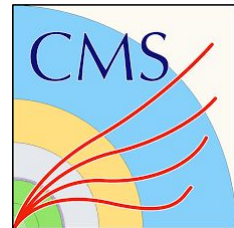


Introduction to GPU programming



WISCONSIN
UNIVERSITY OF WISCONSIN-MADISON

HSF-India HEP Software Workshop

Charis Kleio Koraka

Tuesday December 18th-22nd 2023

Who am I?

- Grew up in Athens, Greece
- Joined the CMS experiment in 2016 as a master student, stayed with CMS ever since
- Did my PhD in the University of Athens on the measurement of the ttH associated production

But what do I do right now?

- Postdoc at University of Wisconsin-Madison
- Interested in searches for very heavy fermions like vector-like leptons and quarks
- Working with cms offline software trying to convince people that software should be written with GPUs and parallel computing in mind :D

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

Central processing unit (CPU)

Silicon-based micro-processor

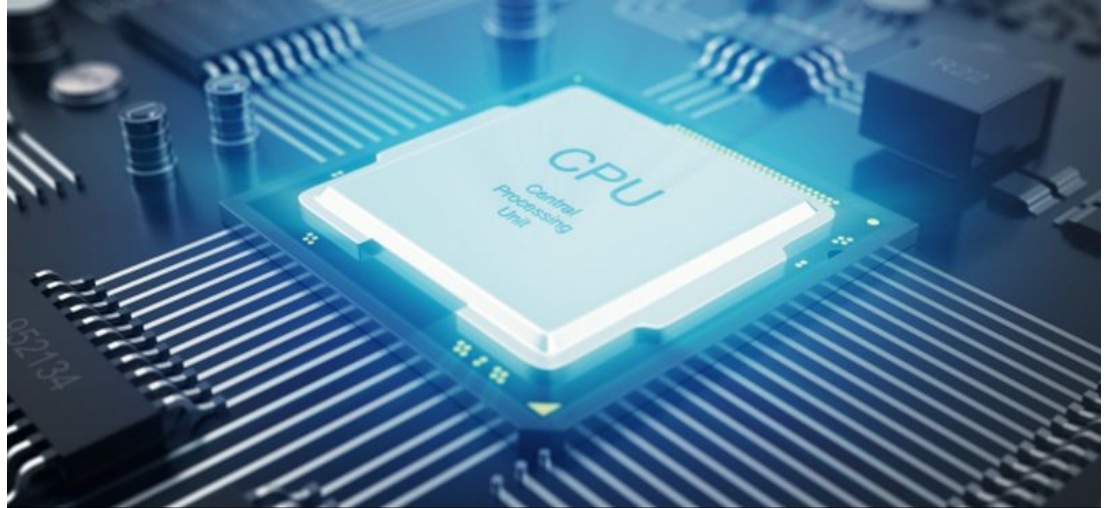
Used in most of our computers since it can handle a variety of tasks.

Performs certain types of operations **serially** :

- Arithmetic (+,*)
- Logical functions (AND, OR, NOT)
- Input/Output (I/O) operation

Is able to execute a sequence of instructions, which constitutes the “program”

The CPU is the brain of our computer, that reads information, performs calculations and moves it where it needs to go

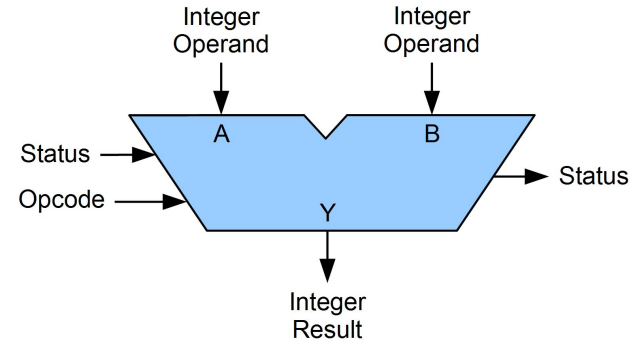


How does a CPU work ? (1)

Principal components of a CPU :

- Arithmetic Logic Unit (ALU) :
 - Used to perform arithmetic and logic operations on integer binary numbers
- Processor registers :
 - A quickly accessible location available to a computer's processor
 - Is used to supply operands to the ALU and store the results of the ALU operations
- Control Unit (CU)
 - Is in charge of orchestrating fetching from memory / decoding / execution of instructions etc.

* Image taken from [\[1\]](#)



* Schematic representation of an ALU

How does a CPU work ? (2)

CPUs are implemented on integrated circuit (IC) microprocessors :

- A single IC chip can have one or more CPU cores
- Microprocessor chips with multiple CPUs are **multi-core processors**
- Processor cores can also be multithreaded to create additional virtual CPUs

Schematic representation of principal components that form a CPU

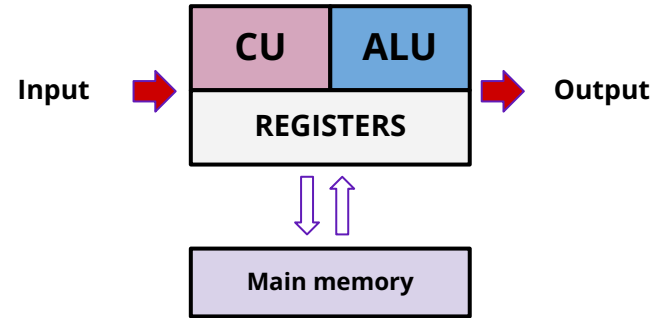
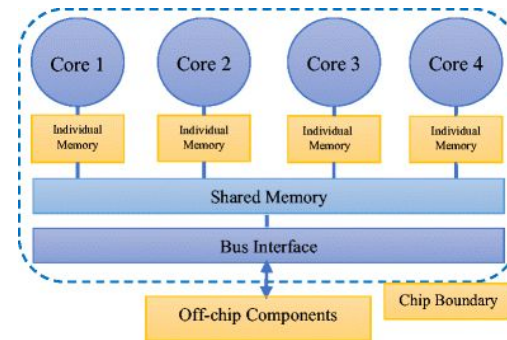
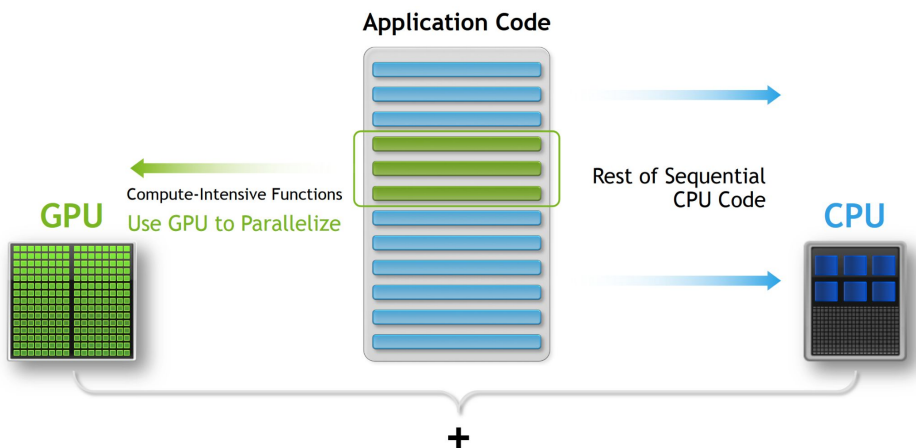


Image taken from [1]



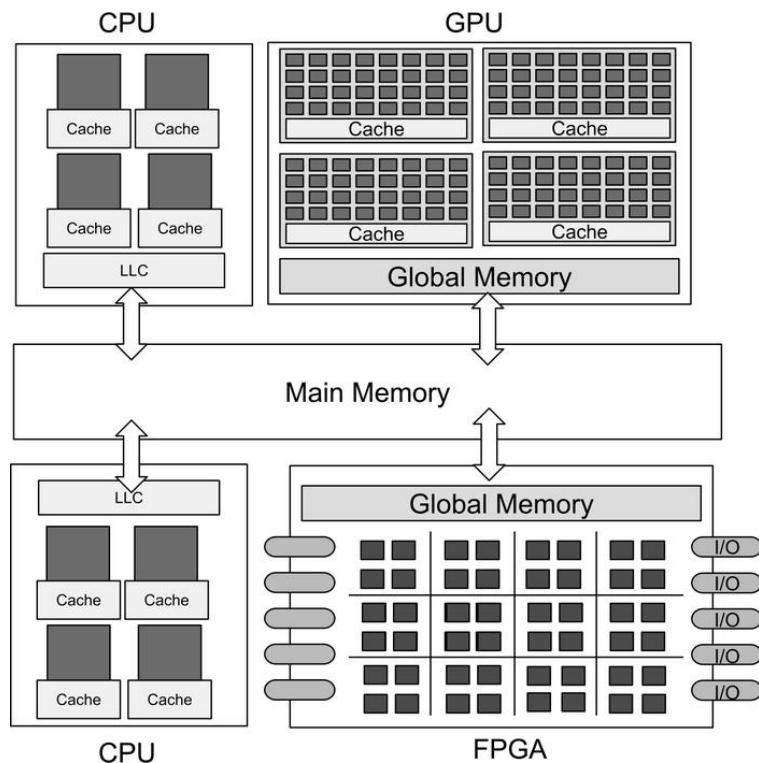
How are hardware accelerators used?



- 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

Heterogeneous computing

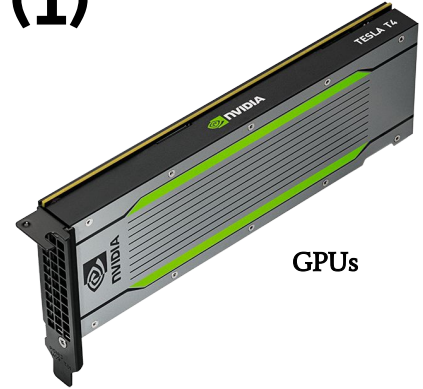
- Heterogeneous computing involves using multiple different types of processors to accomplish a task
- Code can run on more than one platform concurrently
- A heterogeneous system can consist of :
 - Different types of CPUs (i.e. combine compute powerful with less compute powerful but more power efficient CPU cores)
 - Hardware accelerators



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

- **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



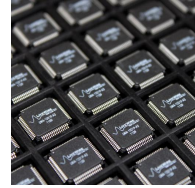
GPUs



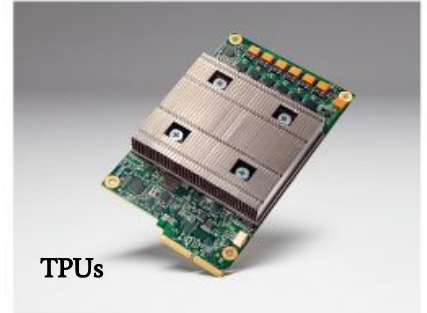
FPGAs

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
- **TPU** (Tensor Processing Unit)
 - Optimised to perform matrix-multiplication operations / used in e.g. NN and RF training
- **VPU** (Vision Processing Unit)
 - Used to accelerate machine vision algorithms, i.e. CNNs , AI etc.



ASIC



TPUs



VPUs

The GPU

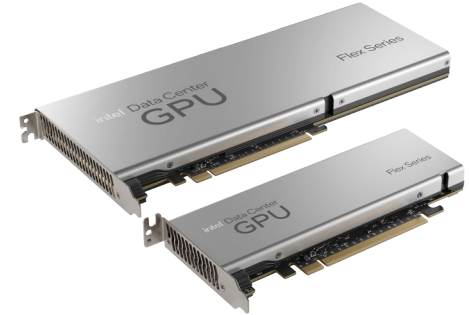
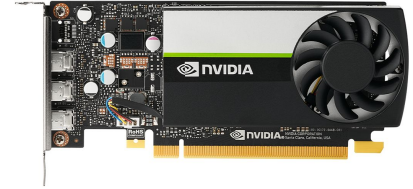
The Graphic Processing Unit (GPU)

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.



Multi-core vs many-core architectures

Multi-core processors

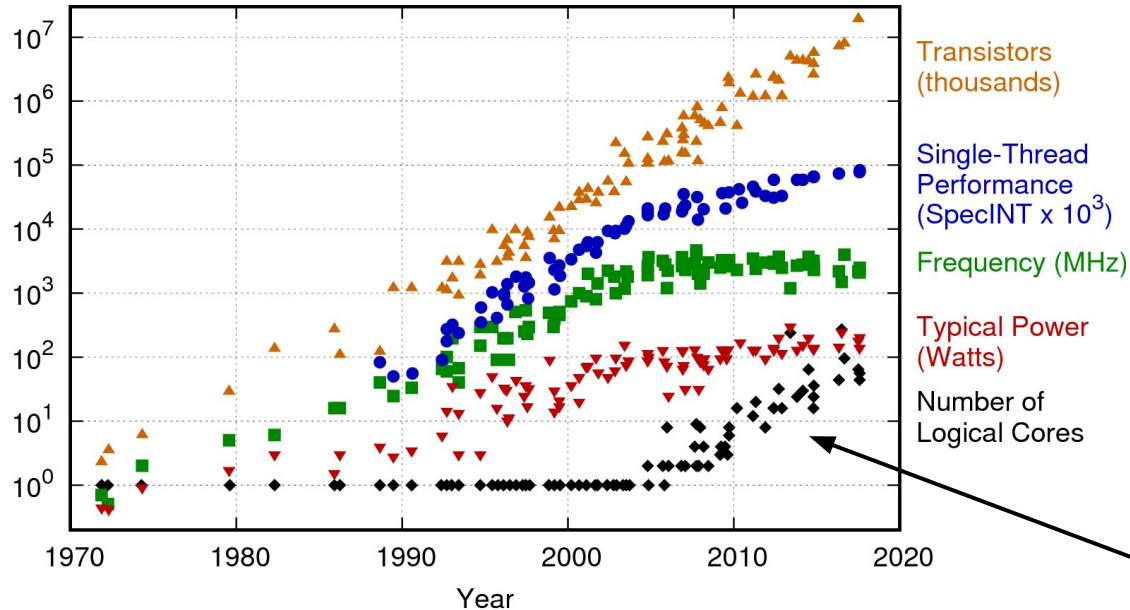
- Built on a single IC with two or more processing units (**cores**)
- Emphasis on high single-thread performance
- Better latency
- Can be complemented by a many-core system

Many-core processors

- Much higher degree of parallelism compared to a multi-core processors
- Emphasis on maximizing throughput
- Lower single-threaded performance and worse latency compared to multi-core processors

Why GPUs?

42 Years of Microprocessor Trend Data



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

Image source [\[1\]](#)

- Moore's law states that the number of transistors in a dense IC doubles every ~2 years.
- Since ~2010 there seems to be a plateauing in single-thread performance
- Gains expected through exploiting parallelization

Number of cores has started to increase

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

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

SM

The NVidia GPU architecture



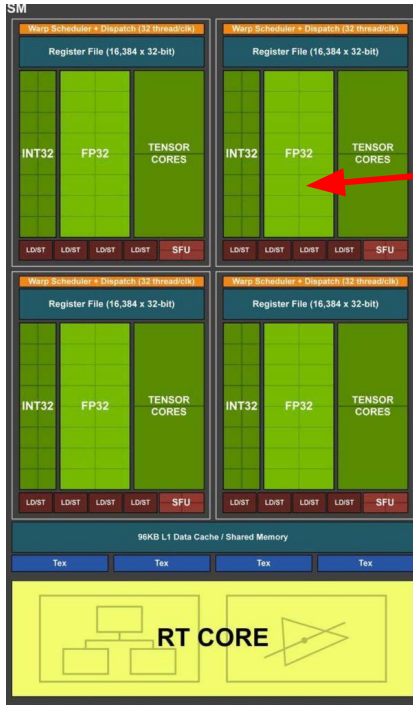
- The GPU architecture is built around a scalable array of **Streaming Multiprocessors (SM)**.
- Each SM in a GPU is designed to support concurrent execution of hundreds of threads

PCIe (*Peripheral Component Interconnect Express*): Can be used for connecting GPU to host CPU

SM

NVlink : Can be used to connect to additional GPUs

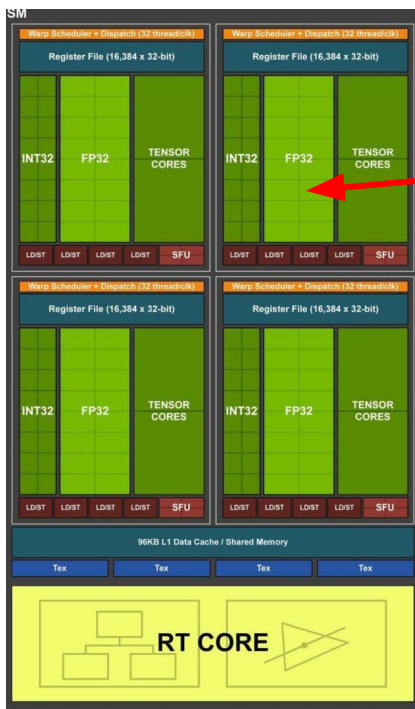
The Streaming Multiprocessor



The SM consists of :

- **Execution cores**
 - e.g. single precision floating-point, special function units etc.

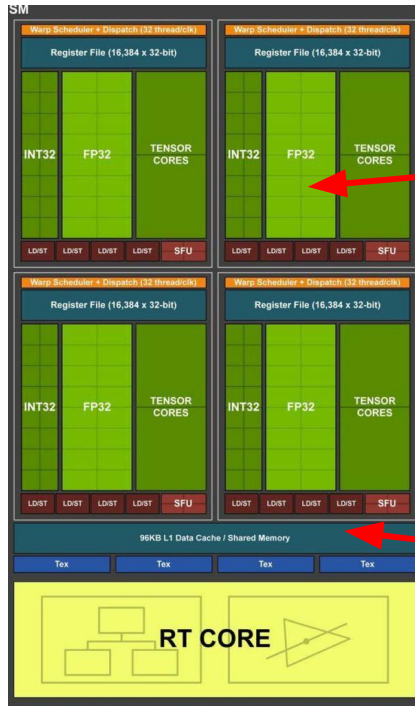
The Streaming Multiprocessor



The SM consists of :

- **Execution cores**
 - e.g. single precision floating-point, special function units etc.
- **Schedulers for warps**
 - These are used for issuing instructions to warps based on a particular scheduling policies.

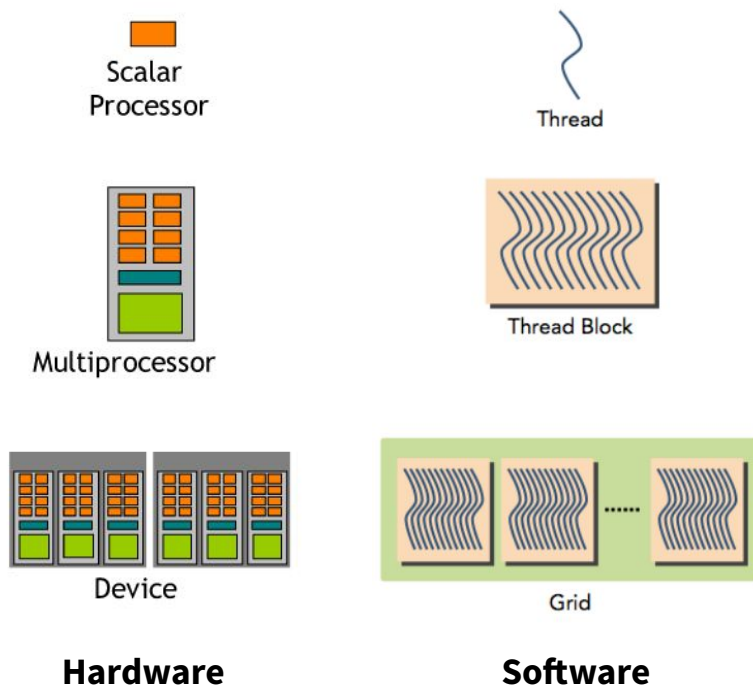
The Streaming Multiprocessor



The SM consists of :

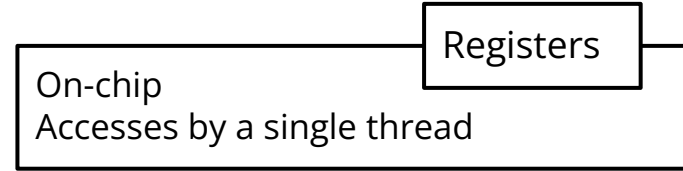
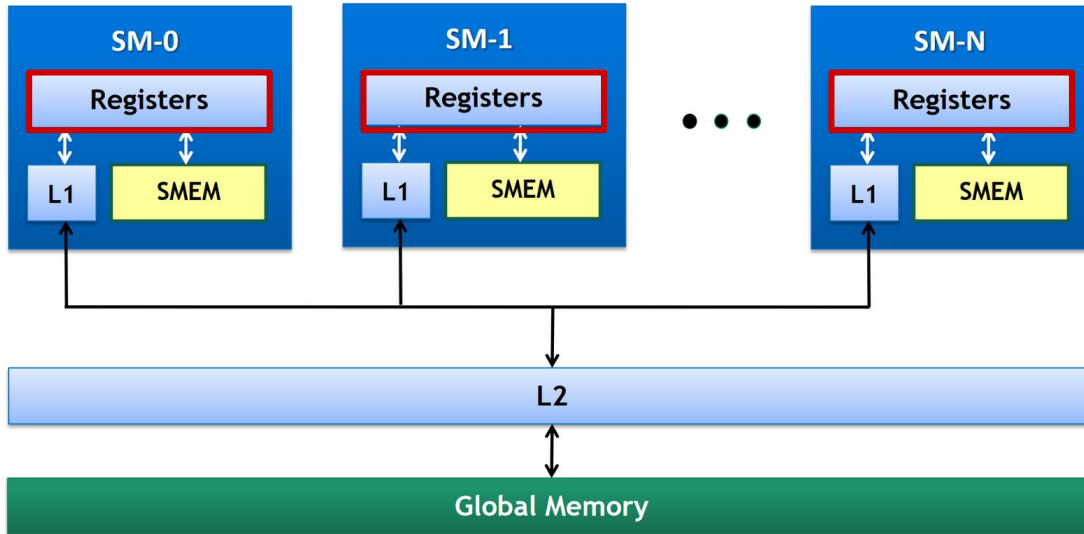
- **Execution cores**
 - e.g. single precision floating-point, special function units etc.
- **Schedulers for warps**
 - These are used for issuing instructions to warps based on a particular scheduling policies.
- **Registers**
 - fast on-chip memory used to store operands for the operations executed by the GPU cores
- **Caches**
 - Intermediate high-speed storage resources between the processor and memory
 - L1/constant/texture cache, Shared memory

Hardware to software mapping

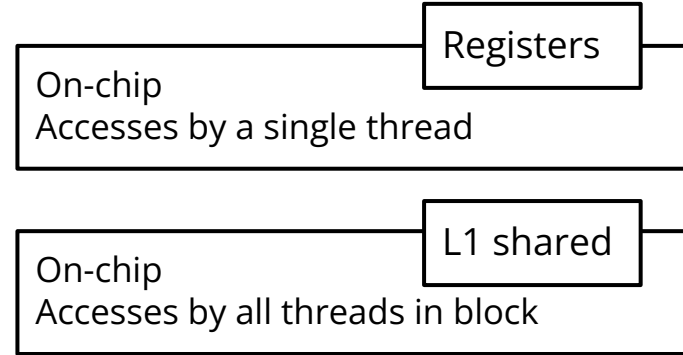
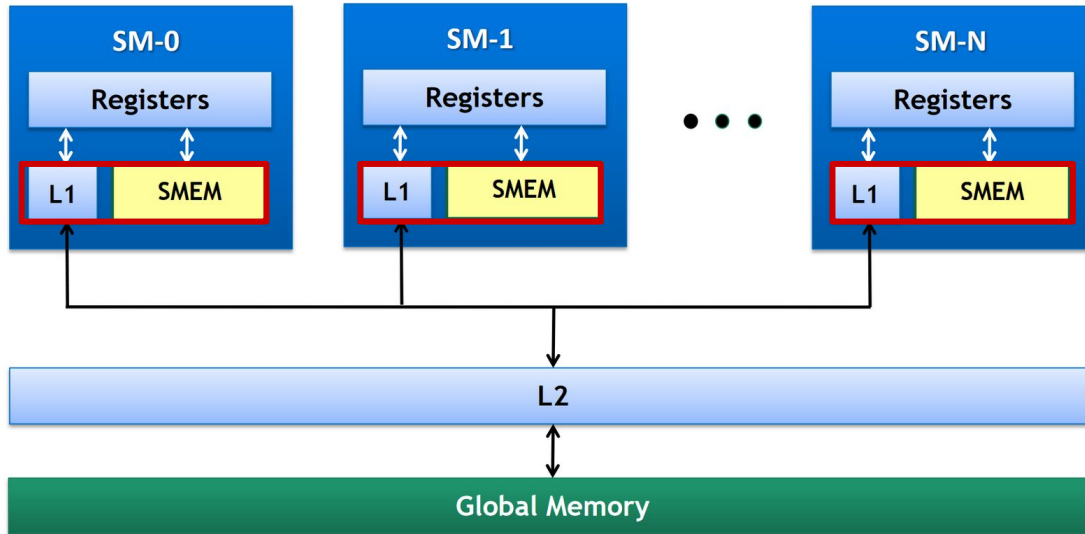


- A **scalar processor** or **CUDA core** is **equivalent to a software thread**
- Scalar processors are grouped into a SM
- Each execution of a GPU function is done concurrently on a number of threads referred to as a **thread block**
- Each thread block is executed by one SM and cannot be migrated 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**

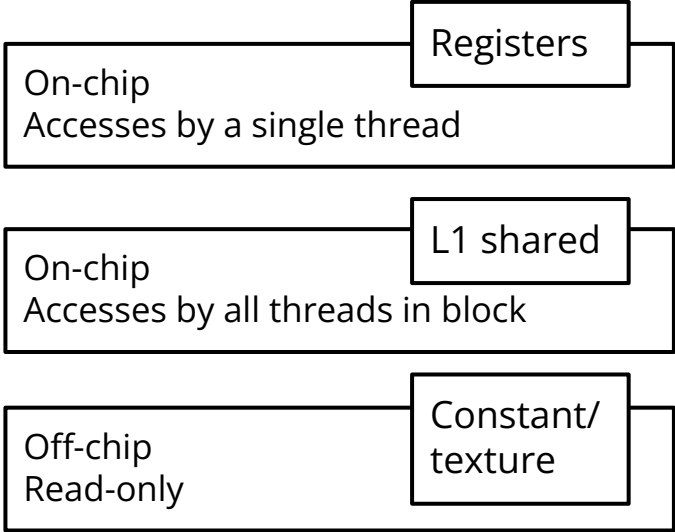
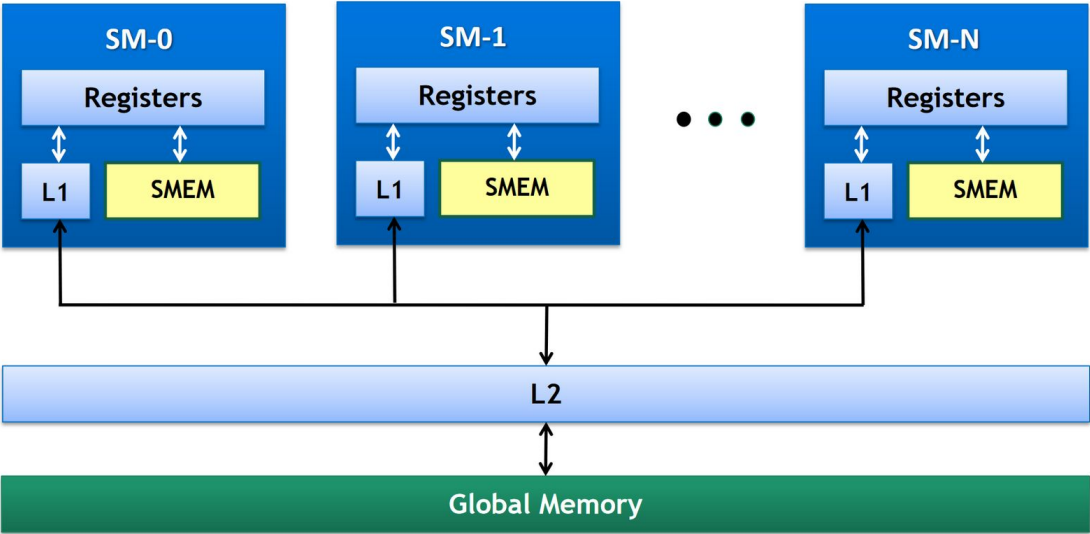
Memory hierarchy of the GPU



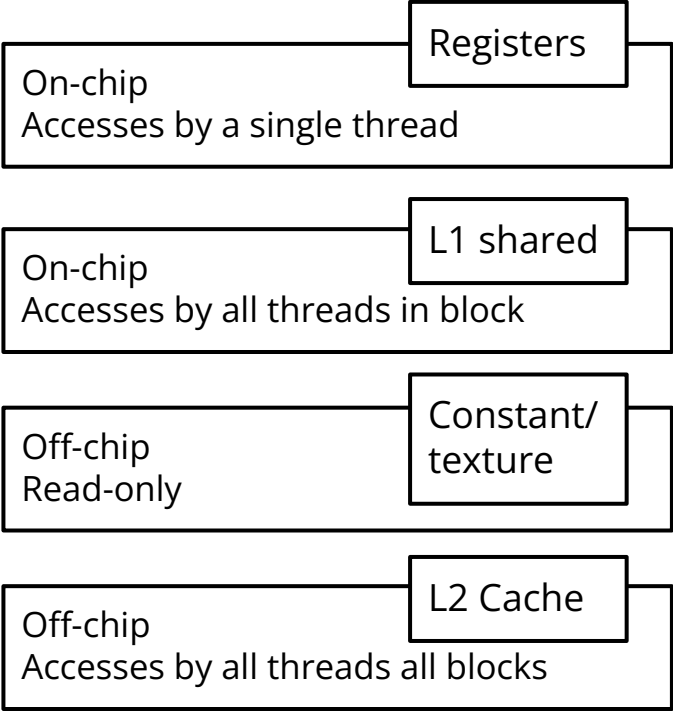
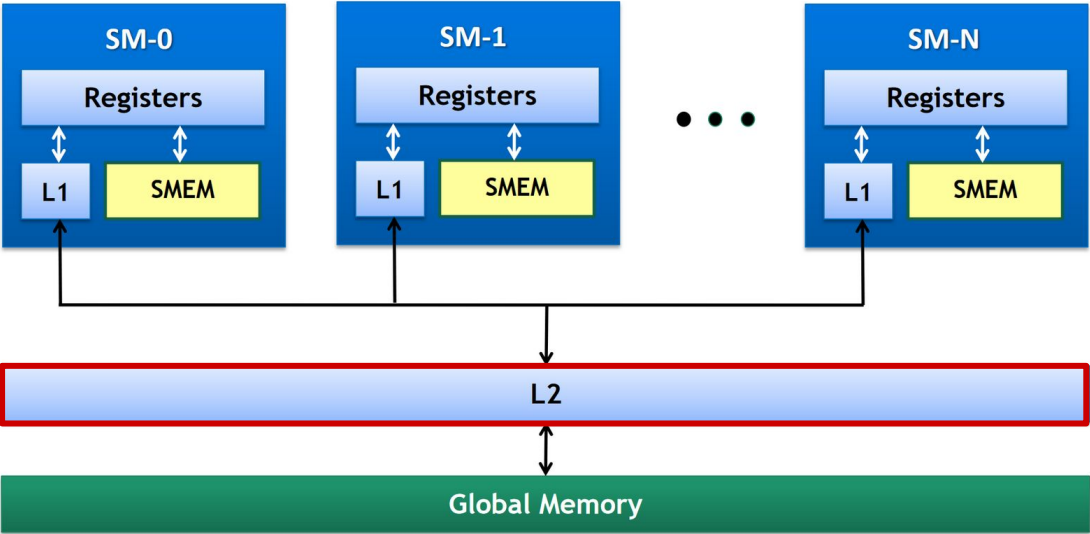
Memory hierarchy of the GPU



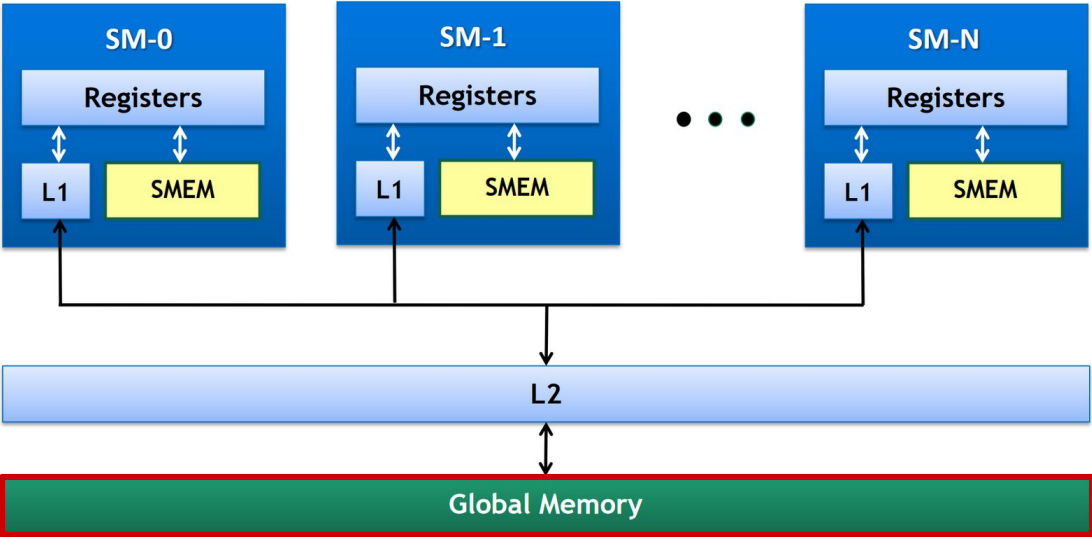
Memory hierarchy of the GPU



Memory hierarchy of the GPU



Memory hierarchy of the GPU



- Registers

On-chip
 Accesses by a single thread
- L1 shared

On-chip
 Accesses by all threads in block
- Constant/
texture

Off-chip
 Read-only
- L2 Cache

Off-chip
 Accesses by all threads all blocks
- Global

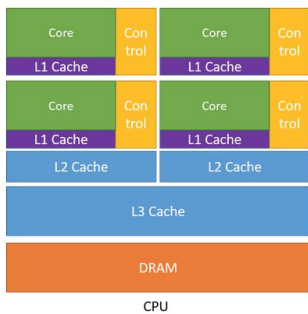
DRAM / large
 Accessed by device & host

Image source [1]

CPU vs GPU - overview of main differences

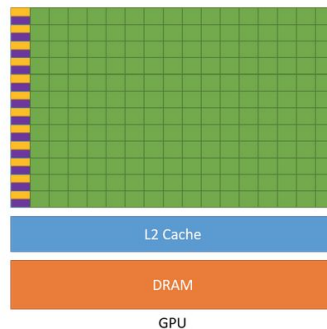
CPU

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



GPUs

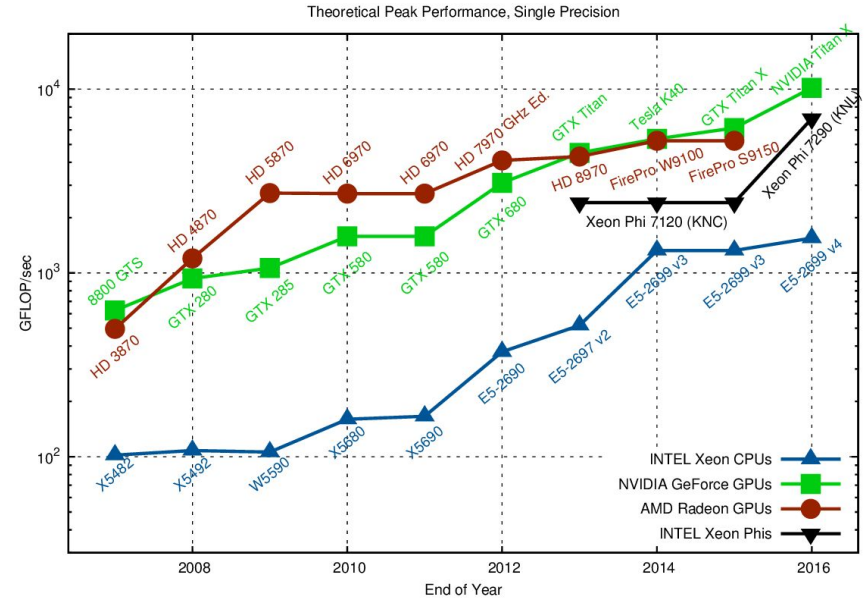
- ~O(1000) of less powerful cores
 - Smaller instruction set
- High throughput
- Parallel processing
- Simple operations
- Better per-watt performance



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



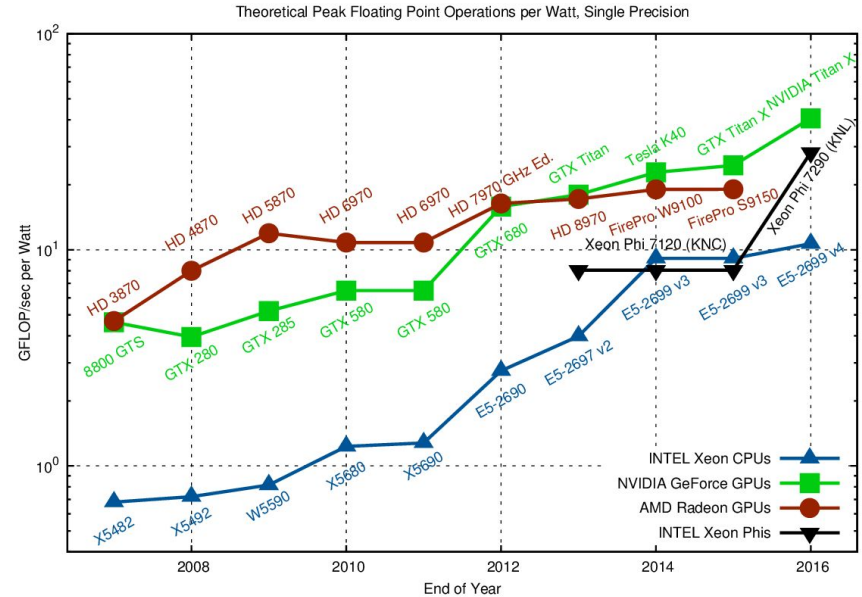
Performance comparison of CPUs and GPUs (2)

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





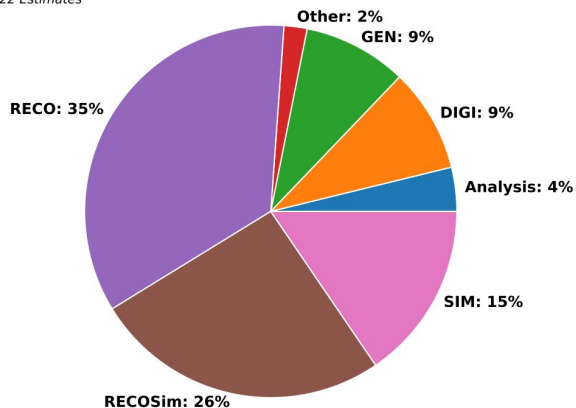
GPUs in High Energy Physics



Computing needs in HEP

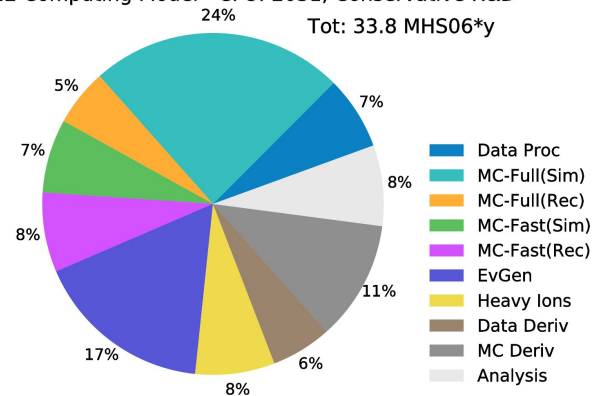
- Event generation
- Simulation
- Event reconstruction
- Event post-processing
- Data analysis

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



[\[Link\]](#)

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



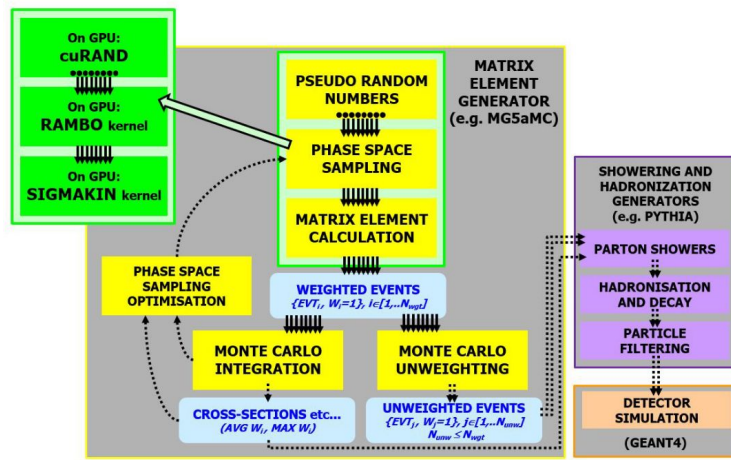
[\[Link\]](#)

How can GPUs help?

- Event generation
- Simulation
- Event reconstruction
- Event post-processing
- Data analysis

- GPU enabled event generator i.e. Madgraph

[ii](#)

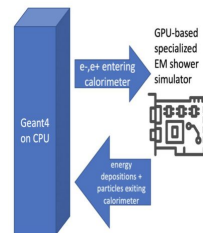


How can GPUs help?

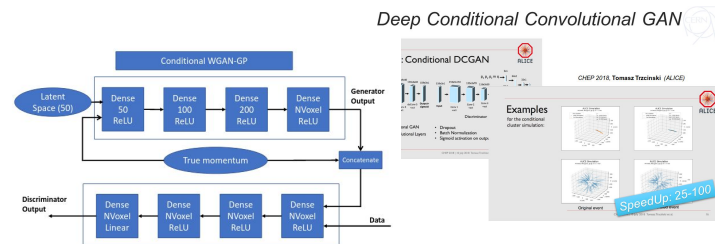
- Event generation
- Simulation
- Event reconstruction
- Event post-processing
- Data analysis

- GPU based Geant4 application (i.e. AdePT)

[ii](#)



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



How can GPUs help?

- Event generation
- Simulation
- Event reconstruction
- Event post-processing
- Data analysis

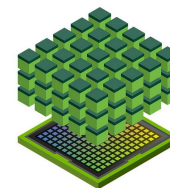
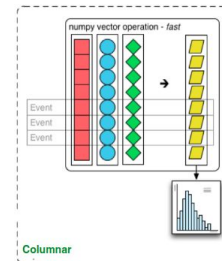
- Track reconstruction, primary vertex reconstruction, raw data unpacking, clustering etc.
- Various efforts in different experiments (Patatrack track reconstruction [\[i\]](#), Allen project [\[ii\]](#), ALICE TPC track reconstruction [\[iii\]](#) etc.)



How can GPUs help?

- Event generation
- Simulation
- Event reconstruction
- Event post-processing
- Data analysis

- Training and inference of ML models
- Perform HEP analysis using columnar analysis paradigm tools (i.e. coffea [\[1\]](#))





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



nvidia-smi

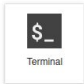
nvidia-smi: NVIDIA System Management Interface program

- Command line utility
- Aids in the management and monitoring of NVIDIA GPU devices

Lets try this out!

Click on the following link to access a GPU

https://binderhub.ssl-hep.org/v2/gh/research-software-collaborations/courses-hsf-india-december2023/gpu_onlycudagpu_true

Click on the Terminal icon (will look like this → )

Then type the following in the terminal:

```
nvidia-smi
```

What do you see ? Let's now try running a small utility script :

```
cd hsf-india-gpus
```

```
nvcc deviceInfo.cu -o deviceInfo
```

```
./deviceInfo
```

What do you see now?

nvidia-smi

nvidia-smi: NVIDIA System Management Interface program

- Command line utility
- Aids in the management and monitoring of NVIDIA GPU devices

Lets try this out!

Click on the following link to access a GPU

https://binderhub.ssl-hep.org/v2/gh/research-software-collaborations/courses-hsf-india-december2023/gpu_onlycudagpu_true

Click on the T

Then type the

nvidia-smi

What do you see

cd hsf-india-g

nvcc deviceInfo

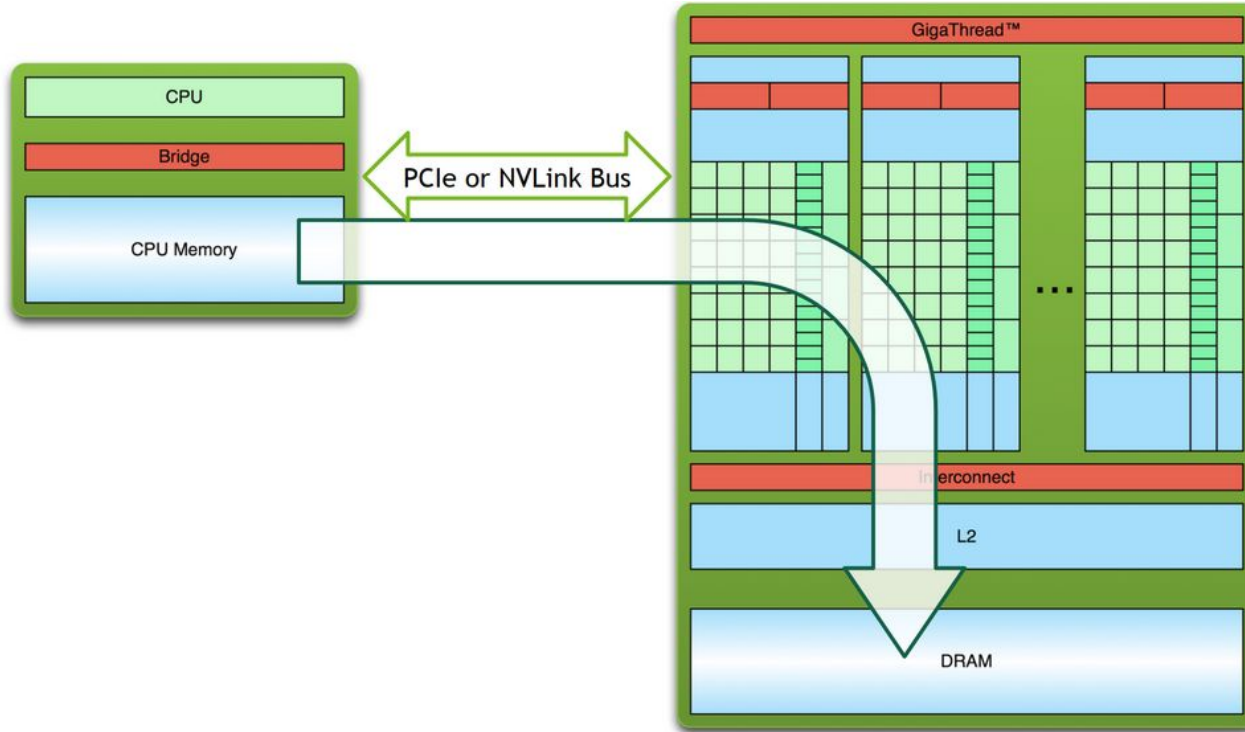
./deviceInfo

What do you see now?

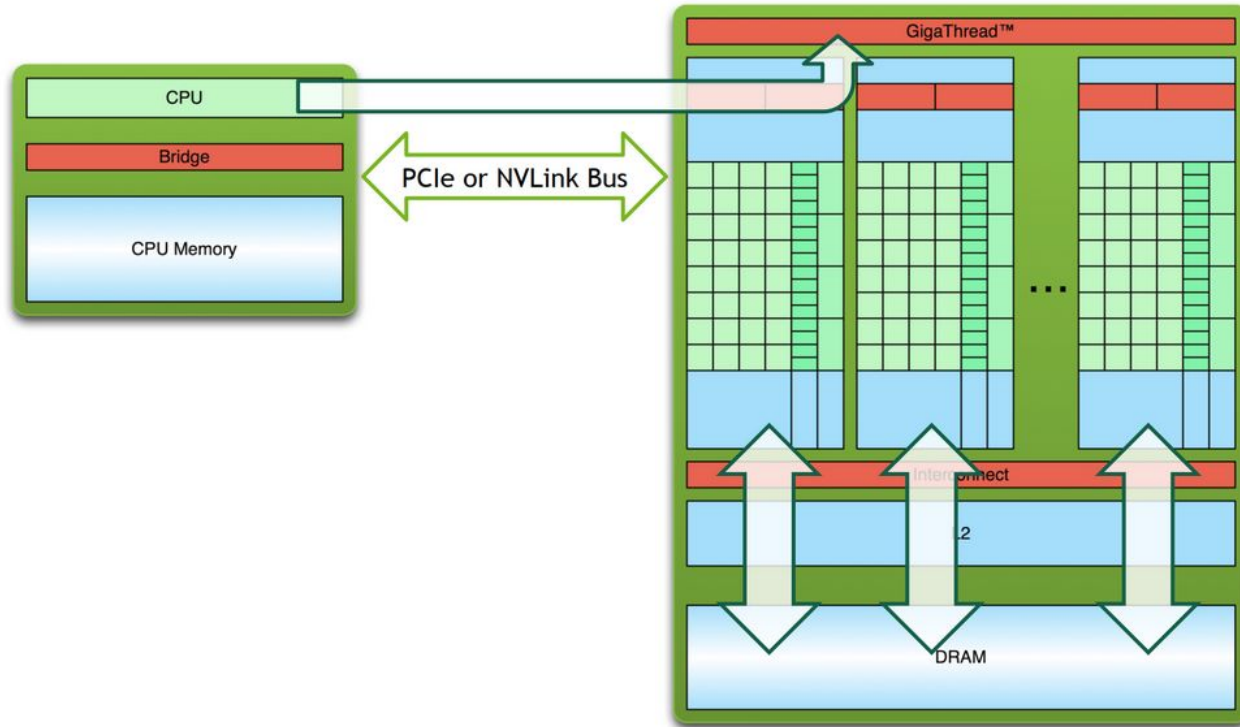
Can you answer some questions?

- How many devices are found?
- What type of GPUs are they?
- How many SMs per device?
- What is the warp size?
- How many threads are allowed per block?

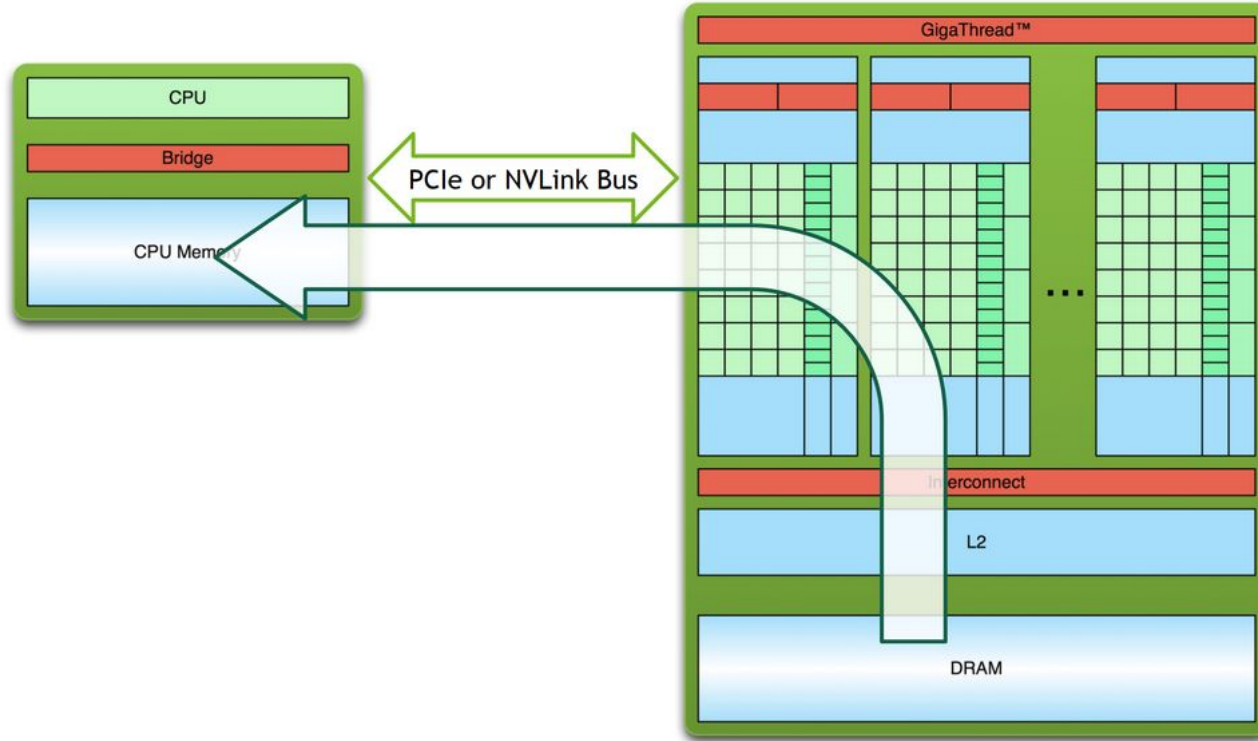
1. Copy data for host to device



2. Execute the CUDA program

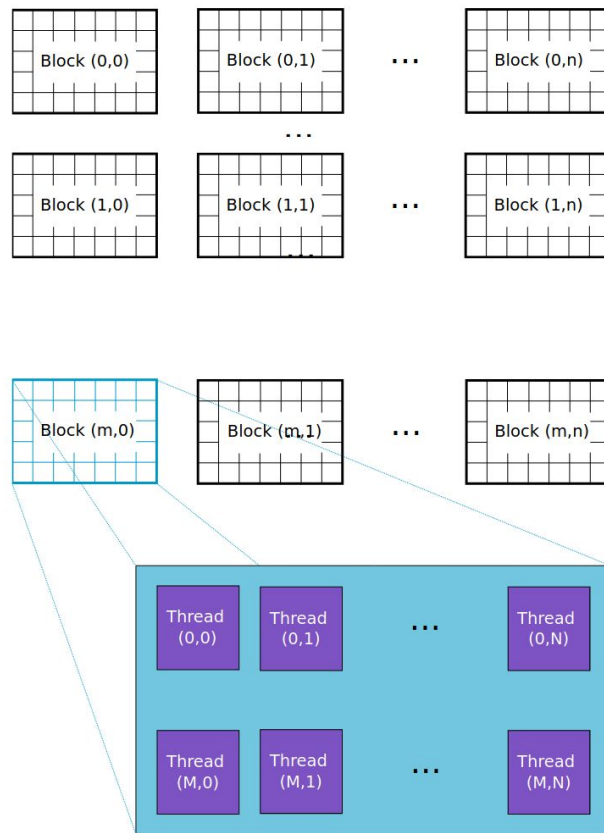
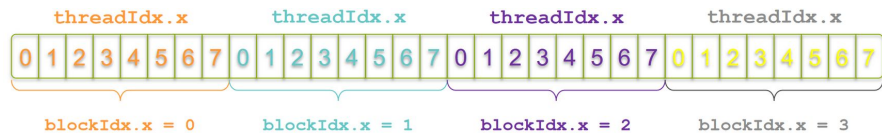


3. Copy data from device back to host



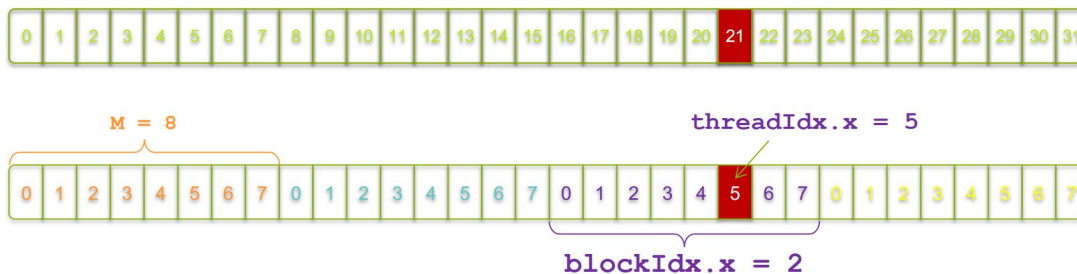
Threads & blocks

- 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 :
 - $\text{index} = \text{threadIdx.x} + \text{blockIdx.x} * M;$

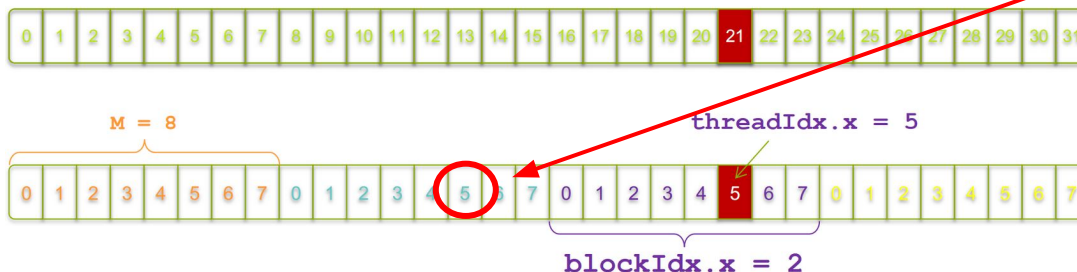


```
int index = threadIdx.x + blockIdx.x * M;  
          =      5      +      2      * 8;  
          = 21;
```

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 :
 - $\text{index} = \text{threadIdx.x} + \text{blockIdx.x} * M;$

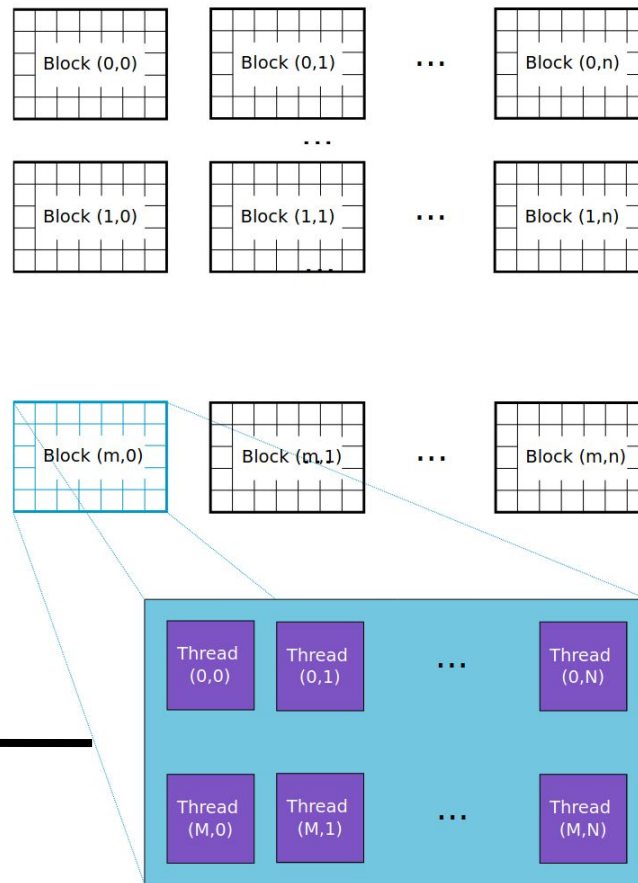
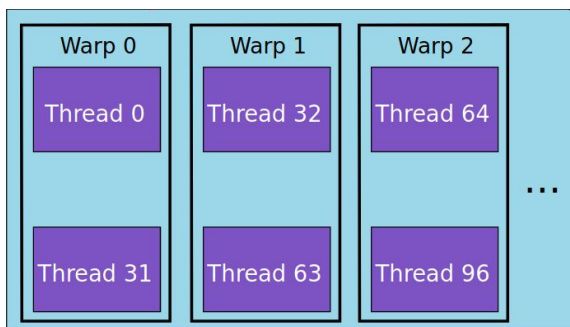
Q: What is the index of this element??



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

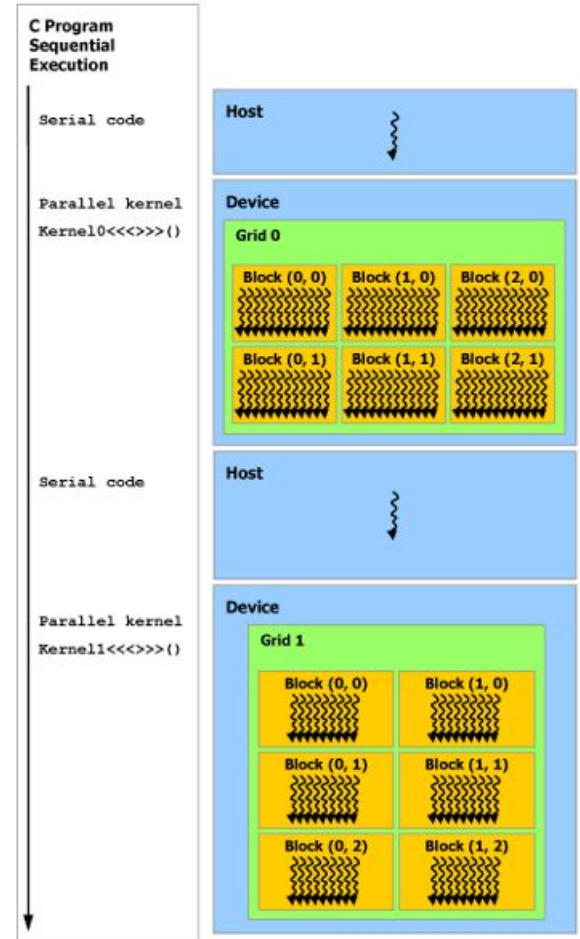
Warps

- Within a thread block, threads are executed in groups \rightarrow **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)
```

Memory management

- **Host pointers :**
 - Typically not passed to device code
 - Typically not dereferenced in device code
- **Device pointers :**
 - Typically passed to device code
 - Typically not dereferenced in host code

```
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);
```

Let's take a look at the syntax of `cudaMalloc`

For transfers between host and device memory the direction can be :

- Copying data from CPU to GPU
- Copying data from GPU to CPU

Memory management

- **Host pointers :**
 - Typically not passed to device code
 - Typically not dereferenced in device code
- **Device pointers :**
 - Typically passed to device code
 - Typically not dereferenced in host code

```
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)); ← Let's take a look at the syntax of cudaMalloc
// 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);
```

For transfers between host and device memory the direction can be :

- Copying data from CPU to GPU
- Copying data from GPU to CPU

Remember the order for copying variables from host \longleftrightarrow device!

Synchronization

- **CUDA kernel calls are asynchronous :**
 - Once the kernel is launched the main program that is executed on the CPU continues normally !
- Additionally, execution order of blocks on a SMs is arbitrary
 - **We need a way to synchronise!**
- We can call **CudaDeviceSynchronize()** from host
 - blocks the CPU execution until all work launched on the device has finished.
- Includes both:
 - kernel launches
 - memory copies



Grid level synchronization

Synchronization

For each kernel launch with N threads/block & M blocks :

- Execution order of threads within one block is arbitrary :
 - Only exception are threads in the same warp which are processed simultaneously
- We might have a problem, where we require all threads in a specific block to have completed execution of a specific task before continuing the next task
- To synchronize threads within one block one can call **__syncthreads()** within the kernel

```
__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]  
    }  
}
```



Block level synchronization

Atomic operations

- Useful when **modifying the same value in memory from different threads** :
 - Are used to prevent race conditions in multithreaded applications
 - Read-modify-write cannot be interrupted
 - Appear to be one operation
- Atomics are special hardware instruction on NVIDIA GPUs e.g.:
 - atomicAdd/Sub (Add or subtract)
 - e.g. syntax : `atomicAdd(int* address, int val);`
 - atomicMax/Min (Find max or min)
 - atomicExch/CAS (Swap or conditionally swap variables)
 - e.g. syntax : `atomicCAS (&addr, compare, value)`
 - atomicAnd/Or/Xor (bitwise operations)
 - ...

A[i]



SumA[i]



Adding elements in a vector

Let's start by writing a CUDA kernel that calculated the sum of the elements of a vector :

```
__global__ void add_array(float* A, float* sum) {  
    int idx = threadIdx.x + blockIdx.x * blockDim.x;  
    if (idx < N) {  
        *sum +=A[idx];  
    }  
}
```

- There are 3 instructions that will be executed :
 - **Load** the value of A for each thread
 - **Read** the value of c
 - **Modify** the value of c

A[i]



SumA[i]

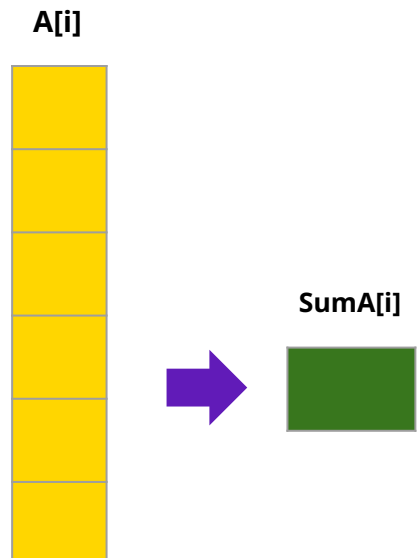


Adding elements in a vector

Let's start by writing a CUDA kernel that calculated the sum of the elements of a vector :

```
__global__ void add_array(float* A, float* sum) {  
    int idx = threadIdx.x + blockIdx.x * blockDim.x;  
    if (idx < N) {  
        *sum +=A[idx];  
    }  
}
```

- There are 3 instructions that will be executed :
 - Load the value of A for each thread
 - Read the value of c
 - Modify the value of c



The behaviour of this kernel will be unpredictable!
The read/writes can happen in random orders.
The sum might be incorrect!!!

Adding elements in a vector

Using **atomicAdd** to sum the vector elements :

```
__global__ void add_array(float* A, float* sum) {  
    int idx = threadIdx.x + blockIdx.x * blockDim.x;  
    if (idx < N) {  
        atomicAdd(sum, A[idx]);  
    }  
}
```

Each read-modify-write access cannot be interrupted

A[i]



SumA[i]



Now the sum will be correct!!

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()` / `CudaDeviceSynchronize()`

or

```
__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);
}
```

Good practices : 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.

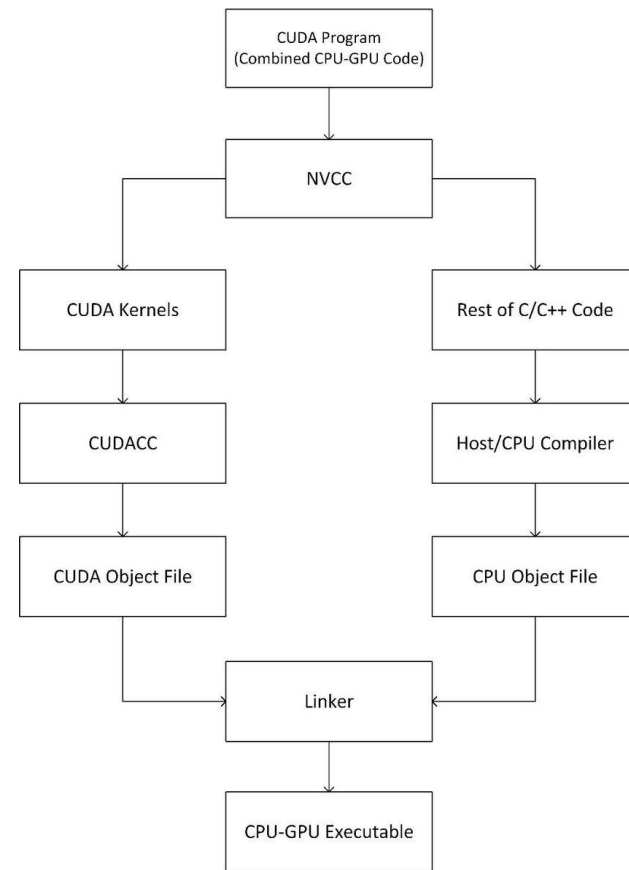
```
cudaError_t err = cudaGetLastError();
```

- 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
```



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>>>()`

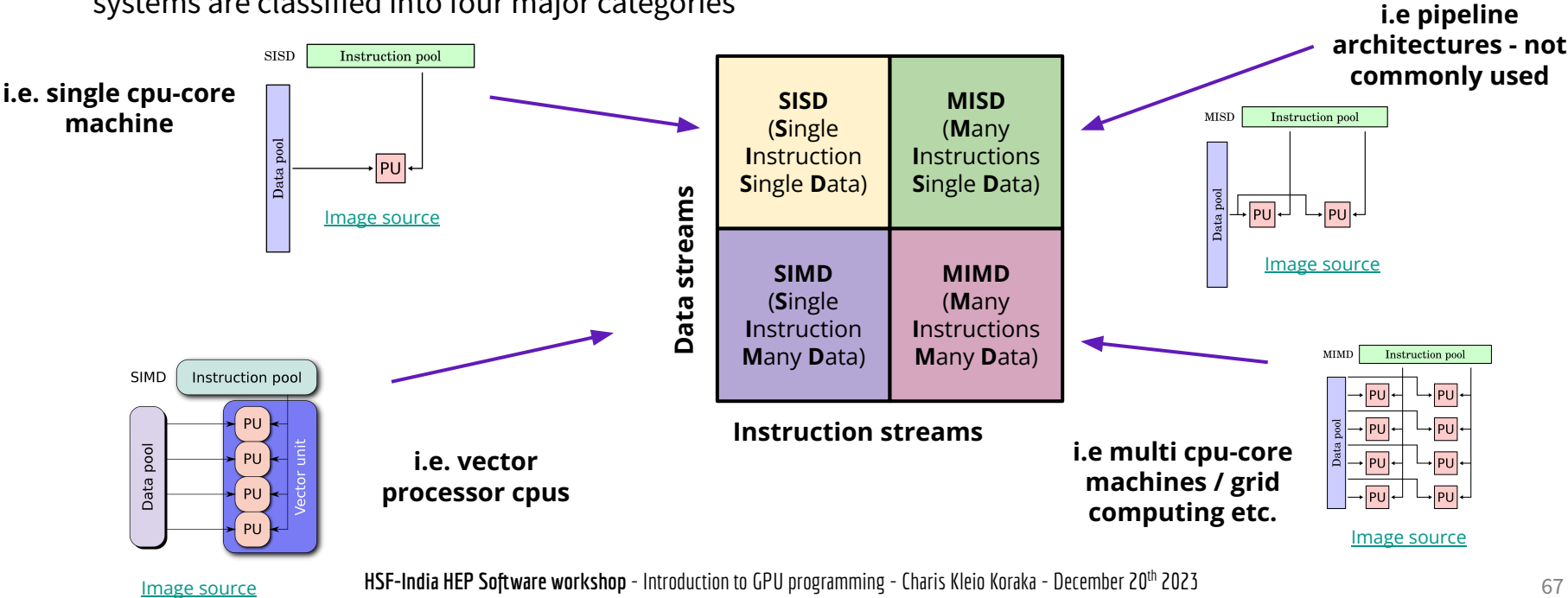
Back-up

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](#)

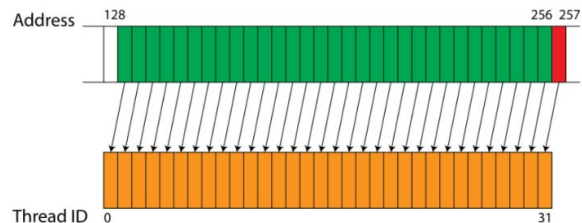
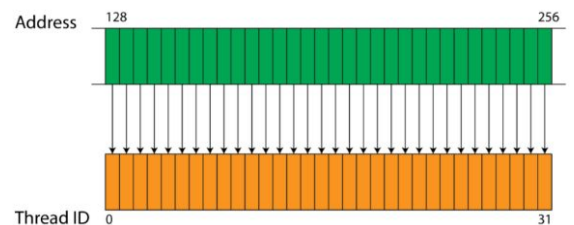
Flynns classification of computer architecture

- Based on the number of instruction and data streams that can be processed simultaneously, computing systems are classified into four major categories

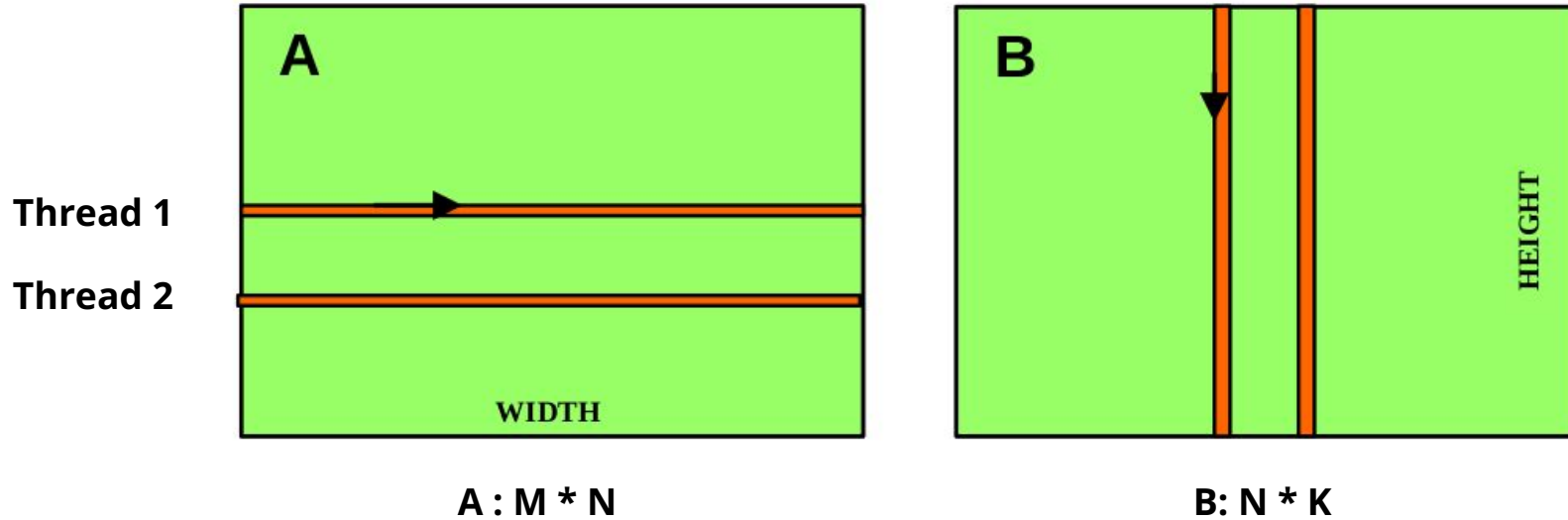


Coalesced global memory access

- 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



Matrix multiplication



Thread synchronization

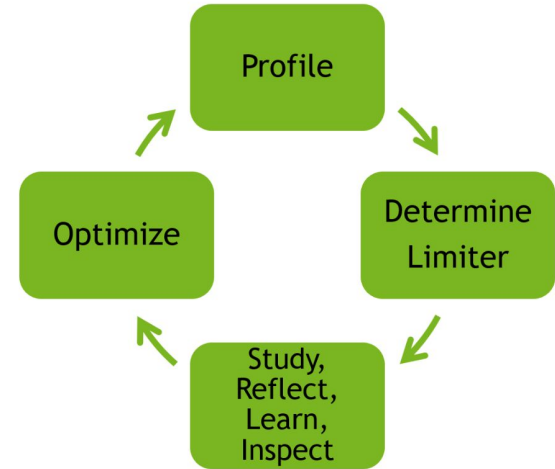
- A kernel call is asynchronous with respect to the host thread :
 - After a kernel is invoked, the program returns to the host side and continues execution.
- There are two levels of synchronization
 - **Block level**
 - **Grid level**
- To synchronize threads within one block :
 - Call `__syncthreads()` 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);
```

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.

```
$ nvidia-smi --query-gpu=memory.free --format=%f --no-header
[Matrix Multiply Using CUDA] - Starting...
==27694== Nsight Systems is profiling process 27694, command: matrixMul
GPU device 0: "GeForce GTX 640M LE" with compute capability 3.0

MatrixA(320,320), MatrixB(640,320)
Computing result using CUDA kernel...
done
Performance: 35.35 GFlop/s, Time: 3.708 msec, Size: 131072000 Ops, WorkgroupSize: 1024 threads/block
Checking computed result for correctness: OK

Note: For peak performance, please refer to the matrixMulCUBLAS example.
==27694== Profiling application: matrixMul
==27694== Profiling result:
Time(s)      Time      Calls      Avg      Min      Max      Name
99.94%  1.15254s      301  3.7959ms  3.6928ms  3.7174ms  void matrixMulCUDA<int>::<float*, float*, float*, int, int>
0.02%   248.29us      2  248.29us  196.33us  279.35us  [CUDA memory H2D]
0.02%   248.29us      1  248.29us  248.29us  248.29us  [CUDA memcpy DtoH]

==27694== API calls:
Time(s)      Time      Calls      Avg      Min      Max      Name
49.61%  205.17ms      3  68.39ms  68.39ms  68.39ms  cudaMalloc
25.90%  140.57ms      1  140.57ms  140.57ms  140.57ms  cudaEventSynchronize
22.23%  127.20ms      1  127.20ms  127.20ms  127.20ms  cudaDeviceReset
1.33%   7.6314ms      301  25.38us  23.05us  143.90us  cudaLaunch
0.20%   1.0342ms      3  478.00us  155.94us  804.30us  cudaMemory
0.11%   601.45us      1  601.45us  601.45us  601.45us  cudaDeviceSynchronize
0.10%   564.40us      1505  375ns    333ns    3.6790us  cudaSetupArgument
0.09%   490.44us      16  6.6500ms  397ns    221.95us  cudaDeviceGetAttribute
0.07%   406.61us      3  135.54us  115.07us  169.99us  cudaFree
0.02%   141.00us      301  470ns    431ns    2.4370us  cudaConfigureCall
0.01%   42.321us      1  42.321us  42.321us  42.321us  cudaEventFinalize
0.01%   31.655us      1  33.655us  33.655us  33.655us  cudaDeviceGetProperty
0.01%   31.900us      1  31.900us  31.900us  31.900us  cudaDeviceGetName
0.00%   21.974us      2  10.937us  8.5805us  13.200us  cudaEventRecord
0.00%   16.513us      2  8.2566us  2.6240us  13.880us  cudaEventCreate
0.00%   13.091us      1  13.091us  13.091us  13.091us  cudaEventElapsedTime
0.00%   9.1410us      1  9.1410us  9.1410us  9.1410us  cudaGetDevice
0.00%   2.6290us      2  1.3145us  998ns    1.1200us  cudaEventCount
0.00%   1.9970us      2  998ns    528ns    1.4770us  cudaDeviceGet
```

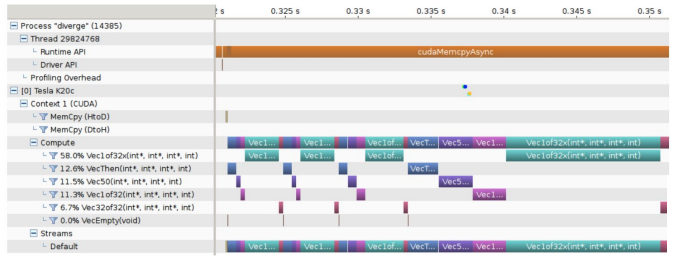
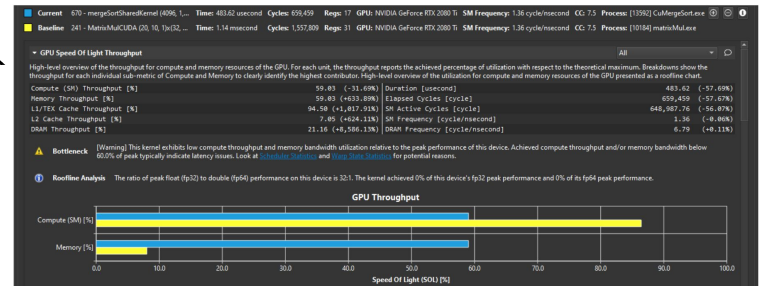


Image source [6]