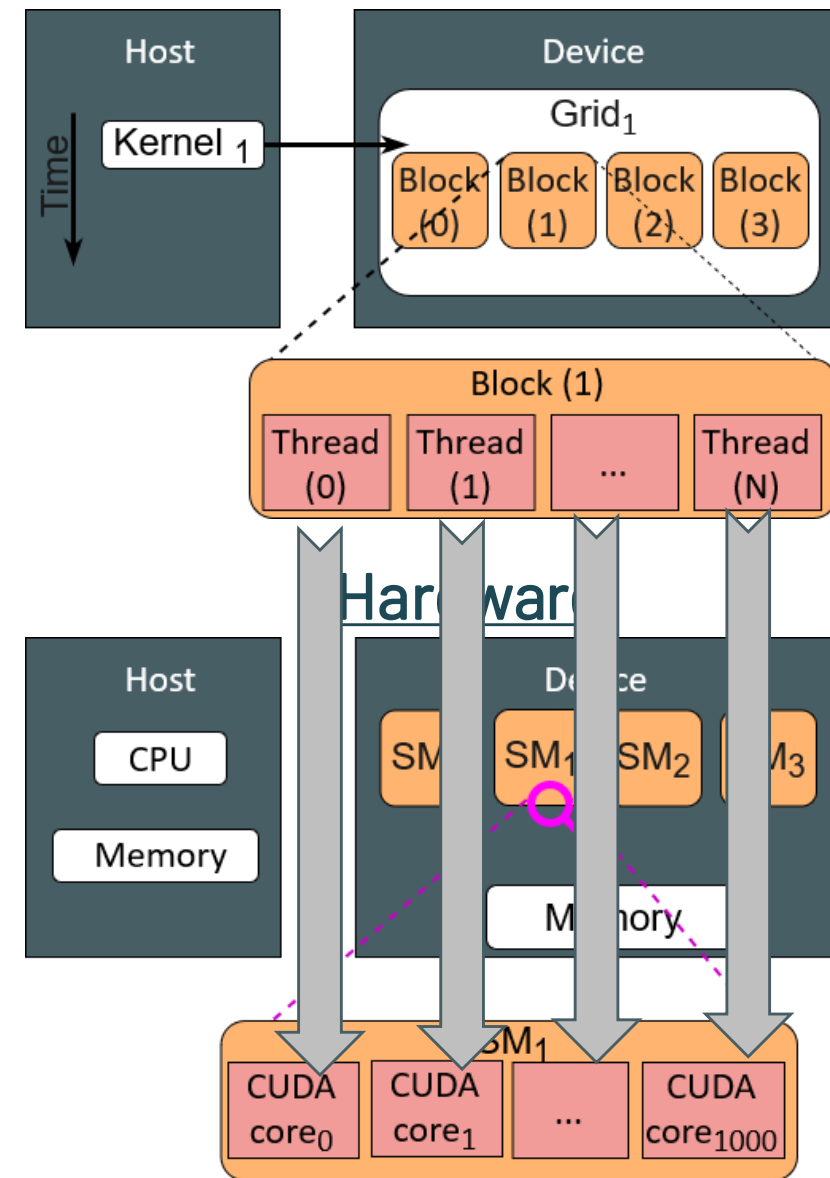    - If we have 4xSMs only 4xBlocks run in parallel

- Why not use CUDA cores !

Hardware

# CUDA Threads → CUDA cores



- A Block can be split into parallel Threads
  - Threads in the same block can cooperate
  - Threads have unique ids (i.e. threadId.x)

# CUDA Threads → CUDA cores

- A Block can be split into parallel Threads
  - Threads in the same block can cooperate
  - Threads have unique ids (i.e. threadId.x)
- Threads are assigned to CUDA cores

# CUDA Threads → CUDA cores

- A Block can be split into parallel Threads
  - Threads in the same block can cooperate
  - Threads have unique ids (i.e. threadId.x)
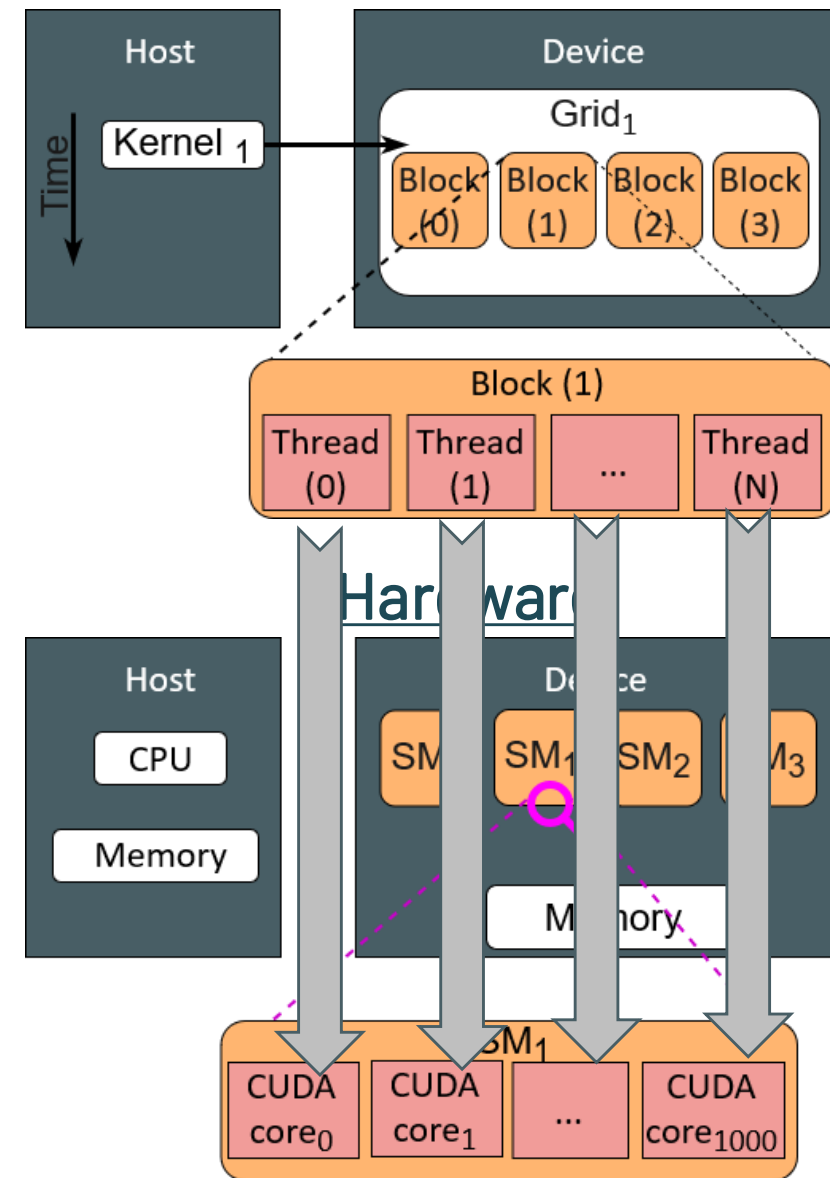- Threads are assigned to CUDA cores
- We have to modify kernel code to use Threads, instead of blocks:

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

- We use the threadId.x instead of blockIdx.x
- In main we have to change the kernel call:
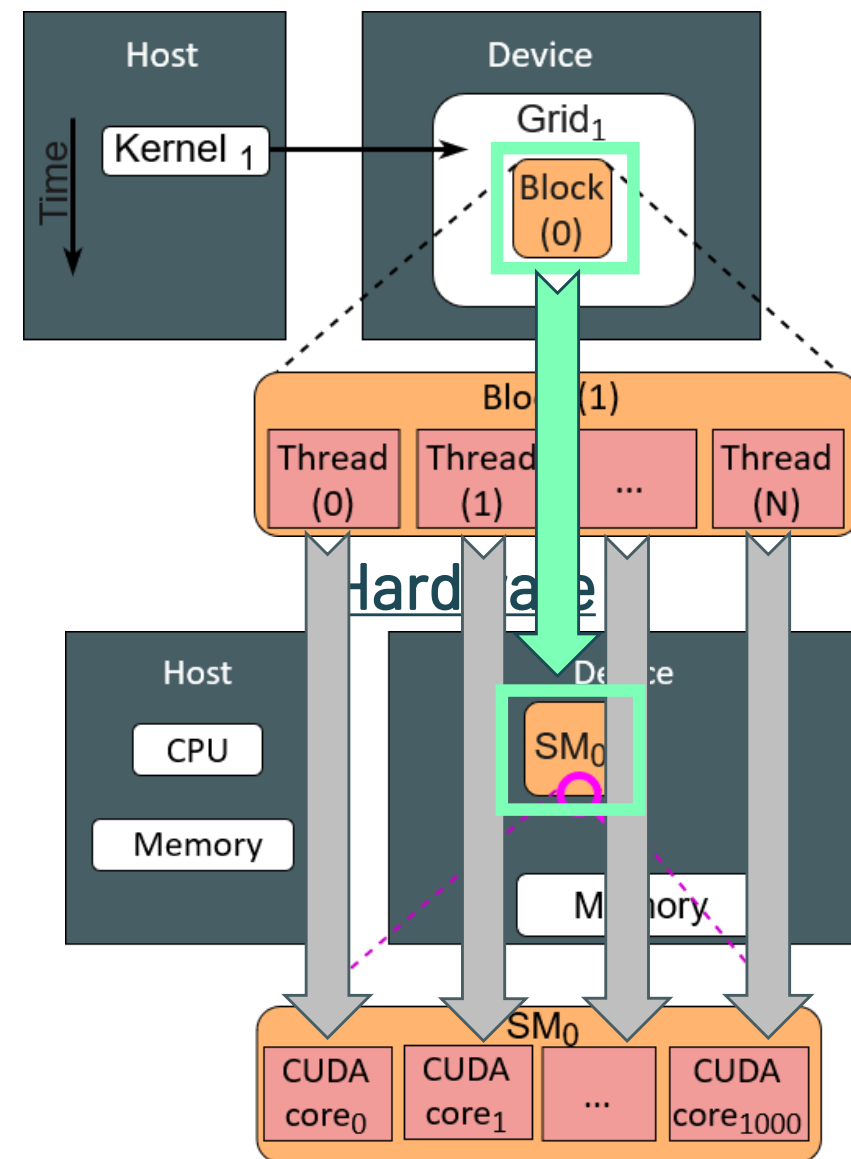
```
add<<<1,N>>>(...);
```

# CUDA Threads → CUDA cores



- A Block can be split into parallel Threads
  - Threads in the same block can cooperate
  - Threads have unique ids (i.e. threadId.x)
- Threads are assigned to CUDA cores
- We have to modify kernel code to use Threads, instead of blocks:

```
__global__ void add(int *a, int *b, int *c){

    c[threadId.x] = a[threadIdx.x] + b[threadIdx.x]; }
```

- We use the threadId.x instead of blockIdx.x
- In main we have to change the kernel call:
```
add<<<1,N>>>(...);
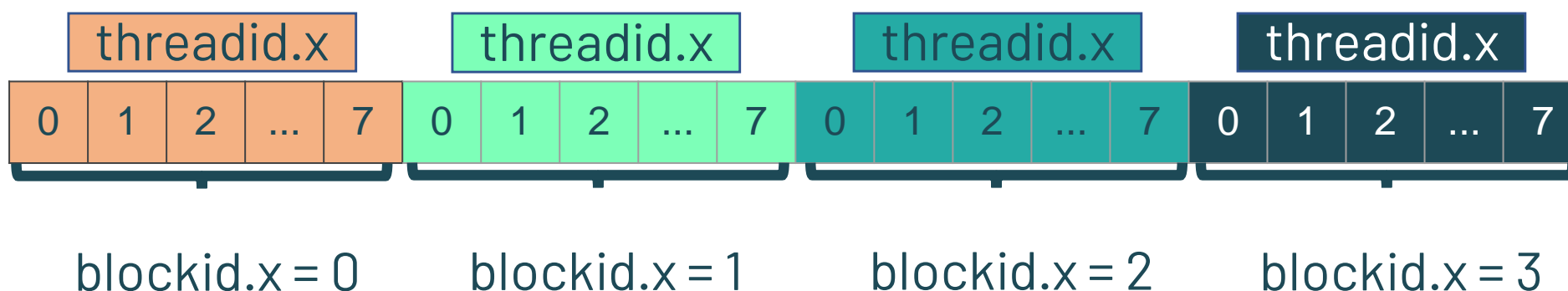```
- Now we call the kernel with 1xBlock and NxThreads

# Combining Threads & Blocks

- Until now we have seen parallel vector addition using
  - Many Blocks with **one Thread** each
  - One **Block** with many Threads

- Now let's get more parallelism (= performance) by using
  - Many Blocks with many Threads

- But before let's discuss data indexing

# Indexing Arrays with Blocks & threads

- Consider indexing an array with one element per thread (8 threads/block)



| threadid.x | threadid.x | threadid.x | threadid.x |

| 0 | 1 | 2 | ... | 7 | 0 | 1 | 2 | ... | 7 | 0 | 1 | 2 | ... | 7 | 0 | 1 | 2 | ... | 7 |

blockid.x = 0          blockid.x = 1          blockid.x = 2          blockid.x = 3

- With **M threads per bloc**k, a unique index for each thread is given by:
  - int index = threadIdx.x + blockIdx.x * M;
  - M = Number of threads per block or block size
  - Use the built-in variable blockDim.x for threads per block
  - int index = threadIdx.x + blockIdx.x * blockDim.x;

# Vector addition with Blocks & Threads

## Host code

```
1  #define THREADS_PER_BLOCK 256
2  const int N=2048;
3
4  int main(void) {
5      int *a, *b, *c; // host copies of a, b, c
6      int *d_a, *d_b, *d_c; // device copies of a, b, c
7      int size = N * sizeof(int);
8
9      // Alloc space for device
10     cudaMalloc((void **)&d_a, size);
11     cudaMalloc((void **)&d_b, size);
12     cudaMalloc((void **)&d_c, size);
13
14     // Alloc spzce on host
15     // space for host copies of a, b, c and setup input values
16     a = (int*)malloc(size); random_ints(a, N);
17     b = (int*)malloc(size); random_ints(b, N);
18     c = (int*)malloc(size);
19
20     // Copy inputs to device
21     cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice);
22     cudaMemcpy(d_b, b, size, cudaMemcpyHostToDevice);
23
24     // Launch add() kernel on GPU
25     add<<<N/THREADS_PER_BLOCK,THREADS_PER_BLOCK>>>(d_a, d_b, d_c);
26
27     // Copy result back to host
28     cudaMemcpy(c, d_c, size, cudaMemcpyDeviceToHost);
29
30     // Cleanup
31     free(a); free(b); free(c);
32     cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
33
34     return 0;
35 }
```

Block size: 256

Array size: 2048

## Kernel code

```
37  __global__ void add(int *a, int *b, int *c) {
38      int index = threadIdx.x + blockIdx.x * blockDim.x;
39          c[index] = a[index] + b[index];
40  }
```

8 x Blocks, 256 x Threads/Block = 2048

add <<< 8, 256 >>>(...);

# Handling arbitrary vector sizes

- Typical problems are not friendly multiples of blockDim.x (e.g. N = 5000)

```
1  #define THREADS_PER_BLOCK 256
2  const int N=5000;                              Array size: 5000
3
4  int main(void) {
5      int *a, *b, *c; // host copies of a, b, c
6      int *d_a, *d_b, *d_c; // device copies of a, b, c
7      int size = N * sizeof(int);
8
9      ...
10
11     // Launch add() kernel on GPU
12     add<<<N/THREADS_PER_BLOCK + 1,THREADS_PER_BLOCK>>>    Update kernel launch
13
14     ...
15
16     return 0;
17 }
18
19 __global__ void add(int *a, int *b, int *c) {
20     int index = threadIdx.x + blockIdx.x * blockDim.x;
21     if (index < N)                                        Avoid accessing beyond the end of arrays
22         c[index] = a[index] + b[index];
23 }
```

# Compile a CUDA program

- nvcc is the CUDA compiler

- With CUDA both Host and Device code can be in the same file
  - With suffix ".cu"

- nvcc separates source code into host and device components
  - Device functions (e.g. add()) processed by NVIDIA compiler
  - Host functions (e.g. main()) processed by standard host compiler (g++/gcc)

- Compile vector addition example:
  - nvcc vectorAdd.cu –o vectorAdd

- Run:
  - ./vectorAdd

# Thank you

# Questions?

Manos Pavlidakis
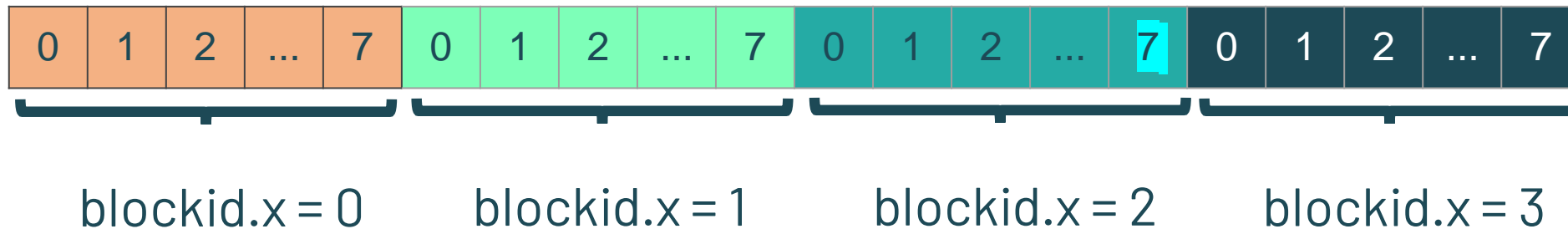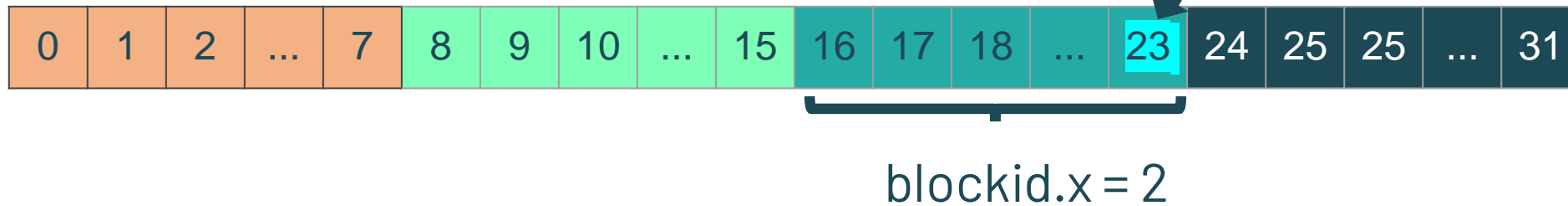manospavl@ics.forth.gr

Find the example in: https://github.com/manospavlidakis/VectorAdditionCUDA.git

# A simple example

- Which thread element will operate on the blue cell?

| 0 | 1 | 2 | ... | 7 | 8 | 9 | 10 | ... | 15 | 16 | 17 | 18 | ... | 23 | 24 | 25 | 25 | ... | 31 |
|---|---|---|-----|---|---|---|----|-----|----|----|----|----|-----|----|----|----|----|-----|----|

# A simple example

- Which thread element will operate on the blue cell?

| 0 | 1 | 2 | ... | 7 | 8 | 9 | 10 | ... | 15 | 16 | 17 | 18 | ... | 23 | 24 | 25 | 25 | ... | 31 |

threadid.x = 7

blockid.x = 2

| 0 | 1 | 2 | ... | 7 | 0 | 1 | 2 | ... | 7 | 0 | 1 | 2 | ... | 7 | 0 | 1 | 2 | ... | 7 |

blockid.x = 0      blockid.x = 1      blockid.x = 2      blockid.x = 3

- With **M threads per block = 8**
  - int index = threadIdx.x + blockIdx.x * M = 7 + 2 * 8 = 23;