# Introduction to accelerated computing

**14th Inverted CERN School of Computing**

**Charis Kleio Koraka**

Tuesday March 7th 2023

# Overview

- Hardware accelerators and heterogeneous computing
- The GPU
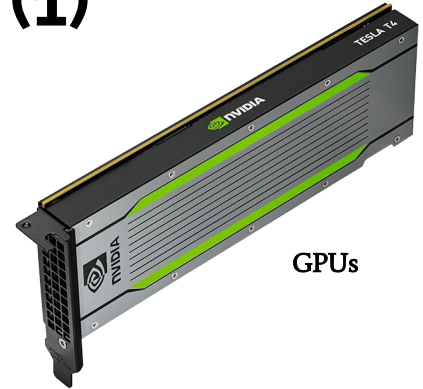- GPU applications in HEP
- The CUDA programming  model

# Hardware accelerators

- Devices built for **executing specific tasks more efficiently** compared to running on the standard computing architecture of a CPU

- Come in many flavors :
    - GPUs / FPGAs / TPUs …

- Part of our everyday lives :
    - Encryption, video stream decoding, 3D graphics acceleration, pattern/object recognition, machine learning, AI and many more

# Some types of hardware accelerators (1)

- **GPU** (Graphic Processing Unit)
  - Initially developed for graphics processing
  - Optimized for parallel processing of floating-point operations & used in a variety of tasks

GPUs

- **FPGA** (Field-Programmable Gate Array)
  - Integrated circuit (IC) configurable by the user and provides interface flexibility
  - FPGAs can be reprogrammed to suit the needs of the application or required functionality
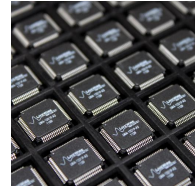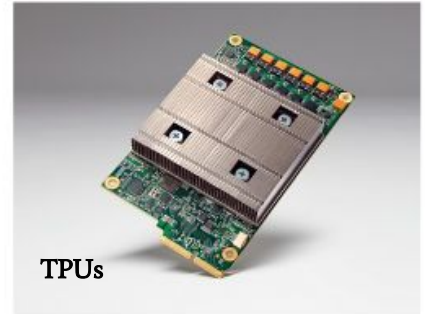
FPGAs

# Some types of hardware accelerators (2)

- **ASIC** (Application-Specific Integrated Circuit)
  - IC chip customized for a particular use

  - i.e. lower precision and/or optimised memory usage to maximize throughput



ASIC

- **TPU** (Tensor Processing Unit)
  - Optimised to perform matrix-multiplication operations / used in e.g. NN and RF training



TPUs

- **VPU** (Vision Processing Unit)
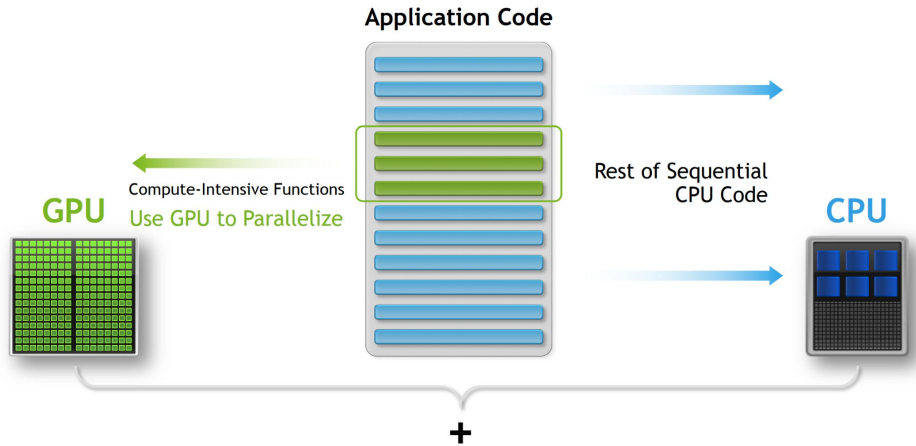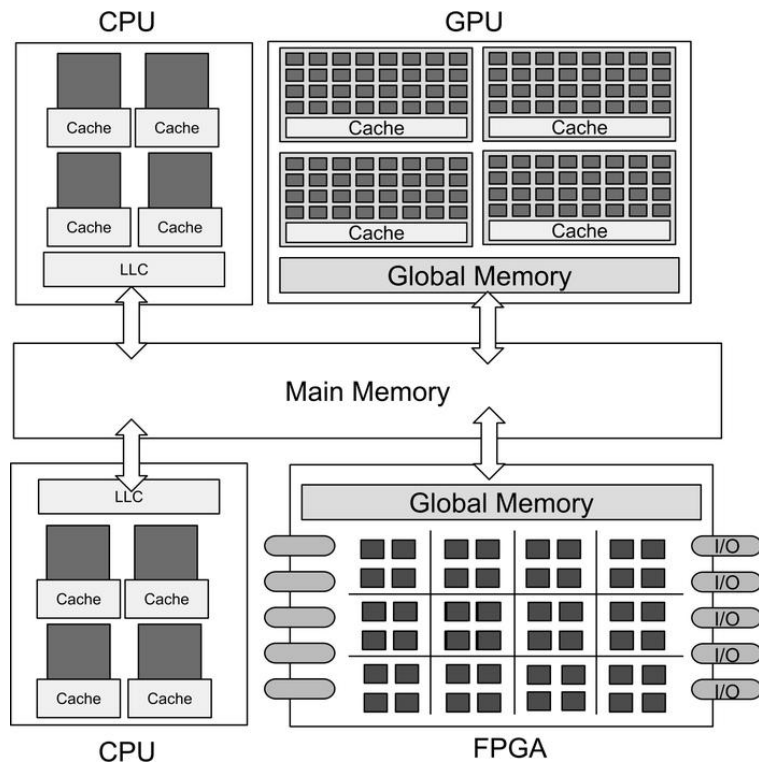  - Used to accelerate machine vision algorithms, i.e. CNNs , AI etc.



VPUs

# How are hardware accelerators used?



**Application Code**

Compute-Intensive Functions
Use GPU to Parallelize

GPU

Rest of Sequential
CPU Code

CPU

+

- In accelerated computing we take the compute intensive parts of the application code and parallelize that for execution on e.g. a GPU
  - Typically integer or floating-point mathematical operations

- The remainder of the code (usually the vast majority) remains on the CPU
  - The part of code that remains on the CPU is ideally serial code

- Data between the CPU and the accelerator has to be transferred:

  - This is performed via interconnect i.e. PCIe (*Peripheral Component Interconnect Express*), NVLink, Ethernet etc.

# Heterogeneous computing



- Heterogeneous computing involves using multiple different types of processors to accomplish a task

- Code can run on more than one platform concurrently

- A heterogeneous system can consist of :
  - Different types of CPUs (i.e. combine compute powerful with less compute powerful but more power efficient CPU cores)
  - Hardware accelerators

# The GPU

# The Graphic Processing Unit (GPU)
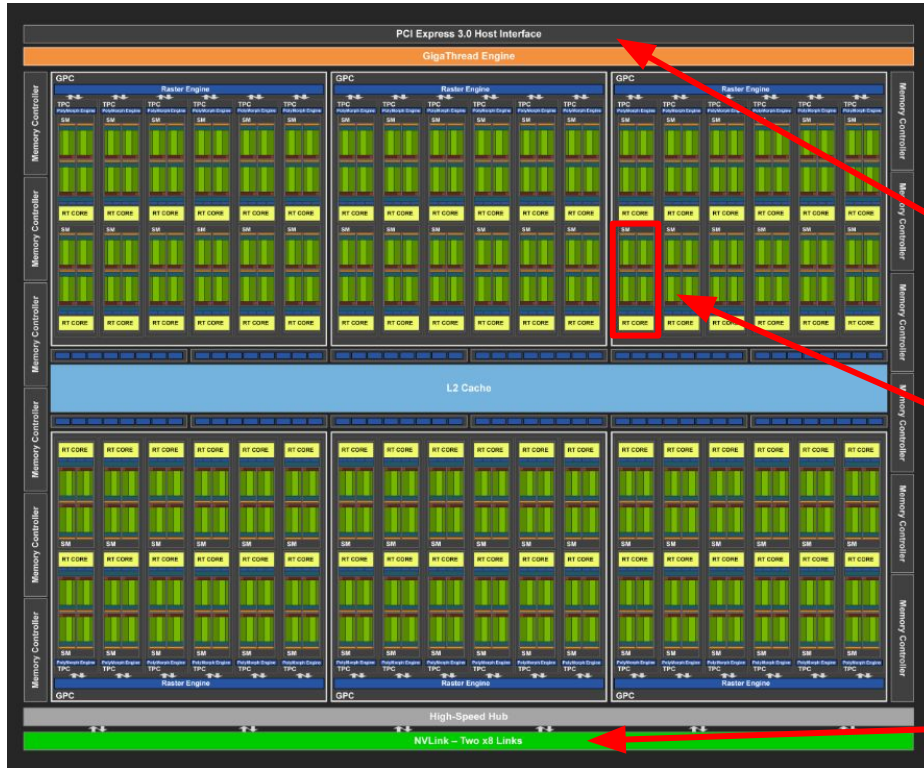
**GPUs are similar to CPUs :**

- Silicon based micro-processor that contain cores, registers, memory, and other components.

**But also very different :**

- **Many-core processor**
- Follows the **Single instruction, multiple threads (SIMT)** execution model
  - Asynchronous programming model where threads are not executed in lockstep
- GPU acceleration emphasizes on :
  - **High data throughput and massive parallel computing:** a GPU consist of hundreds of cores performing the same operation on multiple data items in parallel.
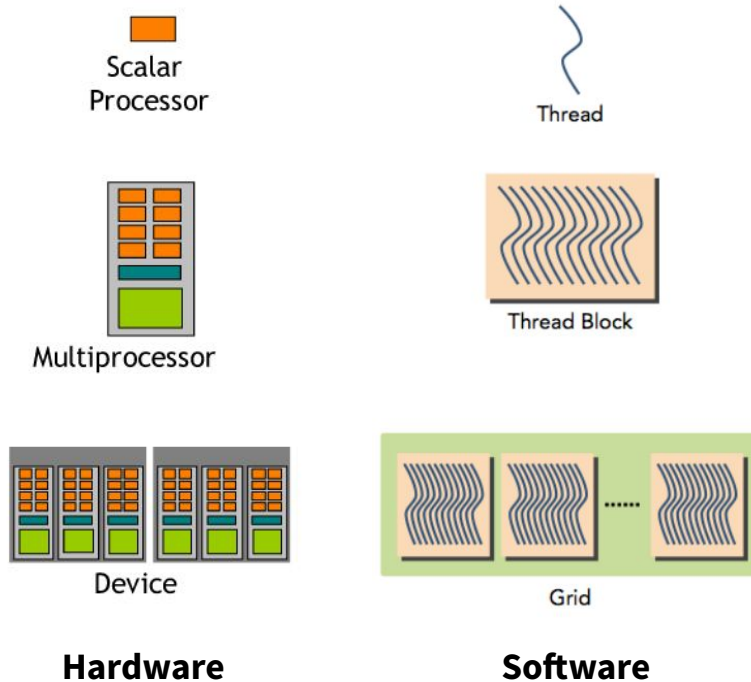
# The NVidia GPU architecture



- The GPU architecture is built around a scalable array of Streaming Multiprocessors (SM).
- Each SM in a GPU is designed to support concurrent execution of hundreds of threads

PCIe interconnect: Can be used for connecting GPU to host CPU

SM

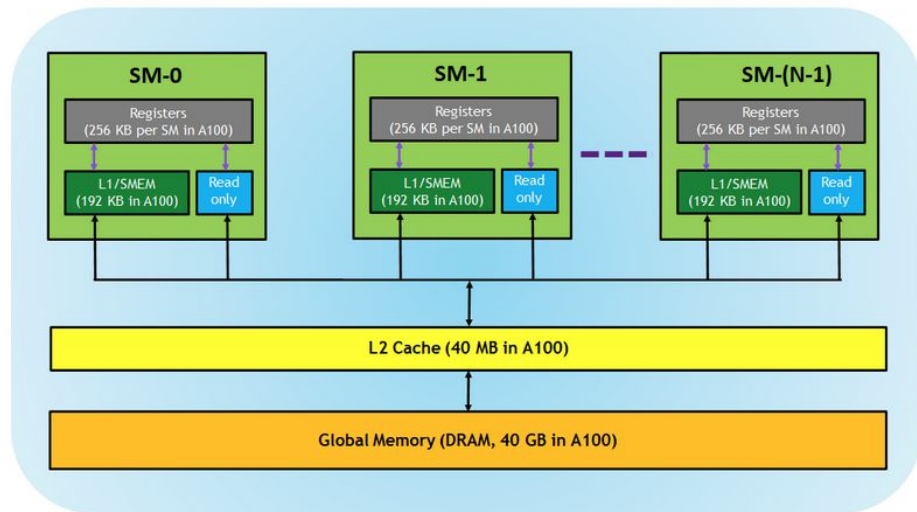NVlink : Can be used to connect to additional GPUs

# Hardware to software mapping



**Scalar Processor**

**Thread**

**Multiprocessor**

**Thread Block**

**Device**

**Grid**

**Hardware**                    **Software**

- A scalar processor or CUDA core is equivalent to a software thread
- Scalar processors are grouped into a SM
- Each execution of a GPU function is done concurrently on a number of threads referred to as a thread block
- Each thread block is executed by one SM and cannot be migrated to other SMs in GPU
- The set of thread blocks executing the GPU function is called a grid.
- In CUDA terminology the GPU is referred to as the device
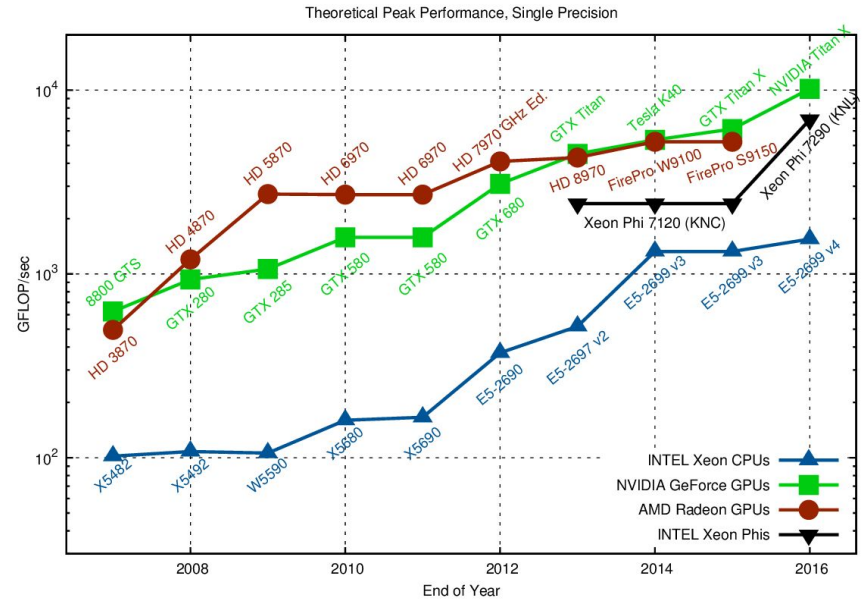
# GPU memory hierarchy

- **Registers**
  - Memory private to each thread
  - Fastest form of memory
- **L1 cache/Shared memory**
  - Fast accessible memory that can be accessed by threads in the same block and threads of different blocks in the same SM
- **Read-only**
  - Each SM has a constant/texture cache memory which is read-only to kernel code. Fast but limited in size
- **L2 Cache**
  - Memory that all threads in all blocks can access. Fast but limited in size.
- **Global memory**
  - GPUs DRAM memory
  - Slow but large

# Performance comparison of CPUs and GPUs (1)

FLOPS : Floating-Point Operations per Second

- Measure of computing performance useful in fields that require floating-point calculations (such as HEP)

- GPUs can deliver more FLOPS compared to CPUs
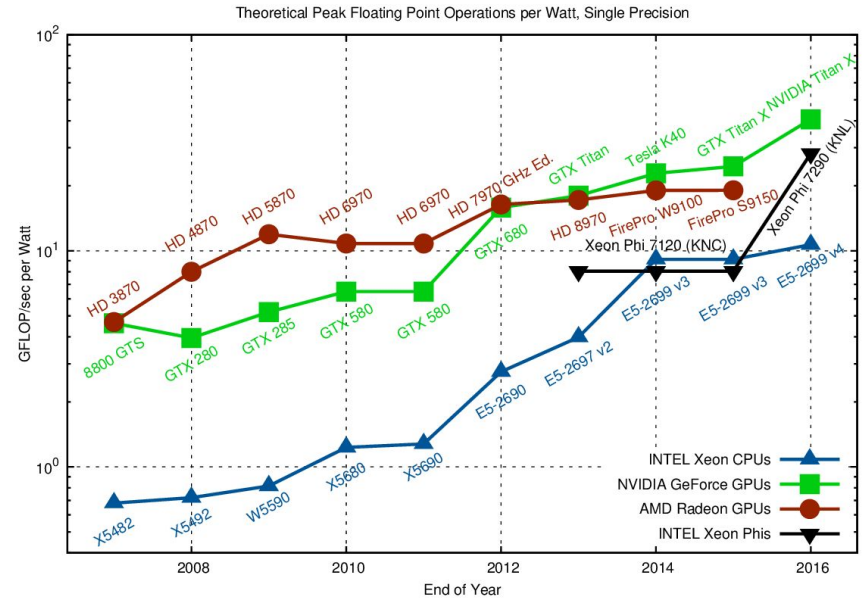


Theoretical Peak Performance, Single Precision

# Performance comparison of CPUs and GPUs (2)

FLOPS per Watt :

- Rate of floating-point operations performed per watt of energy consumed

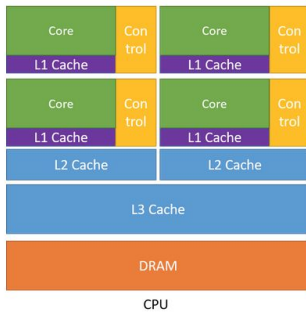Important since power consumption is limiting factor in hardware manufacturing/usage:

- Peak performance constrained by the amount of power it can draw and the amount of heat it can dissipate



Theoretical Peak Floating Point Operations per Watt, Single Precision

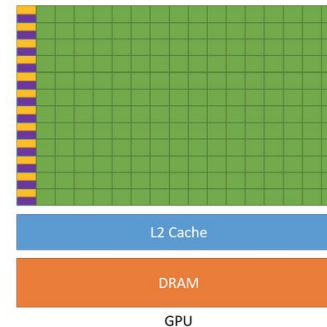# CPU vs GPU - overview of main differences

**CPU**

- ~O(10) powerful cores
  - Larger instruction set
- Low latency
- Serial processing
- Complex operations
- Higher clock speeds

**GPUs**

- ~O(1000) of less powerful cores
  - Smaller instruction set
- High throughput
- Parallel processing
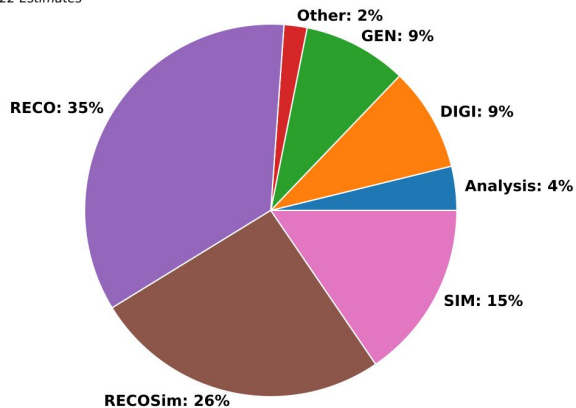- Simple operations
- Better per-watt performance

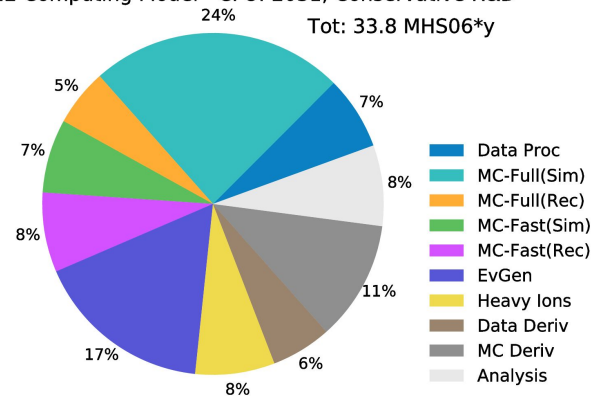# Applications of GPUs in HEP

# Computing needs in HEP

- Event generation

- Simulation

- Event reconstruction

- Event post-processing

- Data analysis

[link]

**CMS** *Public*
Total CPU HL-LHC (2031/No R&D Improvements) fractions
*2022 Estimates*

- Other: 2%
- GEN: 9%
- DIGI: 9%
- Analysis: 4%
- SIM: 15%
- RECOSim: 26%
- RECO: 35%

[link]

**ATLAS** *Preliminary*
2022 Computing Model - CPU: 2031, Conservative R&D
Tot: 33.8 MHS06*y

- 24%
- 5%
- 7%
- 7%
- 8%
- 8%
- 17%
- 8%
- 6%
- 11%
- 8%

- Data Proc
- MC-Full(Sim)
- MC-Full(Rec)
- MC-Fast(Sim)
- MC-Fast(Rec)
- EvGen
- Heavy Ions
- Data Deriv
- MC Deriv
- Analysis
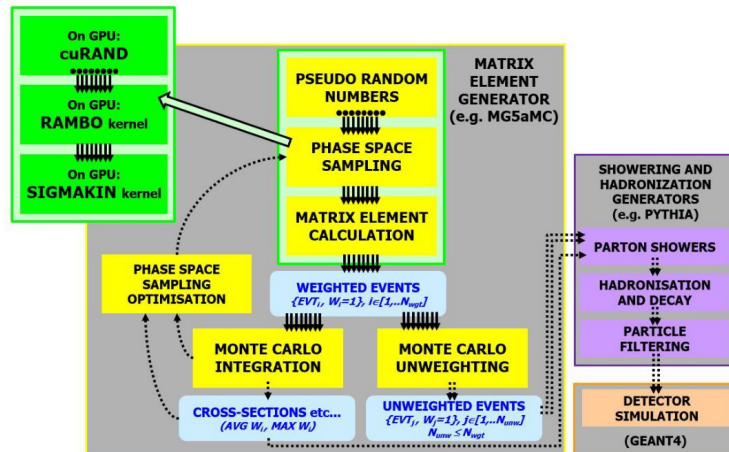
# How can GPUs help?

- Event generation

- Simulation

- Event reconstruction

- Event post-processing

- Data analysis

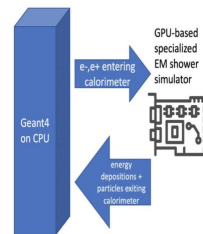- **GPU enabled event generator i.e. Madgraph** [i]

# How can GPUs help?
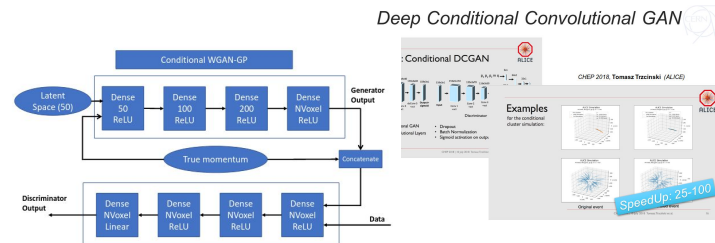
- Event generation

- Simulation

- Event reconstruction

- Event post-processing

- Data analysis

- **GPU based Geant4 application (i.e. AdePT) [i]**



- **AI/ML enabled Fast Simulation (i.e. AltFast3 in ATLAS [ii] , DC-GAN in ALICE [iiii])**

*Deep Conditional Convolutional GAN*

# How can GPUs help?

- Event generation

- Simulation

- Event reconstruction
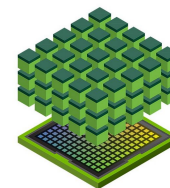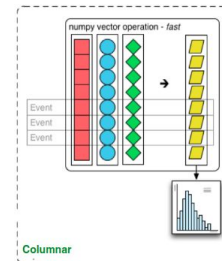
- Event post-processing

- Data analysis

- **Track reconstruction, primary vertex reconstruction, raw data unpacking, clustering etc.**

- **Various efforts in different experiments (Patatrack track reconstruction [i], Allen project [ii], ALICE TPC track reconstruction [iii] etc.)**

# How can GPUs help?

- Event generation

- Simulation

- Event reconstruction

- Event post-processing

- Data analysis

- **Training and inference of ML models**
- **Perform HEP analysis using columinar analysis paradigm tools (i.e. coffea [i])**

# Introduction to CUDA

# The CUDA programming model

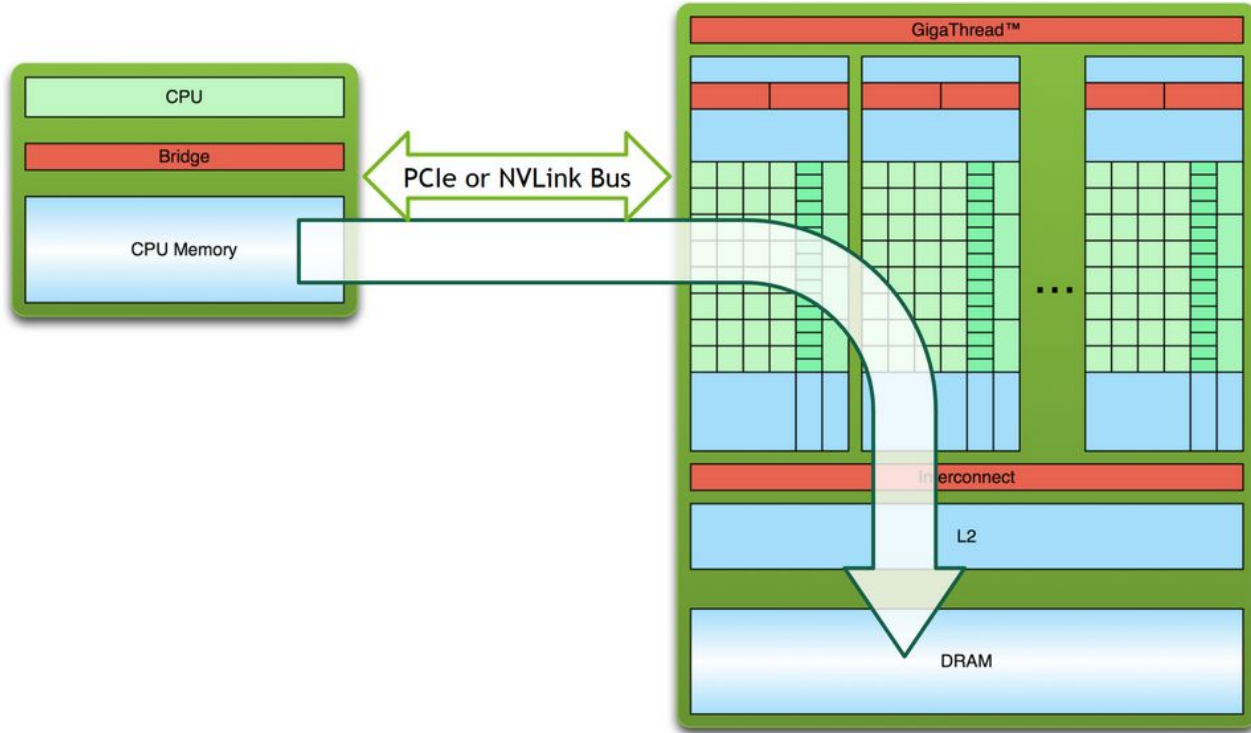**CUDA** → **C**ompute **U**nified **D**evice **A**rchitecture.

- It is an extension of C/C++ programming
- Developed by Nvidia and is used to develop applications executed on NVidia GPUs

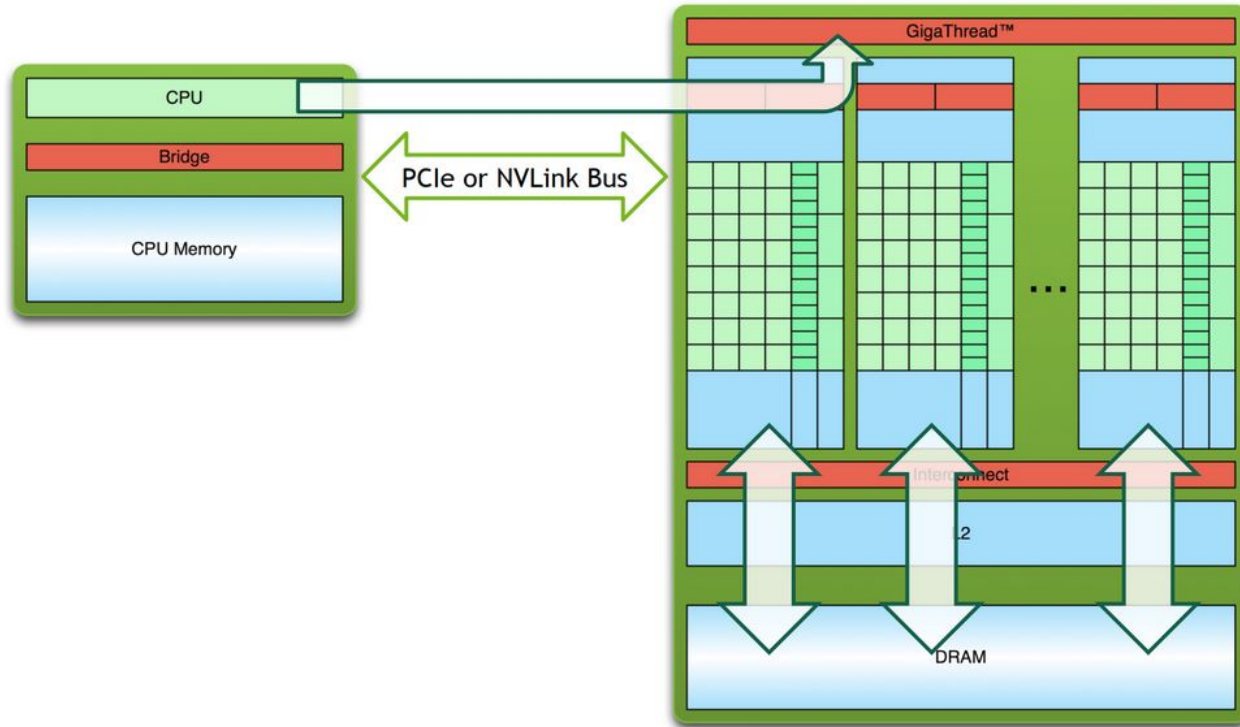To execute any CUDA program, there are three main steps:

- Copy the input data from CPU or host memory to the device memory
- Execute the CUDA program
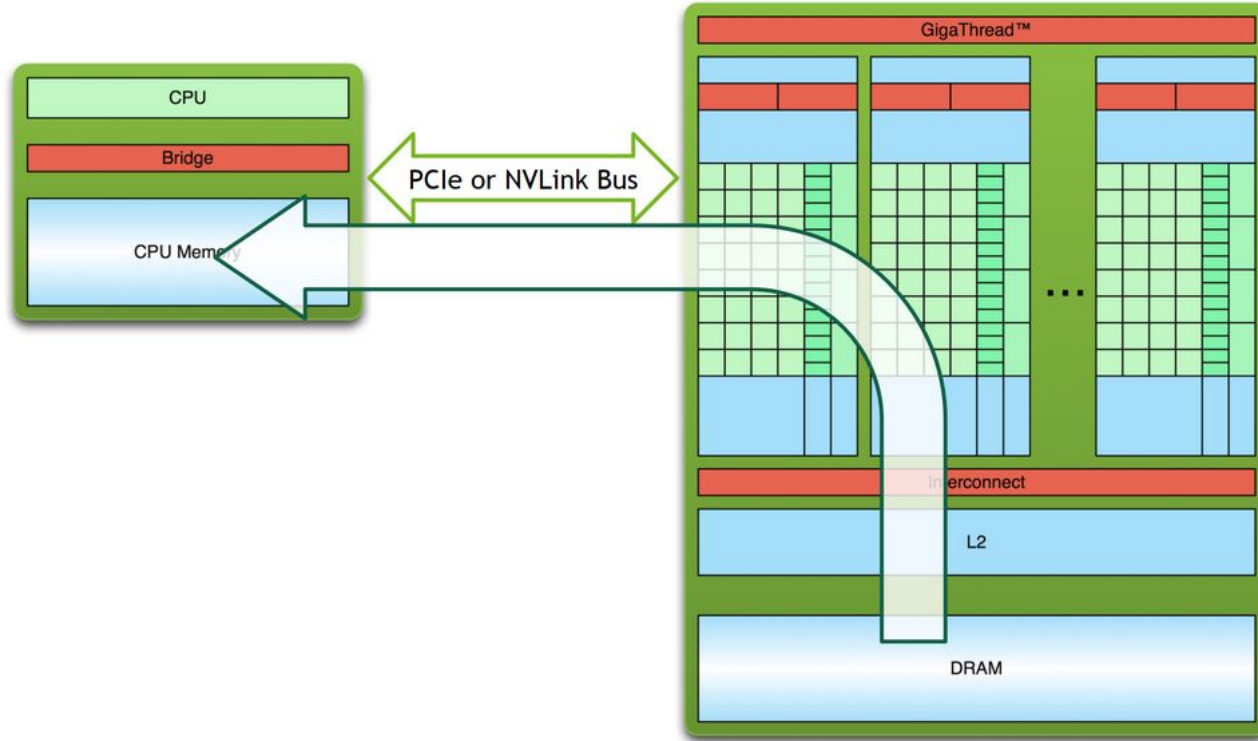- Copy the results from device memory to host memory

# 1. Copy data for host to device
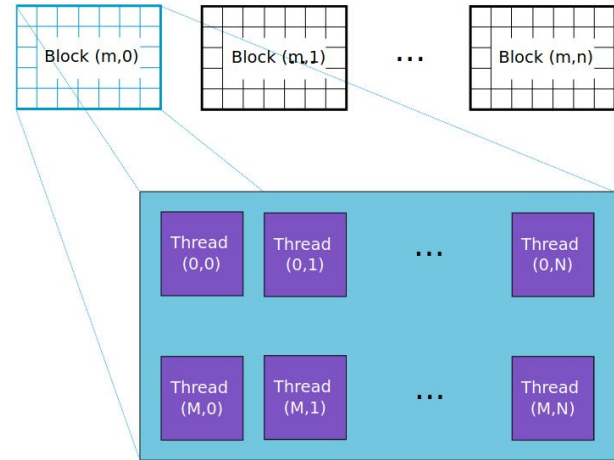
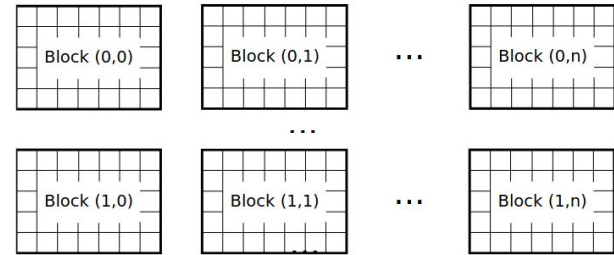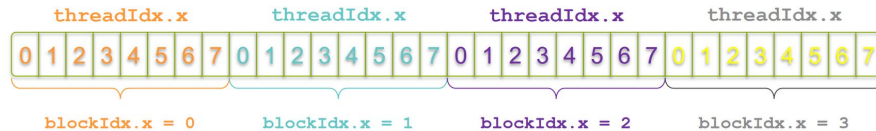# 2. Execute the CUDA program
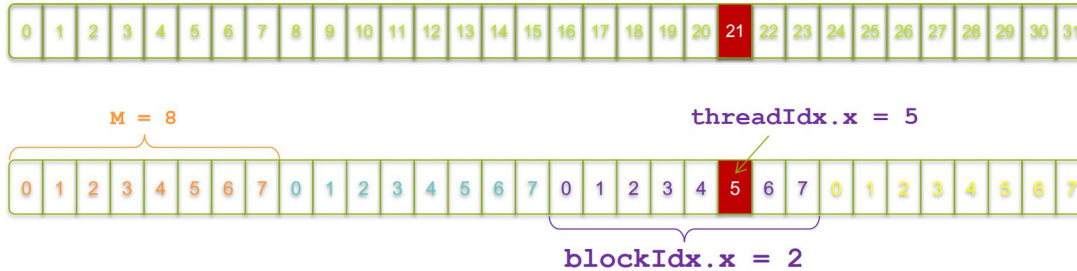
# 3. Copy data from device back to host

# Threads & blocks

- In CUDA, built-in variables are available in order to express threads and blocks :
  - `threadIdx & blockIdx`

- The variables have 3-dimensional indexing & provide a natural way to express elements in vectors and matrices :
  - `threadIdx.x , threadIdx.y threadIdx.z`

- CUDA architecture limits the numbers of threads per block (1024 threads per block limit).
- The dimension of the thread block is accessible within the kernel through the built-in `blockDim` variable.
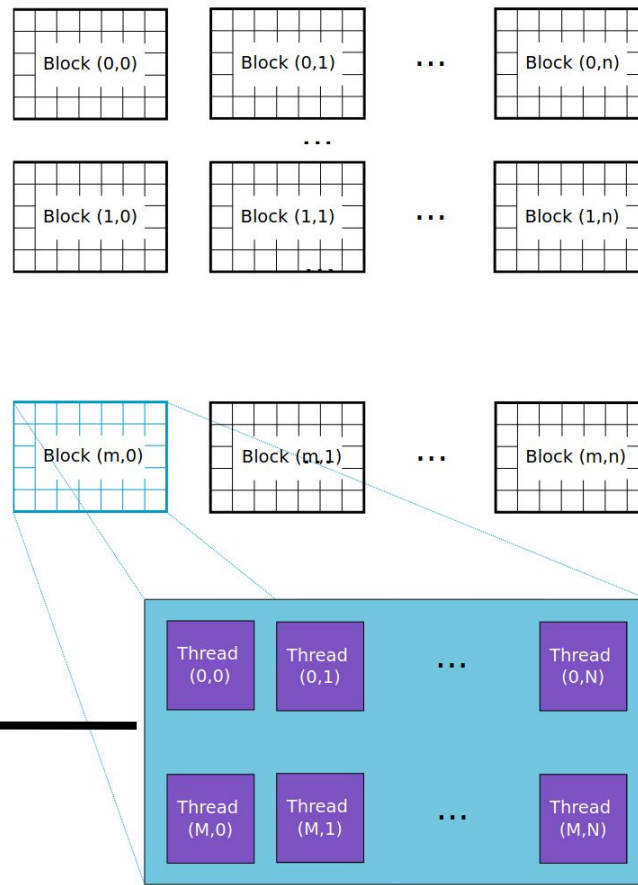
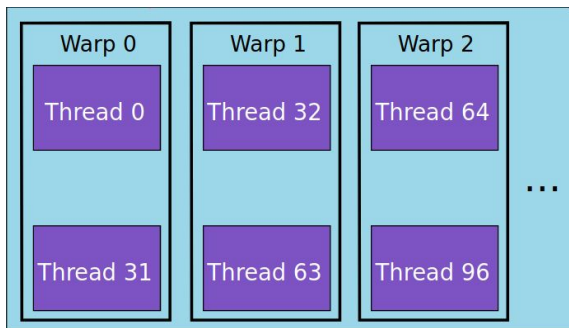# Indexing using blockIdx and threadIdx

- The `threadIdx` & `blockIdx` variables can be used to express the unique index of an element in an array/matrix etc.
- Assuming that each block consists of a number of M threads :
  - `index = threadIdx.x + blockIdx.x * M;`



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

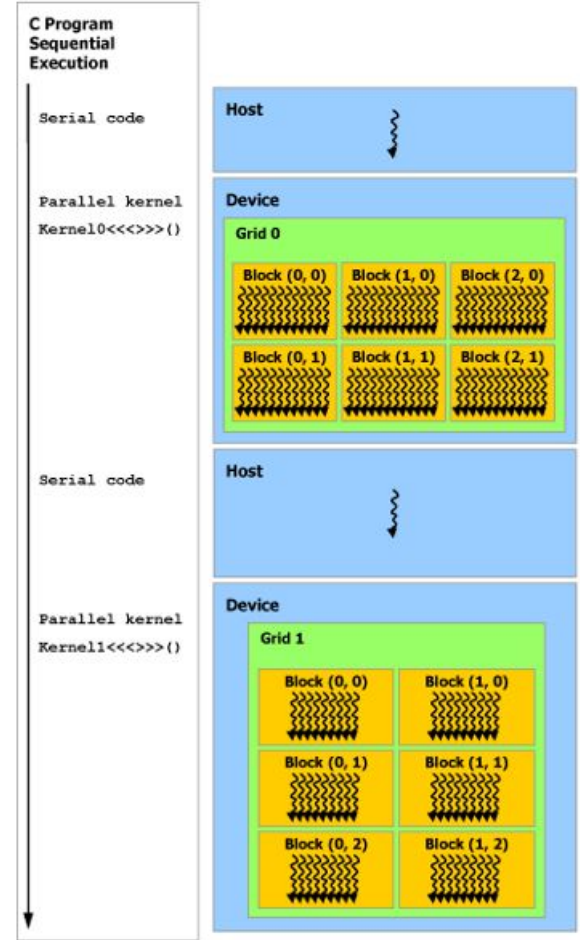# Warps

- Within a thread block, threads are executed in groups → **Warps**
- A warp is an entity of 32 threads on Nvidia GPUs
- If the block size is not divisible by 32, some of the threads in the last warp will remain idle :
  - block size should be chosen to be a multiple of the warp size
- Threads in the same warp are processed simultaneously

# CUDA kernel

- CUDA kernel is a function that gets executed on the GPU
- The kernel expresses the portion of the application that is parallelizable
  - It will be executed multiple times in parallel by different CUDA threads

# CUDA function declarations

| Declaration | Callable from: | Executed on: |
|---|---|---|
| _ _global_ _ | host | device |
| _ _device_ _ | device | device |
| _ _host_ _ | host | host |

- _ _**global**_ _ keyword defines a kernel function:
  - Is launched by host and executed on the device
  - Must return void
- _ _**device**_ _ and _ _**host**_ _ can be used together
- _ _**host**_ _ declaration, if used alone, can be omitted

# Launching a CUDA kernel

- Let's assume we have the following kernel :

```
__global__ void mykernel() {
    …Do something…
}
```

**This is the grid dimension i.e. the number of blocks that will be launched**

**This is the block dimension i.e. the number of threads within a block**

- How do we launch it?

```
myKernel<<<nBlocks,nThreads>>>();
```

- The above command will launch the kernel with **nBlocks**, each of which has **nThreads**.
- The kernel is executed multiple times concurrently by different threads
- The total number of invocations of the kernel body is now **nBlocks** * **nThreads**.

# Memory management

- The host and device have their own separate memory:
    - Device pointers point to GPU memory
    - Host pointers point to CPU memory
- CUDA kernels operate out of device memory
- CUDA provides functions to allocate device memory, release device memory, and transfer data between the host memory and device memory :

```
cudaMalloc(&ptr, size_in_bytes_to_allocate)

cudaFree(ptr)

cudaMemcpy(destination_ptr,source_ptr, size_in_bytes, direction)
```
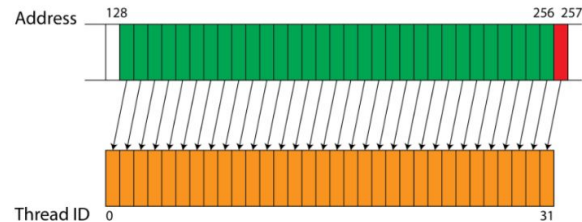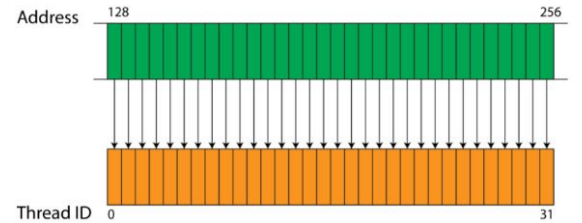
where direction can be :

- For copying data from CPU to GPU
- For copying data from GPU to CPU

```cpp
int* a;
int* d_a;
// Host copy of variable a
a = (int*) malloc(sizeof(int));
// Device copy of variable a
cudaMalloc(&d_a, sizeof(int));
// Set the host value of a
*a = 1;
// Copy the value of a to the device
cudaMemcpy(d_a, a, sizeof(int), cudaMemcpyHostToDevice);
// Launch the kernel to set the value
do_something<<<1,1>>>(d_a);
cudaDeviceSynchronize();
// Copy the value of a back to the host
cudaMemcpy(a, d_a, sizeof(int), cudaMemcpyDeviceToHost);
// Free the allocated memory
free(a);
cudaFree(d_a);
```

# Coalesced global memory access

- Global memory loads and stores data in as few as possible transactions → coalesced memory access

- Important performance consideration as it can affect the time needed to access data

- Every successive 128 bytes (DRAM burst) can be accessed by a warp

- If the data accessed by the threads in a warp are not in the same burst section, the data access will take twice as long

# Putting together a CUDA program

The main components of a CUDA program are :

- Declarations of functions :
  - These can be __host__ / __global__ / __device__ functions

- Copying data to/from host :
  - Use cudaMalloc / cudaMemcpy / cudaFree

- Kernel launch <<<grid size, block size >>>(<arguments>)

- Concurrency management
  - Use __syncthreads() or CudaDeviceSynchronize()

```
__global__ void do_something (int* a) {
    *a = 2;
}

int main() {

    int* a;
    int* d_a;
    // Host copy of variable a
    a = (int*) malloc(sizeof(int));
    // Device copy of variable a
    cudaMalloc(&d_a, sizeof(int));
    // Set the host value of a
    *a = 1;
    // Copy the value of a to the device
    cudaMemcpy(d_a, a, sizeof(int), cudaMemcpyHostToDevice);
    // Launch the kernel to set the value
    do_something<<<1,1>>>(d_a);
    cudaDeviceSynchronize();
    // Copy the value of a back to the host
    cudaMemcpy(a, d_a, sizeof(int), cudaMemcpyDeviceToHost);
    // Free the allocated memory
    free(a);
    cudaFree(d_a);

}
```

# Error handling

- Error codes can be converted to a human-readable error messages with the following CUDA run- time function:

```
char* cudaGetErrorString(cudaError_t error)
```

- A common practice is to wrap CUDA calls in utility functions that manage the error returned :

```
int* a;
// Illegal: cannot allocate a negative number of bytes
cudaError_t err = cudaMalloc(&a, -1);
if (err != cudaSuccess) {
    printf("CUDA error %s\n", cudaGetErrorString(err));
    exit(-1);
}
```

- To detect errors in a kernel launch, we can use the API call **cudaGetLastError()** which returns the error code for whatever the last CUDA API call was.

```
cudaError_t err = cudaGetLastError();
```

- For errors that occurs asynchronously during the kernel launch, **cudaDeviceSynchronize()** has to be invoked after the kernel in order to return any errors associated with the kernel launch.
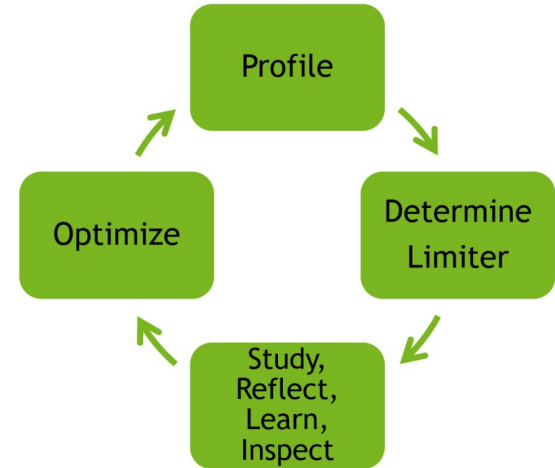
# Compilation

- Compiling a CUDA program is similar to compiling a  C/C++ program.
- Cuda code should be typically stored in a file with extension .cu
- NVIDIA provides a CUDA compiler called **nvcc** :
  - nvcc is called for CUDA parts
  - gcc is called for c++ parts
  - nvcc converts .cu files into C++ for the host system and CUDA assembly or binary instructions for the device
- Usage :

```
nvcc myCudaProgram.cu -o myCudaProgram
```

# Profiling (1)

- Similarly to CPU code, GPU code can be profiled

- Goal of profiling is to identify and optimise performance limiters.

- Common reasons for limited performance include :
    - Portions of the code that run serially on the CPU
    - Memory copies for host to device
    - Latency of launching GPU kernels
    - Uncoalesced memory accesses, lack of cache reuse, not using shared memory, register spilling etc.
    - Low arithmetic intensity (operations computed per byte accessed from memory)

# Profiling (2)

- Many profiling tools exist. Some commonly used ones include
- **nsys** :
  - Command line profiler for CUDA applications
  - Results can be saved for later viewing by the Visual Profiler.
- **nvvp / ncu** :
  - Nvidia Visual Profiler, Nsight Compute
  - Interactive kernel profilers for CUDA applications.
  - Provide detailed performance metrics and API debugging via a user interface and command line tool.
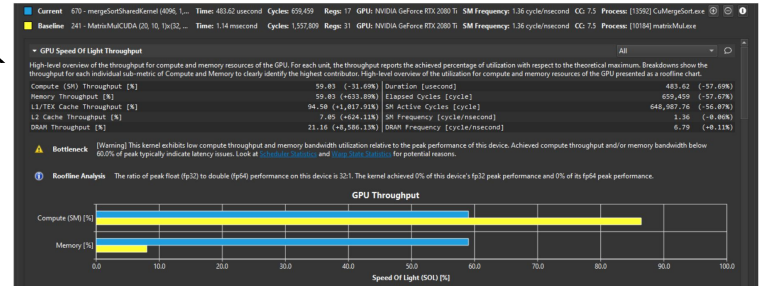
**Image source [6]**

# Wrapping-up

# Summary

- Hardware accelerators are a part of everyday life and are used in heterogeneous computing systems
- GPUs emphasize on high data throughput and massive parallel computing
- GPUs have made their way into HEP and are used for many applications
- The CUDA programming model :
  - Extension of C/C++ programming developed by Nvidia and used for applications executed on Nvidia GPUs
  - CPU and GPU system are referred to as host and device respectively.
    - The host and device have their own separate memory
  - Typically, we run serial workload on the CPU and offload parallel computation to the GPUs
    - CUDA threads are used to execute work in parallel
  - Basic CUDA syntax:
    - __global__ function declaration (kernel) is called from the host and executed on the device
    - Memory management can be performed using cudaMalloc(), cudaFree() & cudaMemcpy()
    - To launch a CUDA kernel with N blocks and M threads/block syntax is <<<N,M>>>()

# Back-up

# Resources

1. NVIDIA Deep Learning Institute material [link](#)
2. 10th Thematic CERN School of Computing material [link](#)
3. Nvidia turing architecture white paper [link](#)
4. CUDA programming guide [link](#)
5. CUDA runtime API documentation [link](#)
6. CUDA profiler user's guide [link](#)

# Thread synchronization

- A kernel call is asynchronous with respect to the host thread :
    - After a kernel is invoked, the program returns to the host side and continues execution.
- There are two levels of synchronization
    - **Block level**
    - **Grid level**
- To synchronize threads within one block :
    - Call `__syncthreads()` within the kernel code
- To synchronize threads at grid level
    - Call to `CudaDeviceSynchronize()` from host code.
    - Program waits until all work launched on the device has finished.

```
__global__ void myKernel () {
    for (int i = threadIdx.x; i < N; i++) {
        Fill variable[threadIdx.x]
    }
    __syncthreads();
    for (int i = threadIdx.x; i < N; i++) {
        Use variable[threadIdx.x]
    }
}
```

```
int* a;
int* d_a;
// Host copy of variable a
a = (int*) malloc(sizeof(int));
// Device copy of variable a
cudaMalloc(&d_a, sizeof(int));
// Set the host value of a
*a = 1;
// Copy the value of a to the device
cudaMemcpy(d_a, a, sizeof(int), cudaMemcpyHostToDevice);
// Launch the kernel to set the value
do_something<<<1,1>>>(d_a);
cudaDeviceSynchronize();
// Copy the value of a back to the host
cudaMemcpy(a, d_a, sizeof(int), cudaMemcpyDeviceToHost);
// Free the allocated memory
free(a);
cudaFree(d_a);
```