

# Introduction to GPU programming (part 1)

By Kevin Fotso

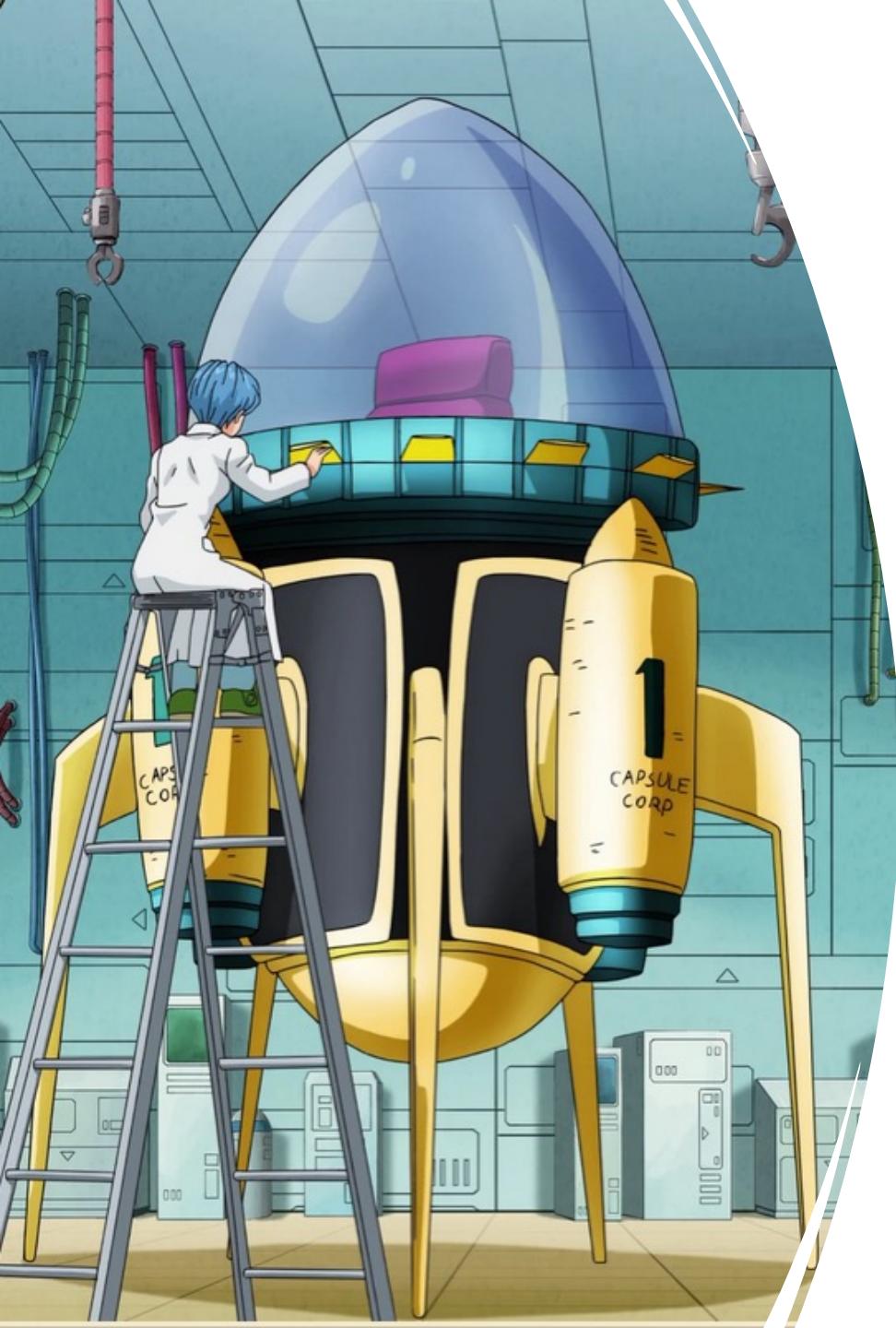
A photograph of a person's hand pointing their index finger towards a computer monitor. The monitor displays a dark-themed code editor with white and yellow text. The code is a Python script for Blender, specifically for creating a mirror modifier. It includes logic for different mirror operations (X, Y, Z) and handling operator classes. The background is a warm, indoor lighting.

```
mirror_mod = modifier_obj
# Set mirror object to mirror
mirror_mod.mirror_object = mirror_obj
operation = "MIRROR_X"
mirror_mod.use_x = True
mirror_mod.use_y = False
mirror_mod.use_z = False
operation = "MIRROR_Y"
mirror_mod.use_x = False
mirror_mod.use_y = True
mirror_mod.use_z = False
operation = "MIRROR_Z"
mirror_mod.use_x = False
mirror_mod.use_y = False
mirror_mod.use_z = True

# Selection at the end - add
modifier_obj.select= 1
modifier_obj.select=1
context.scene.objects.active = modifier_obj
("Selected" + str(modifier_obj))
modifier_obj.select = 0
bpy.context.selected_objects = []
data.objects[one.name].select = 1
print("please select exactly one object")

- OPERATOR CLASSES ---

types.Operator):
    X mirror to the selected object.mirror_mirror_x" or X"
    context):
        context.active_object is not None
        context.active_object.select = 1
        context.active_object.select = 0
        print("please select exactly one object")
```



# I- Part 1) Elements of computing

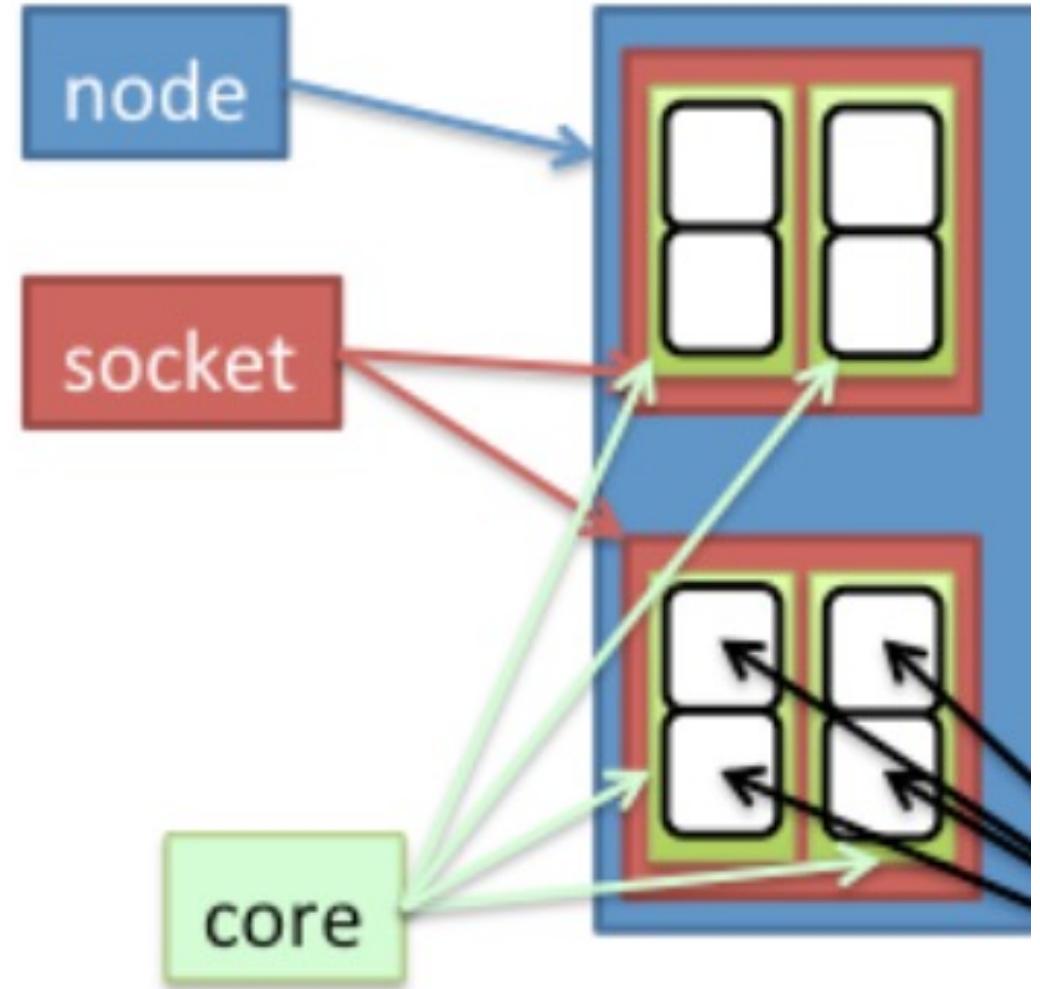
# Parts of the processing chip (1)



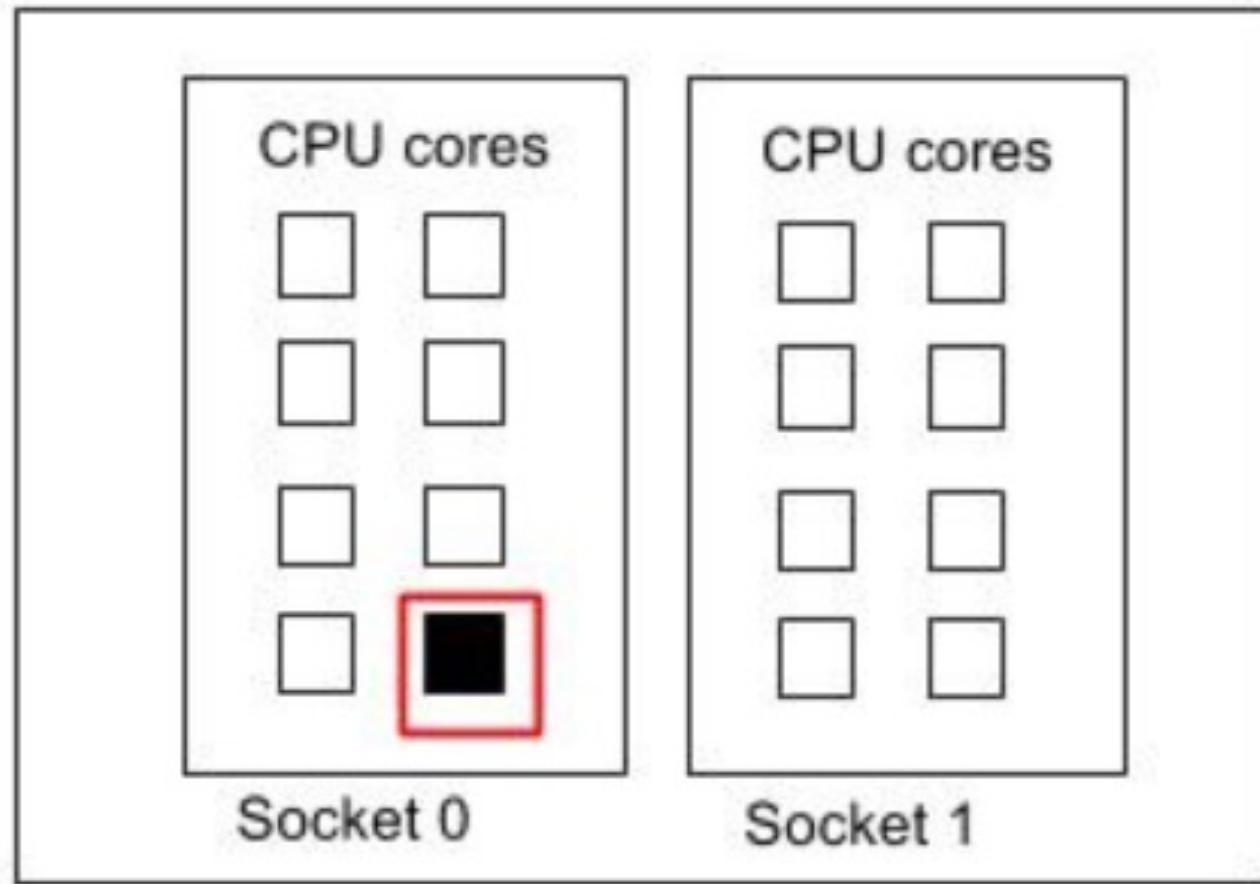
- A socket can be considered an entire computing chip in modern computers.
- Contains more than 1 core.
- On supercomputers a node usually contains more than 1 socket.

# Parts of the processing chip (2)

A core contains an Arithmetic and Logic Unit (**ALU**) + Registers

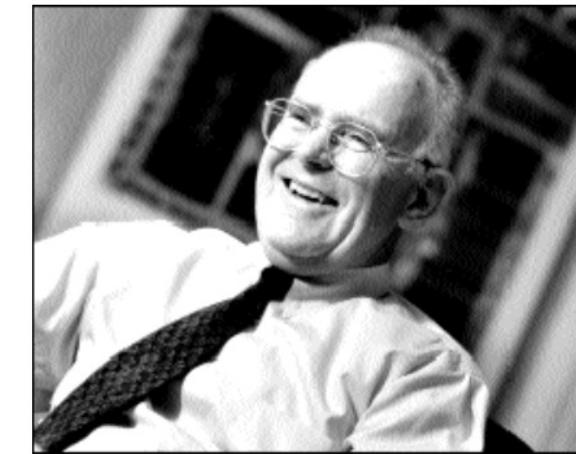
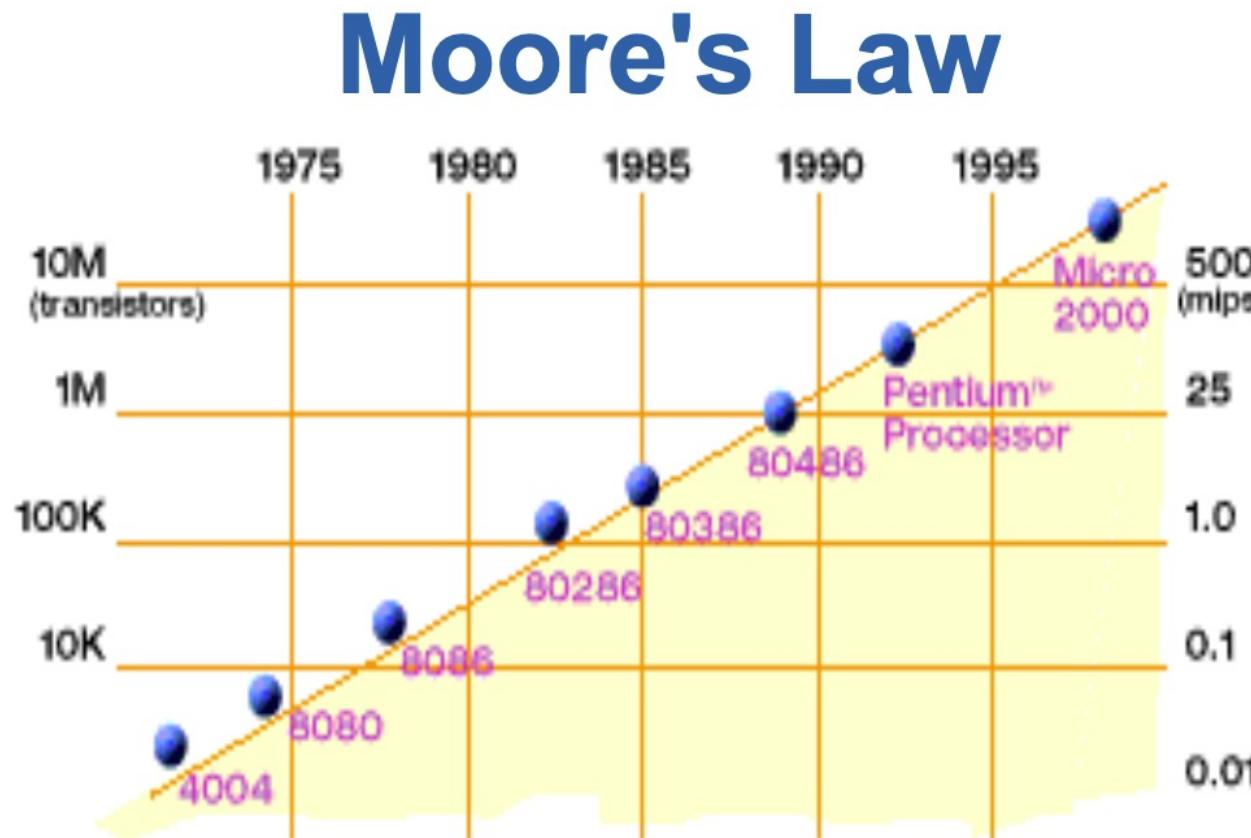


## Parts of the processing chip (2)



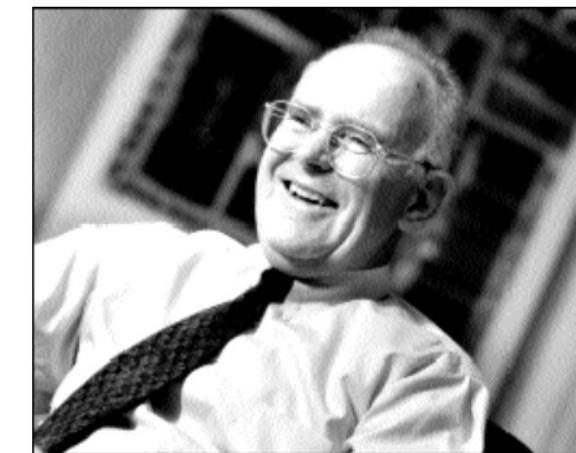
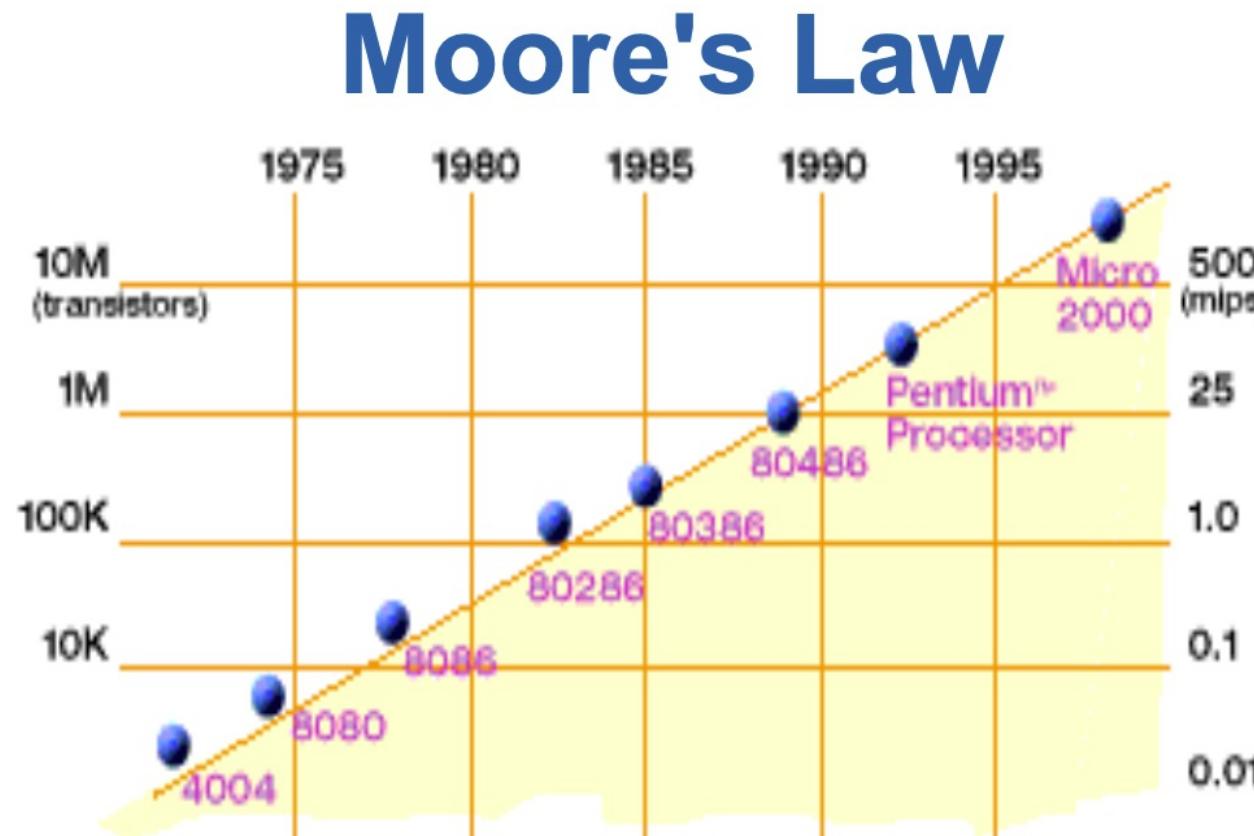
# Moore's law (1)

- Transistor count doubles in a chip every 18 months.

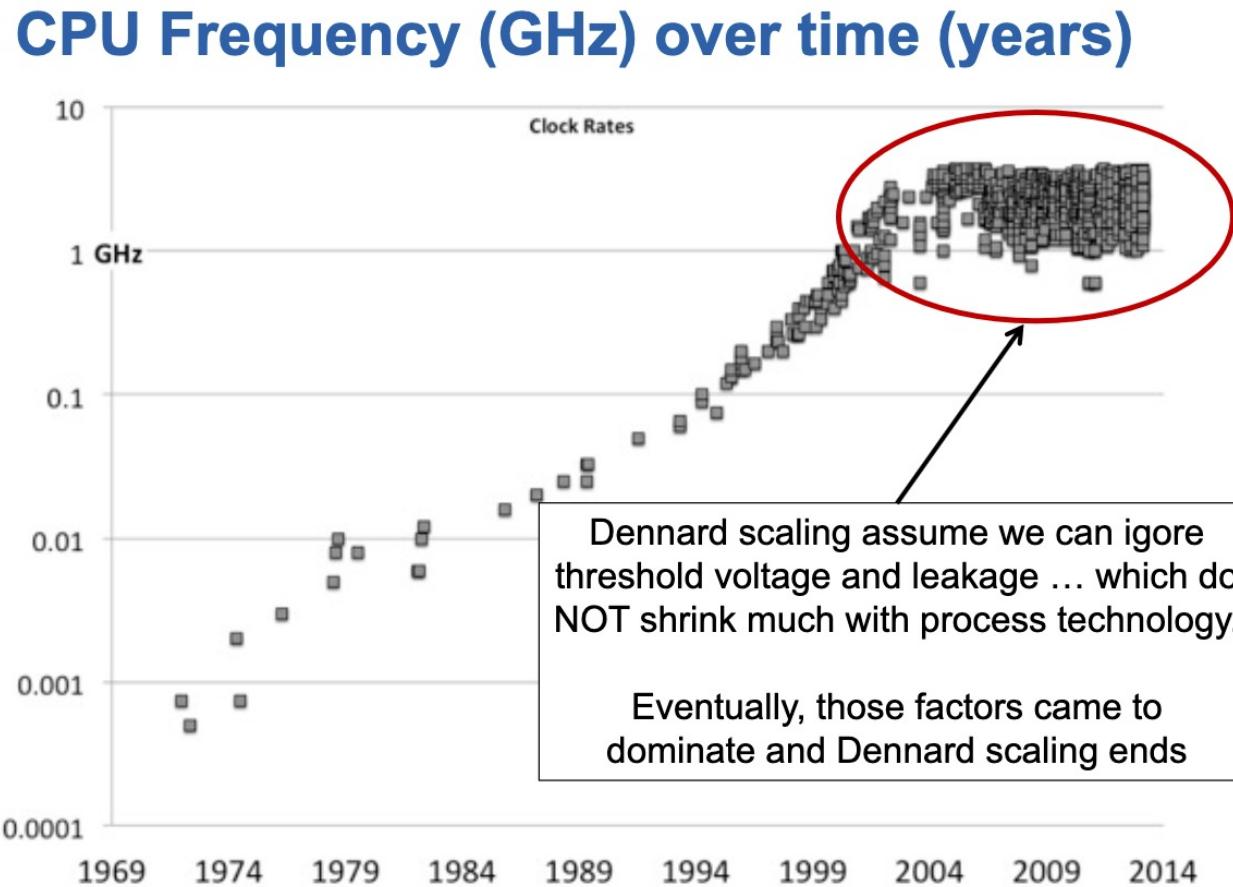


# Moore's law (2)

- Worked until around 2010 where it started to plateau



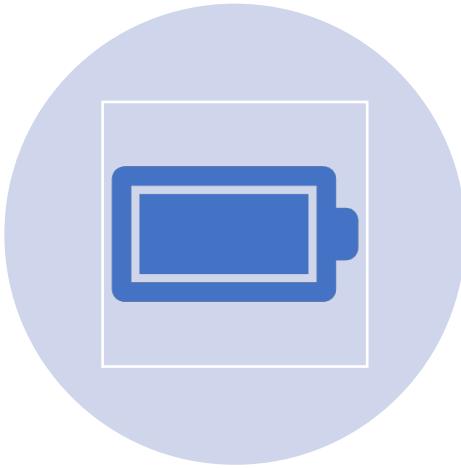
# Moore's law plateau



- Number of cycles per seconds plateaued in the 2010s
- Why?



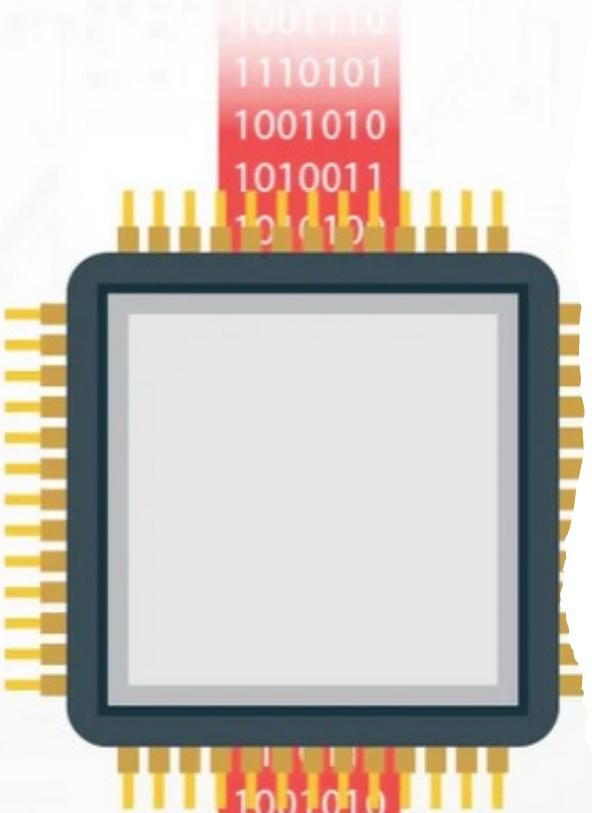
# Why?



TRADITIONALLY, POWER DENSITY STAYS CONSTANT AS CHIPS FEATURE GOT SMALLER.



HOWEVER AROUND 2010, CLOCK FREQUENCY COULD NOT BE RAISED FURTHER BECAUSE THE HEAT DISSIPATION OR POWER WOULD HAVE GONE TOO FAR.

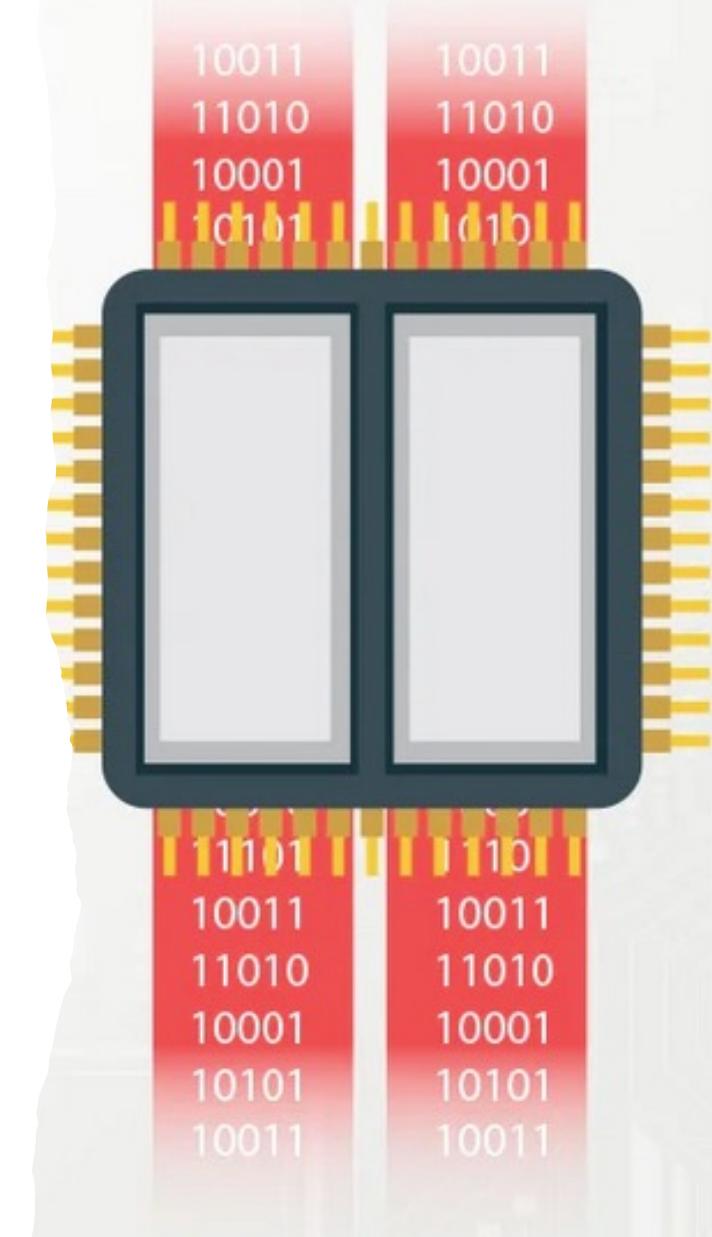


**Single-core**

# Increase need of parallelism

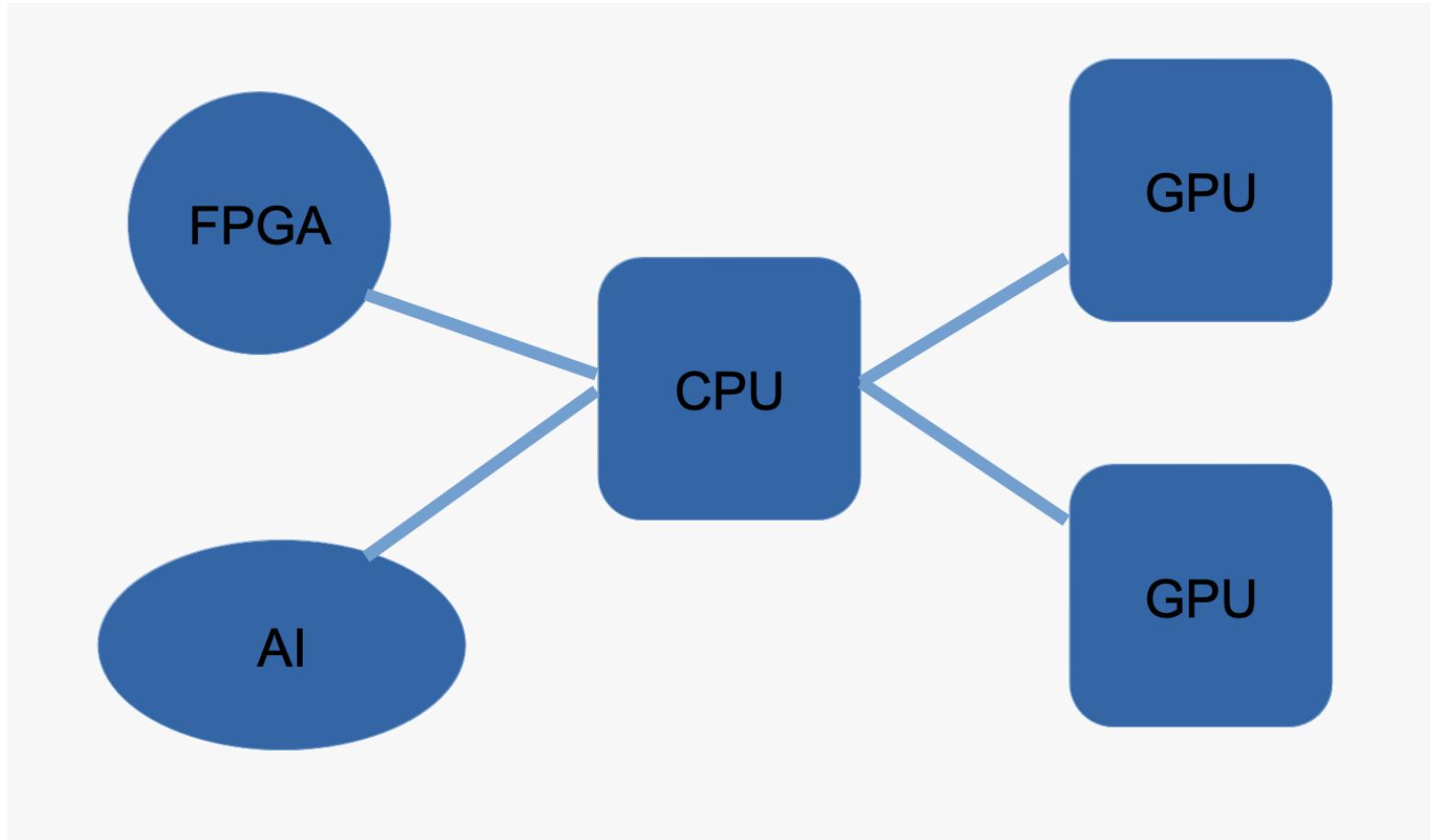
Reduce power  
by adding  
cores of small  
frequency

2 cores at half  
frequency  $f$   
can achieve  
the same  
throughput as a  
single core of  
frequency  $f$



**Dual Core**

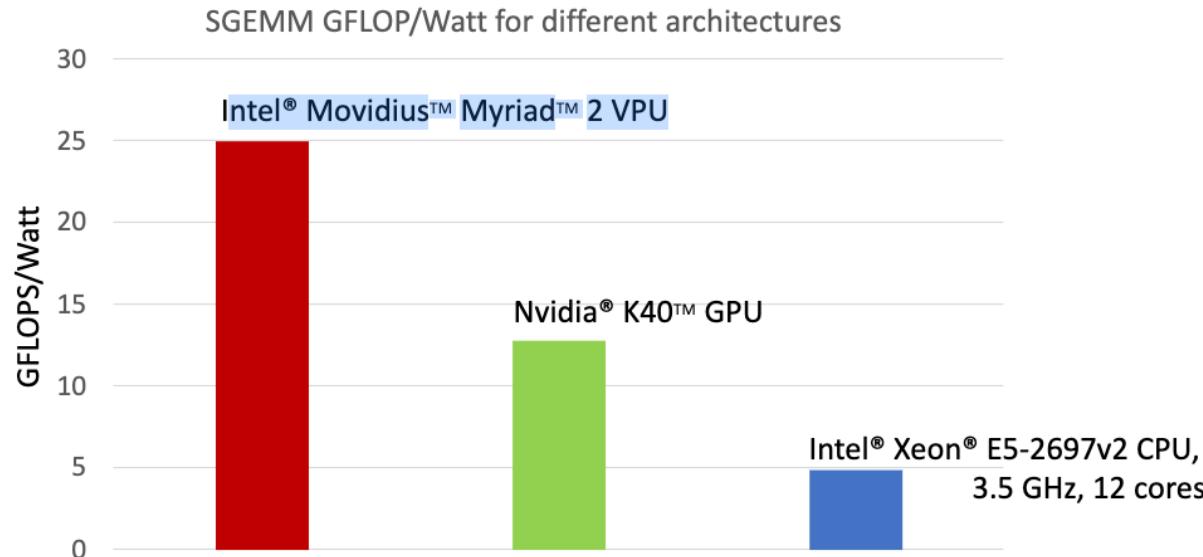
# Hence the rise of heterogeneous computing



# Why care about accelerators?

If you care about power, the world is heterogeneous?

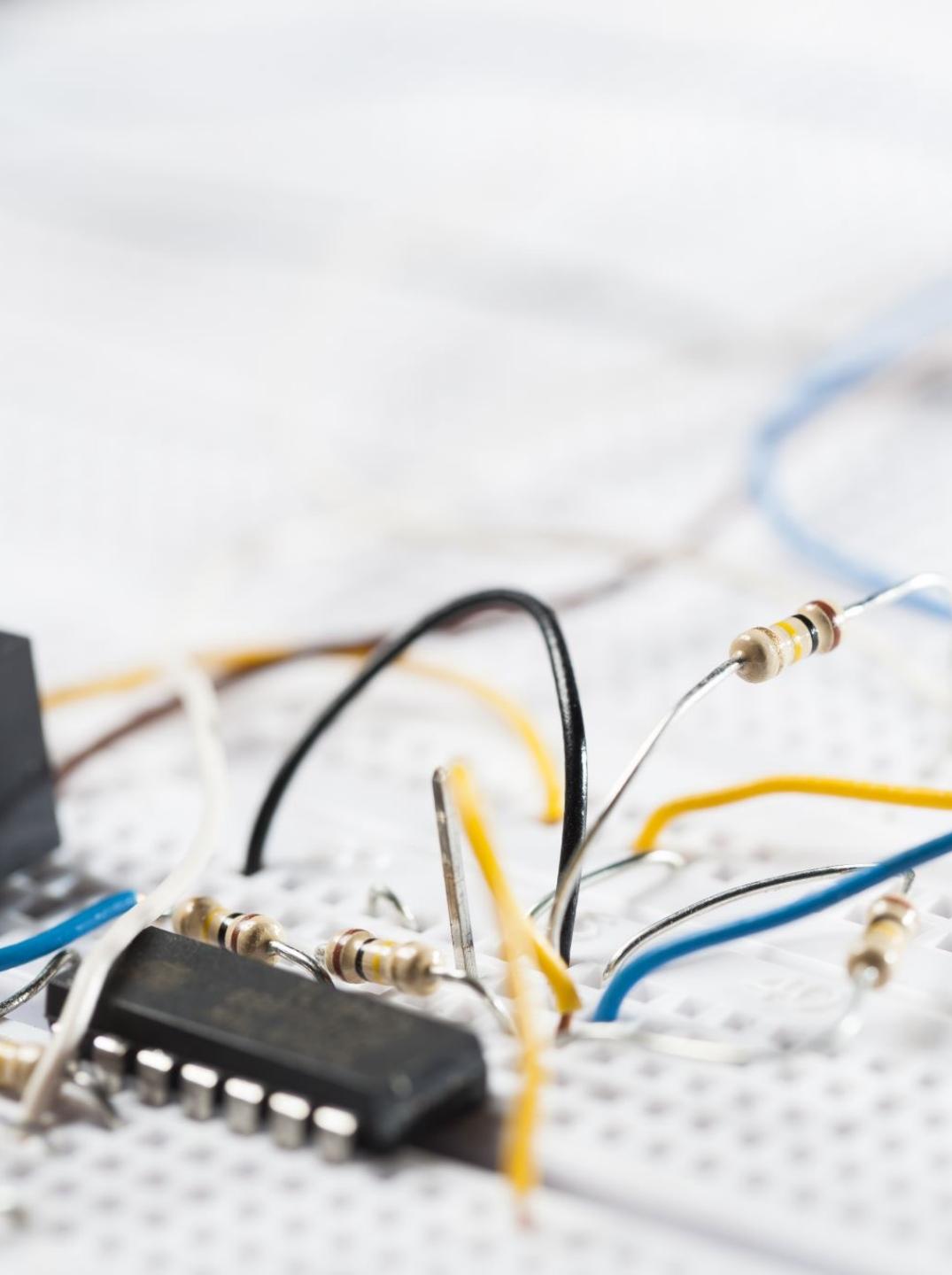
Specialized processors doing operations suited to their architecture are more efficient than general purpose processors.



Hence, future systems will be increasingly heterogeneous ... GPUs, CPUs, FPGAs, and a wide range of accelerators

Source: Suyash Bakshi and Lennart Johnsson, "A Highly Efficient SGEMM Implementation using DMA on the Intel/Movidius Myriad-2. IEEE International Symposium on Computer Architecture and High Performance Computing, 2020

Source: <https://extremecomputingtraining.anl.gov/wp-content/uploads/sites/96/2023/08/ATPESC-2023-Track-2b-Matteson-Openmp.pdf>



---

## Important note!!

- If you are looking for Deep learning acceleration on GPU please refer to the Horovod workshop given on 10/24/23.
- It covered Parallelization across multiple GPUs/nodes using Horovod Tensorflow/Pytorch
- [https://github.com/kf-cuanschutz/CU-Anschutz-HPC-documentation/blob/main/Workshops/Introduction\\_to\\_Horovod\\_102423\\_part1\\_official\\_v2.pdf](https://github.com/kf-cuanschutz/CU-Anschutz-HPC-documentation/blob/main/Workshops/Introduction_to_Horovod_102423_part1_official_v2.pdf)

## Historical context

- GPU are graphics processing units.



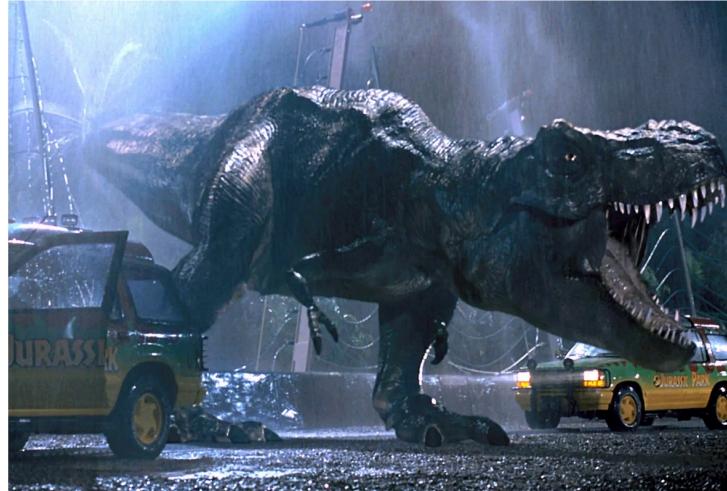
## Historical context

- GPU are graphics processing units.
- Originally came from the video game and movie industry.



# Historical context

- GPU are graphics processing units.
- Originally came from the video game and movie industry.



Jurassic Park

# Historical context

- GPU are graphics processing units.
- Originally came from the video game and movie industry.

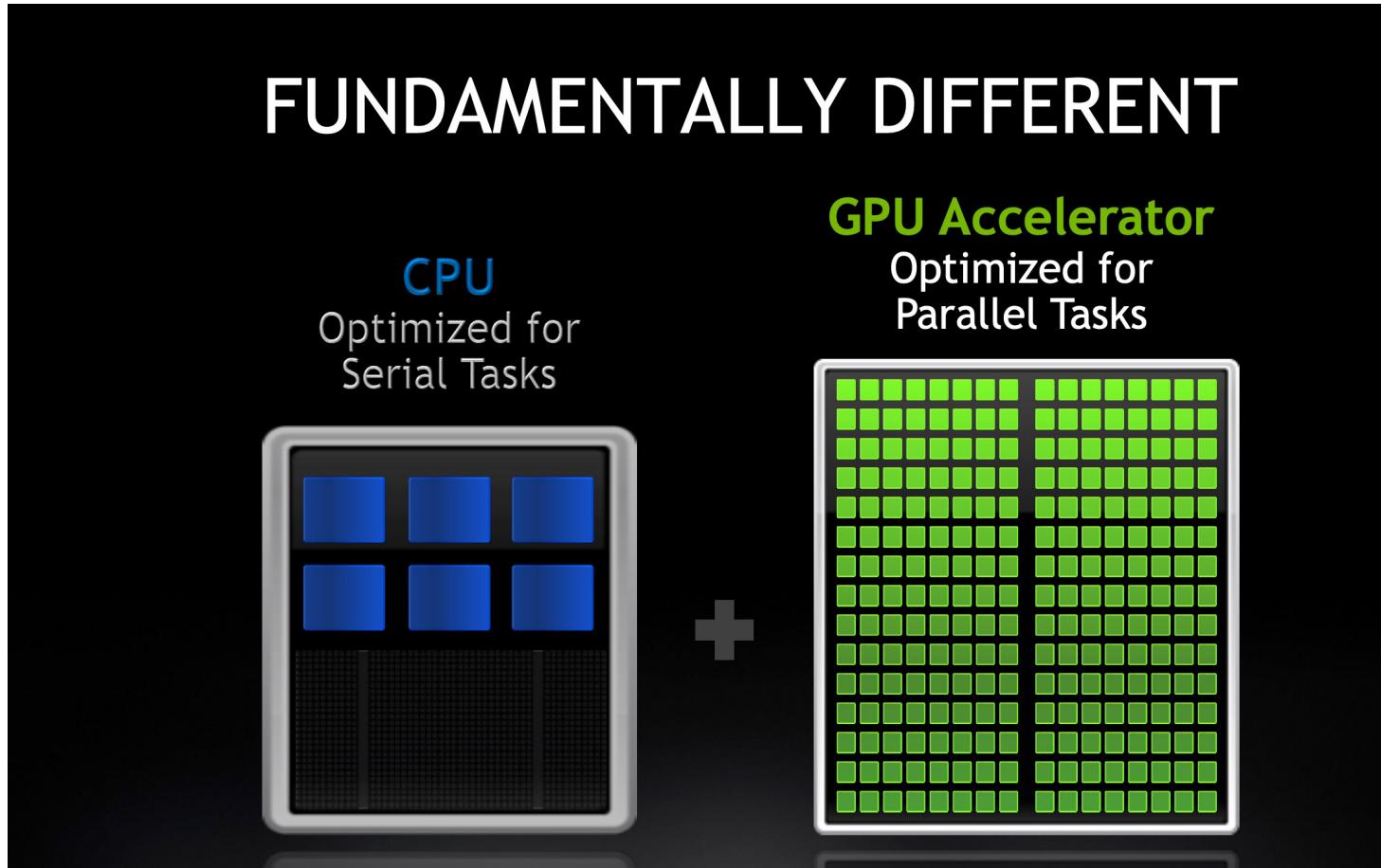


Jurassic Park

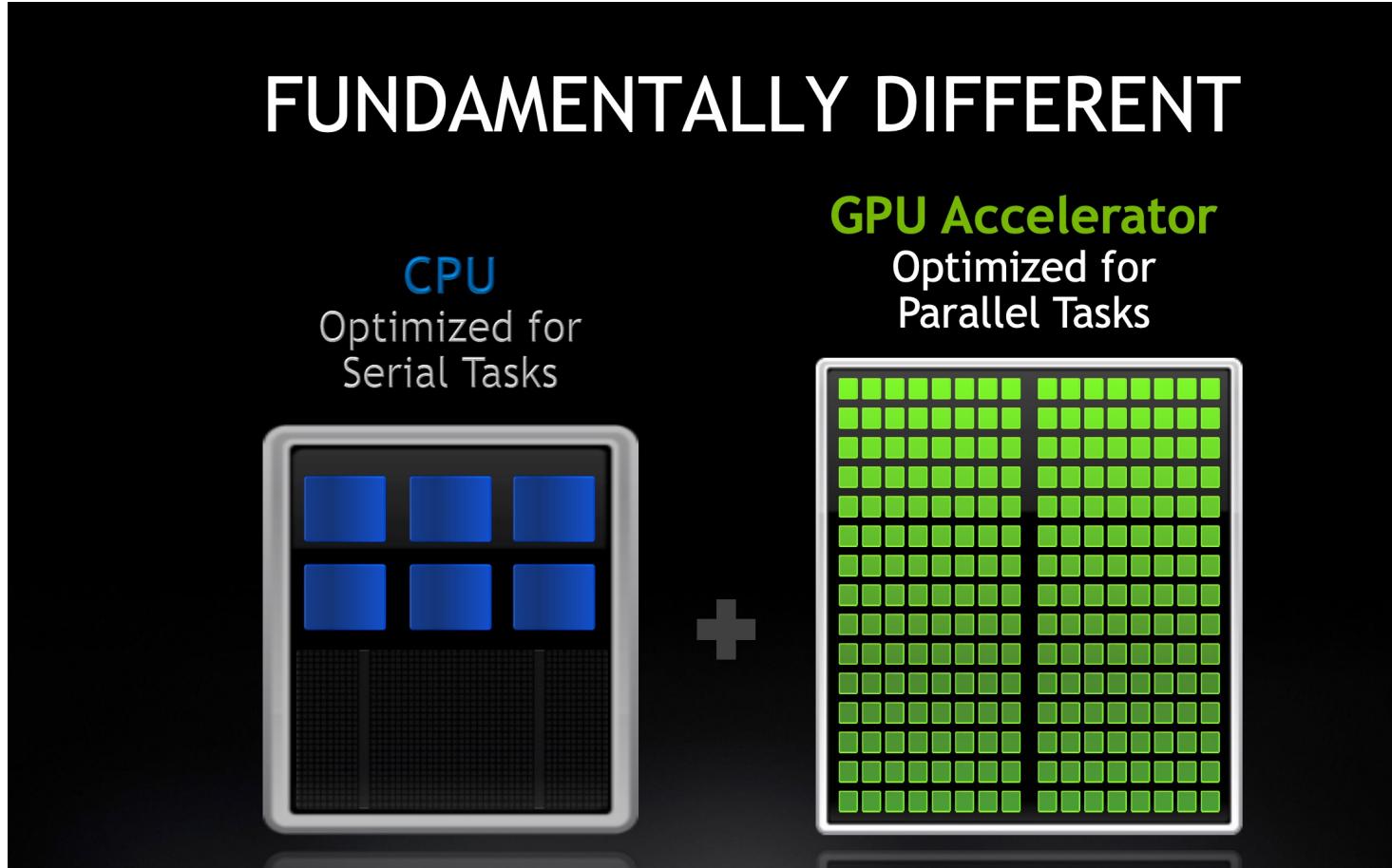


Zelda Ocarina of Time

# CPUs vs GPUs comparison

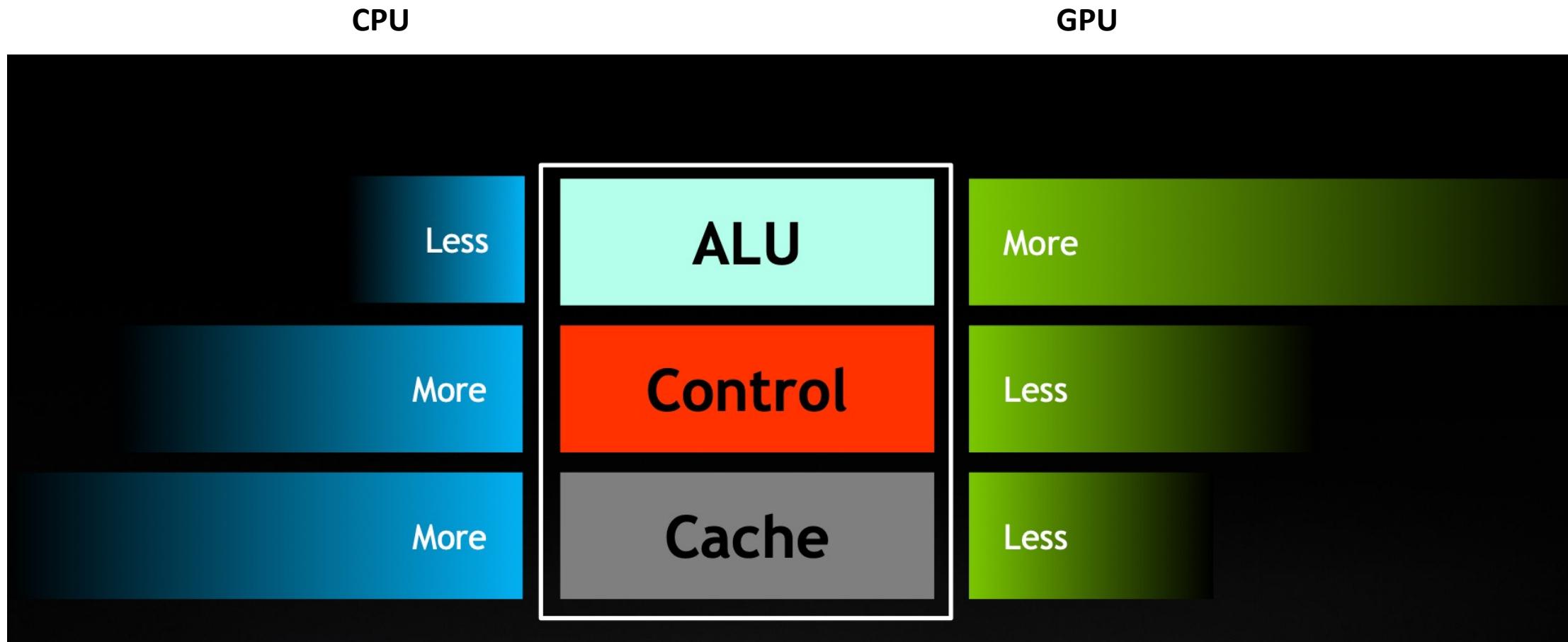


# CPUs vs GPUs comparison

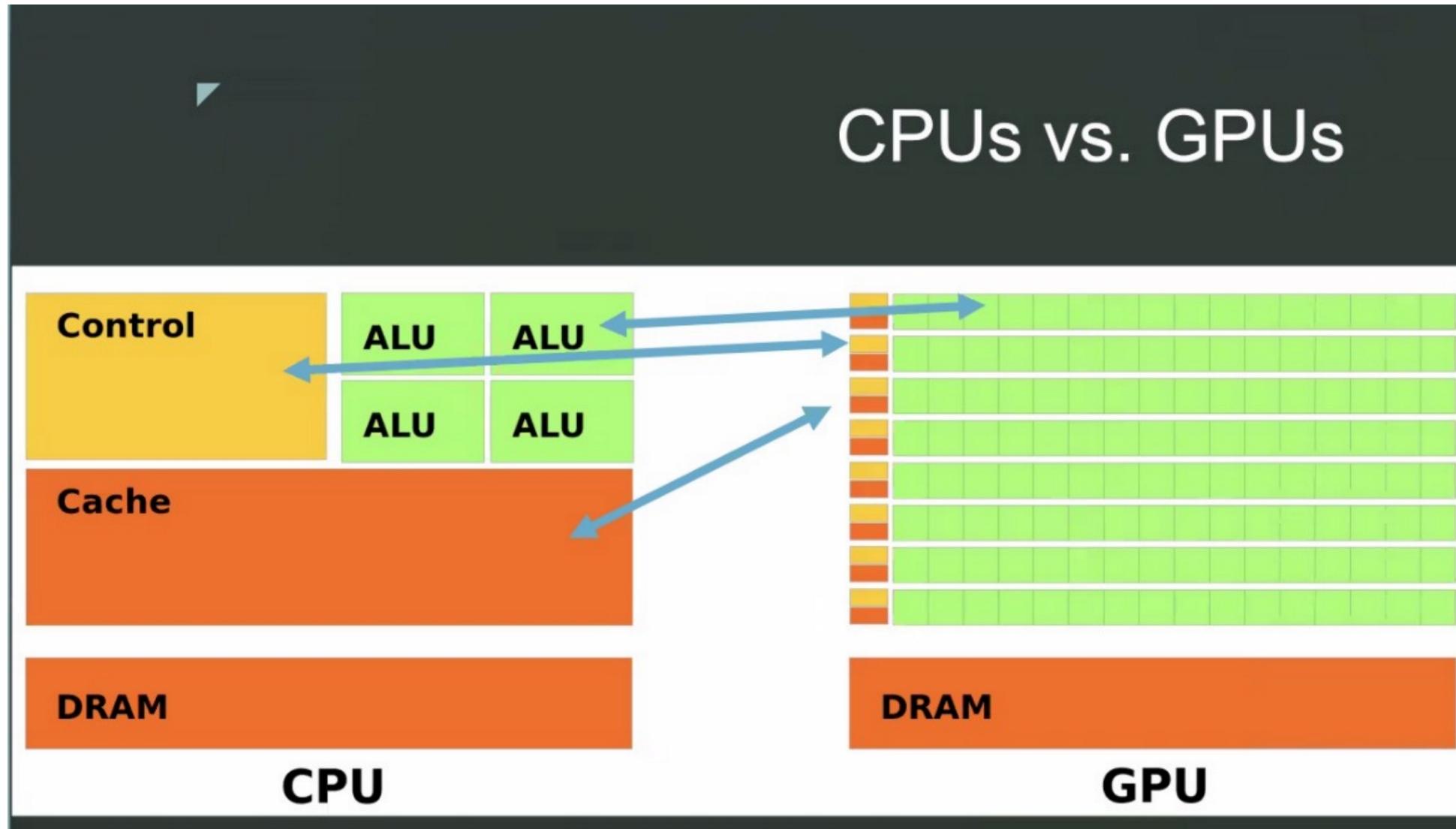


GPU follow SIMD pattern (Single Instruction multiple threads)

# CPUs vs GPUs comparison



# CPU vs GPU comparison



# CPU vs GPU comparison

- DRAM is the global memory
- CPUs have few cores (ALU) when GPUs have a large number of cores.

# CPU vs GPU comparison

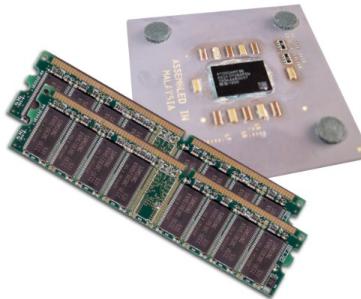
- DRAM is the global memory
- CPUs have few cores (ALU) when GPUs have a large number of cores.

|                     | Base Clock speed | Number of cores |
|---------------------|------------------|-----------------|
| Tesla V100 GPU      | 1.45 GHz         | 5120            |
| Intel Core i9-9900K | 3.6 GHz          | 8               |

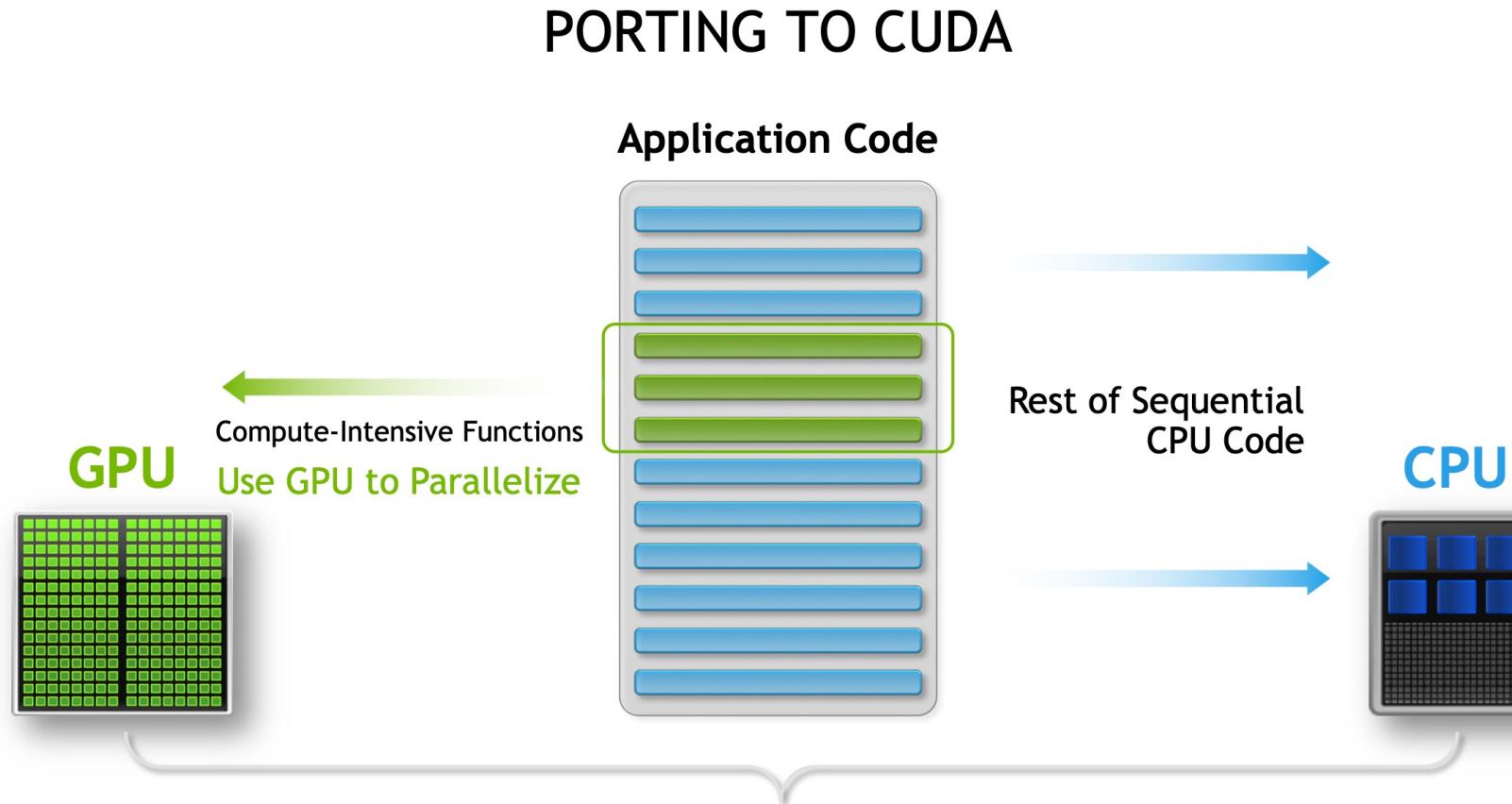
# CPU vs GPU programming model

## HETEROGENEOUS COMPUTING

- ▶ **Host** The CPU and its memory (host memory)
- ▶ **Device** The GPU and its memory (device memory)



# CPU vs GPU programming model

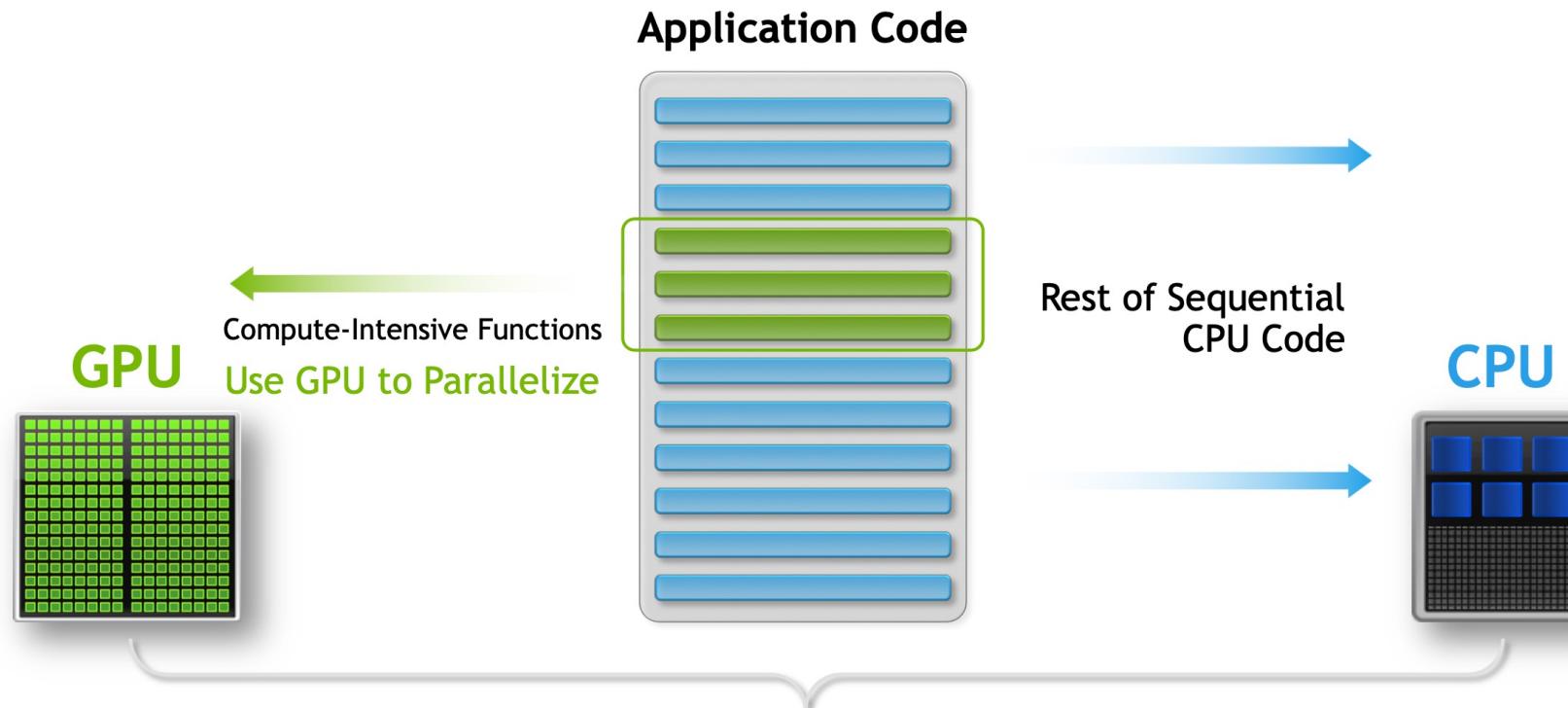


# CPU vs GPU programming model

**GPU:** Matrix operations. (e.g. FFT)

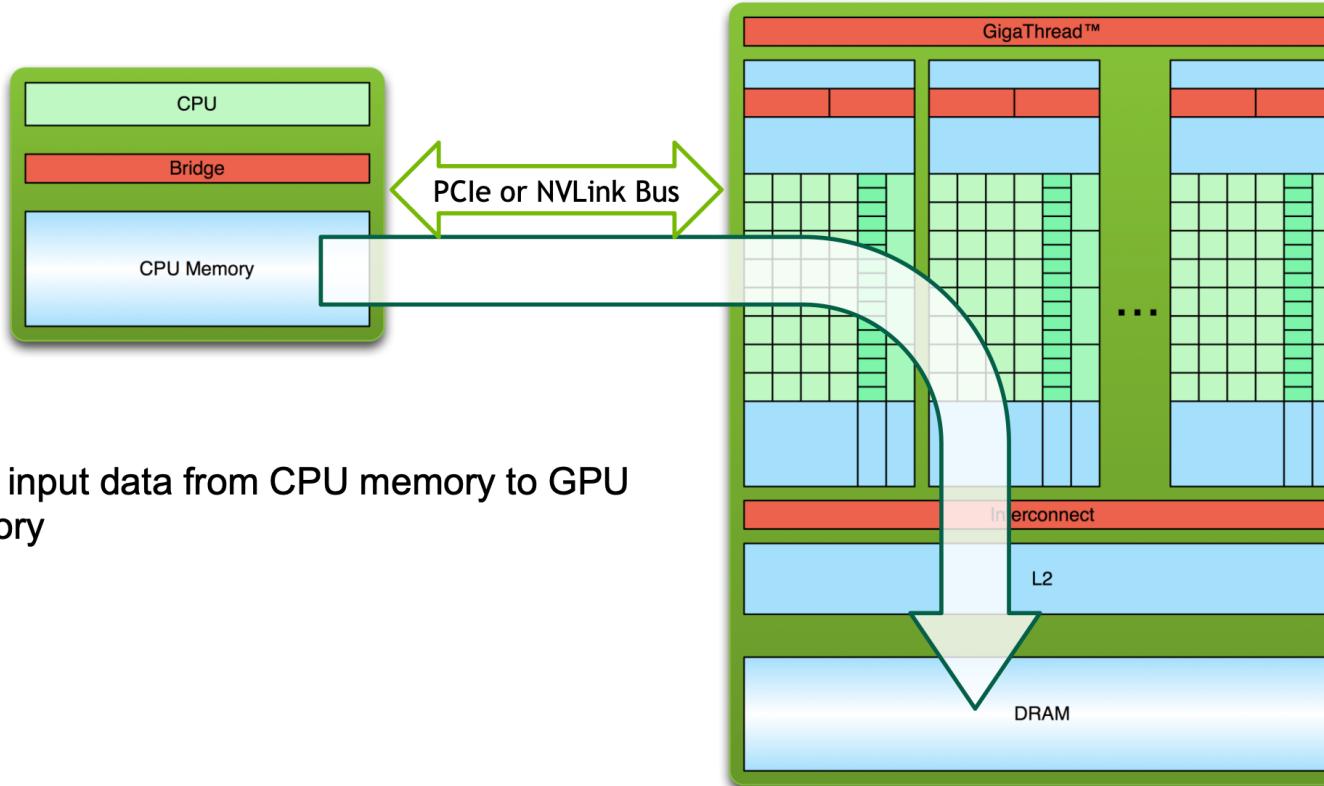
## PORTING TO CUDA

**CPU:**  
OS, security etc ...



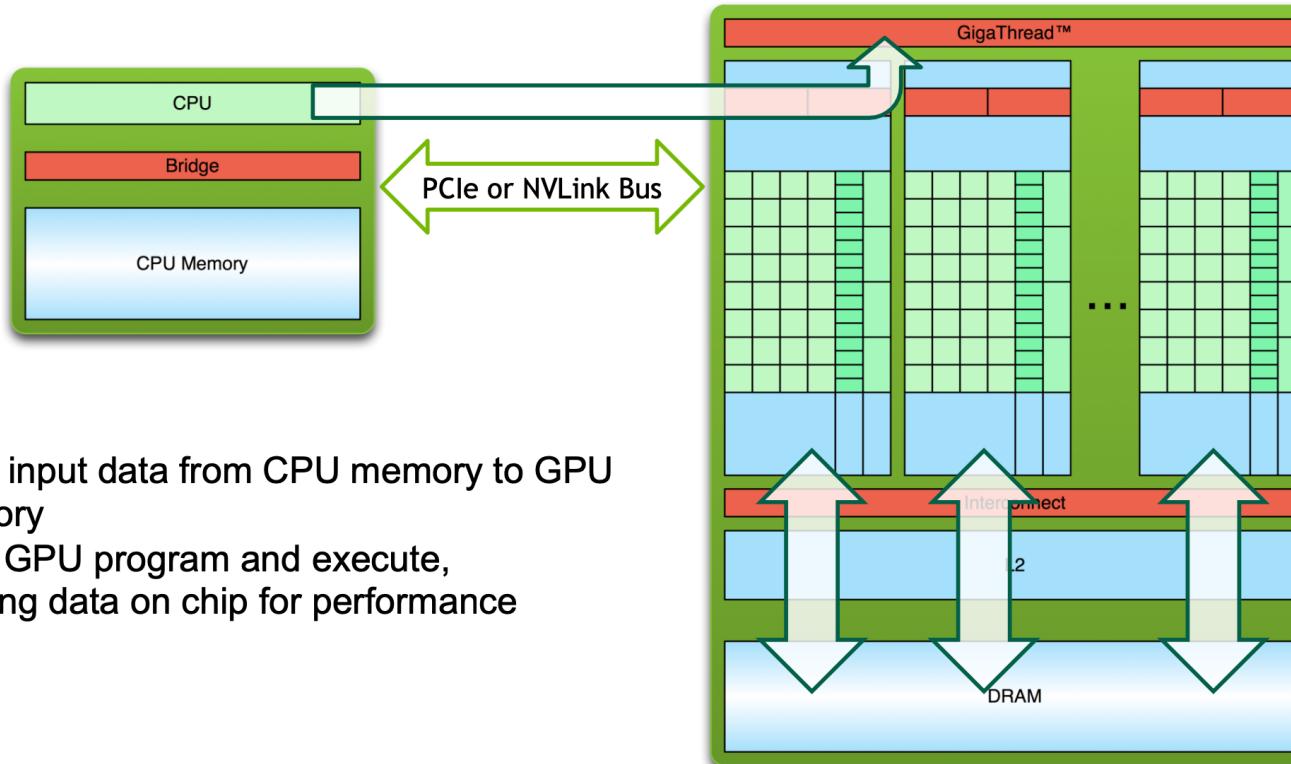
# CPU vs GPU programming model

## SIMPLE PROCESSING FLOW



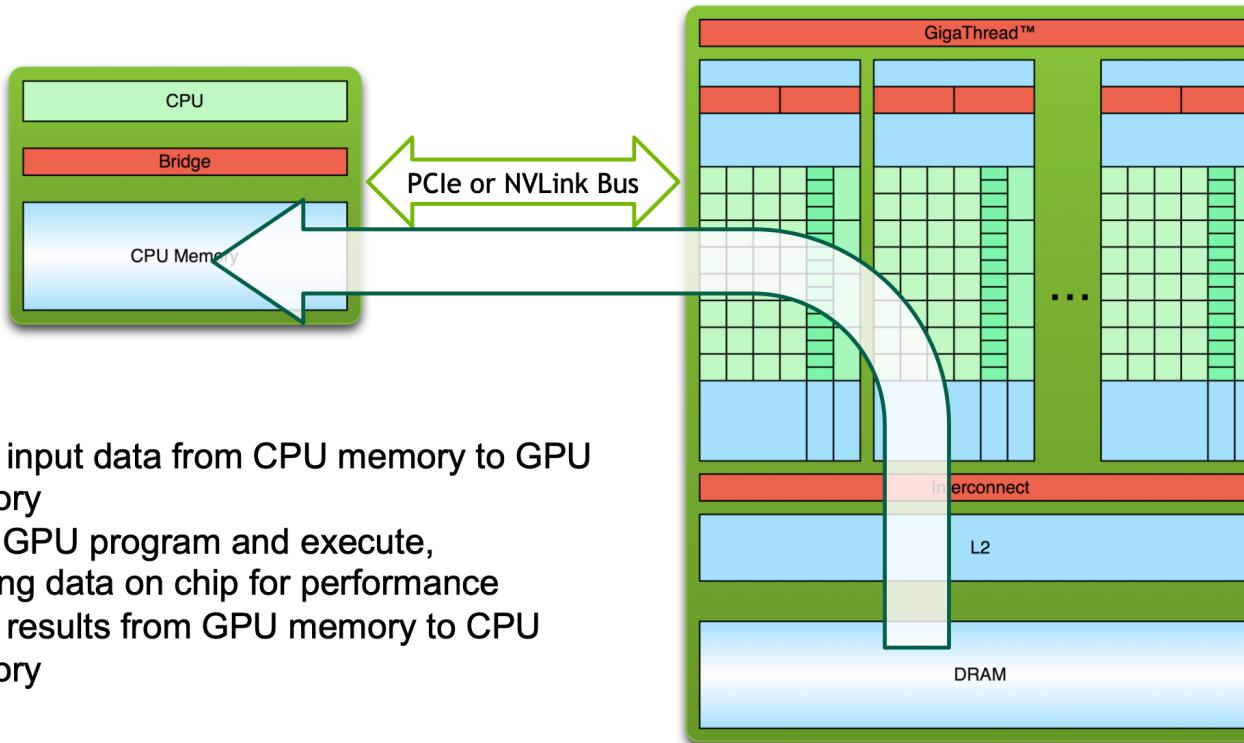
# CPU vs GPU programming model

## SIMPLE PROCESSING FLOW



# CPU vs GPU programming model

## SIMPLE PROCESSING FLOW



# CPU vs GPU programming model

## CPU IS A LATENCY REDUCING ARCHITECTURE



### GPU Accelerator

Optimized for

- Very large main memory
- Very fast clock speeds
- Latency optimized via large caches
- Small number of threads can run very quickly

### CPU Weaknesses

- Relatively low memory bandwidth
- Cache misses very costly
- Low performance/watt

# CPU vs GPU programming model

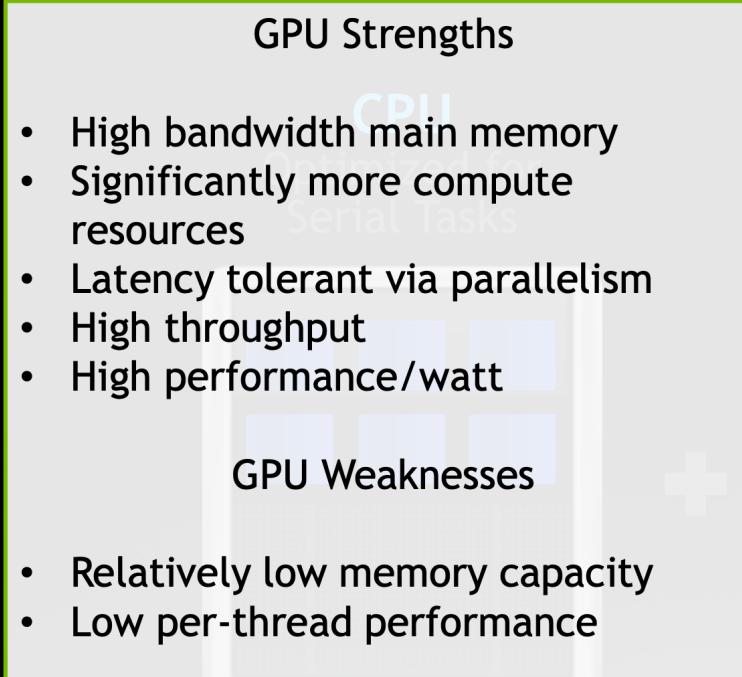
## GPU IS ALL ABOUT HIDING LATENCY

**GPU Strengths**

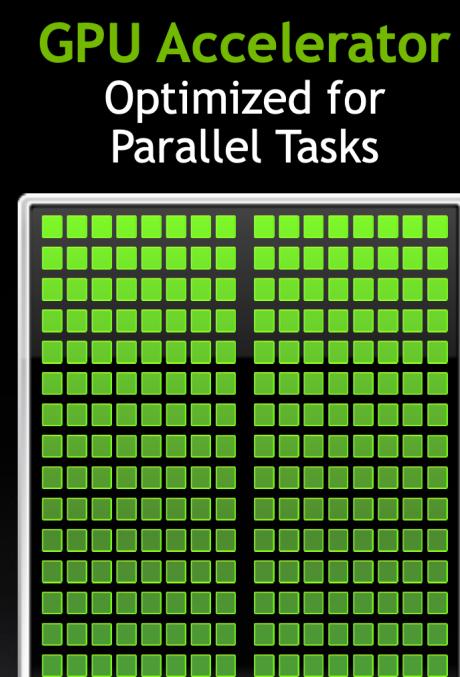
- High bandwidth main memory
- Significantly more compute resources
- Latency tolerant via parallelism
- High throughput
- High performance/watt

**GPU Weaknesses**

- Relatively low memory capacity
- Low per-thread performance



**GPU Accelerator**  
Optimized for Parallel Tasks



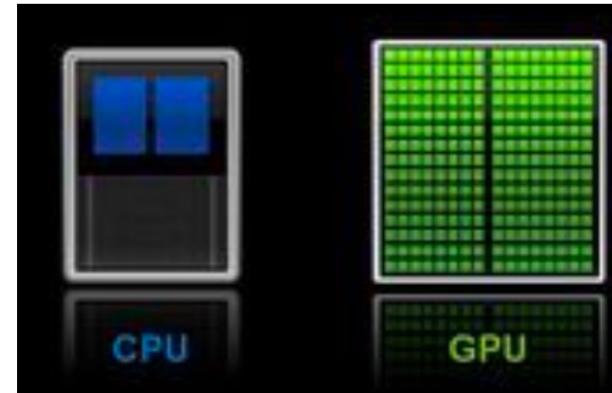
# CPU vs GPU programming model

- You need lots of parallelism with GPU cores to optimize computation.
- GPU memory is very fast but moving data from CPU to GPU is very costly.

# CPU vs GPU programming model

## CPU (Haswell)

- 128GB DDR
- ~120 GB/Sec Memory Bandwidth



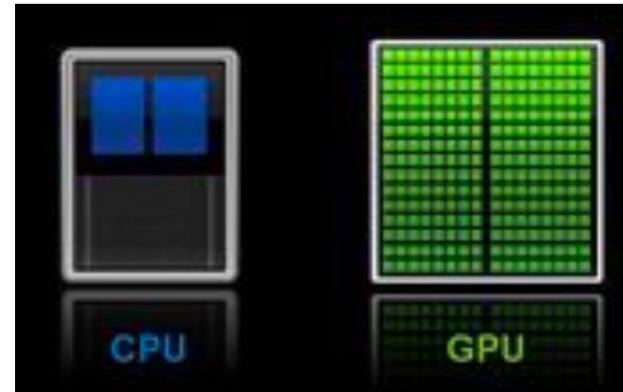
## GPU (A100)

- 40GB HBM
- 1,500 GB/Sec Memory Bandwidth

# CPU vs GPU programming model

## CPU (Haswell)

- 128GB DDR
- ~120 GB/Sec Memory Bandwidth



## GPU (A100)

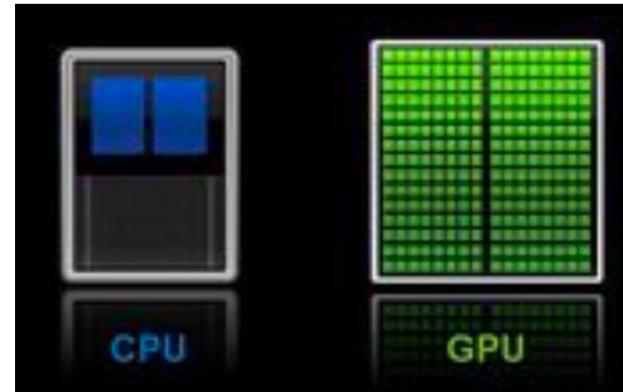
- 40GB HBM
- 1,500 GB/Sec Memory Bandwidth



# CPU vs GPU programming model

## CPU (Haswell)

- 128GB DDR
- ~120 GB/Sec Memory Bandwidth



## GPU (A100)

- 40GB HBM
- 1,500 GB/Sec Memory Bandwidth



CPU - Speed



GPU - Throughput



# NVIDIA A100 specs



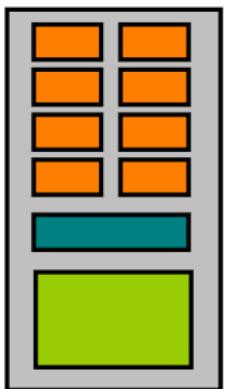
# NVIDIA A100 specs



Scalar  
Processor



cuda core.



streaming  
multiprocessor (SMs).

Multiprocessor

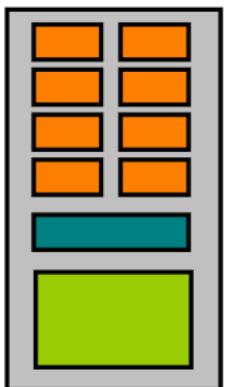
# NVIDIA A100 specs



Scalar  
Processor



cuda core.



streaming multiprocessor  
**: 64 cuda cores per SM**

Multiprocessor

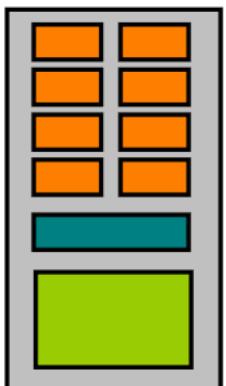
# NVIDIA A100 specs



Scalar  
Processor

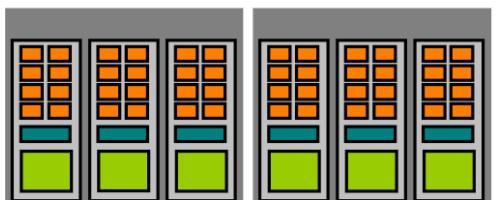


cuda core. **64 per SM**



streaming multiprocessor  
**: 108 SMs**

Multiprocessor



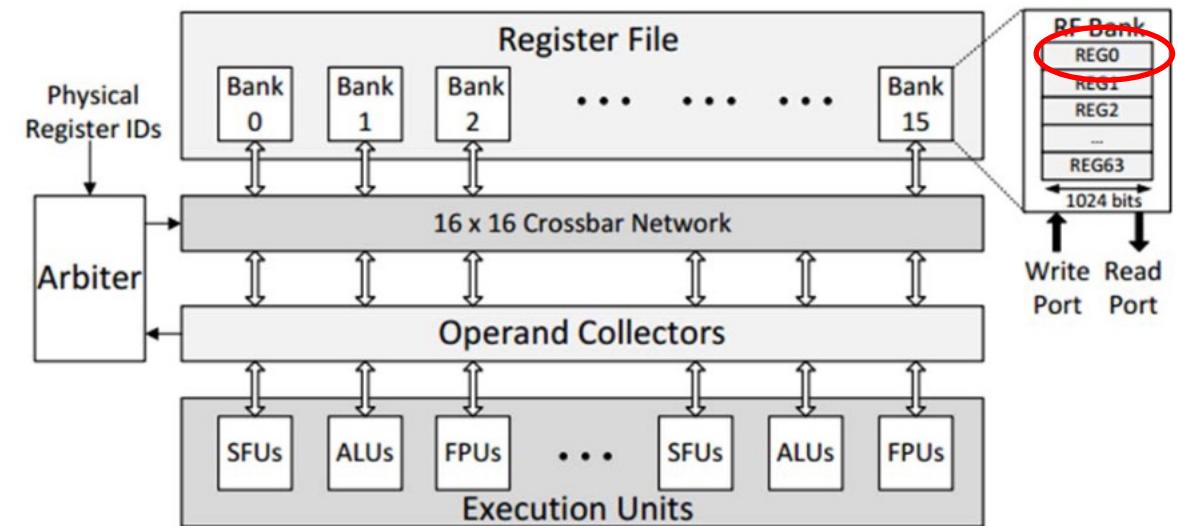
Device



**108 SMs per device**

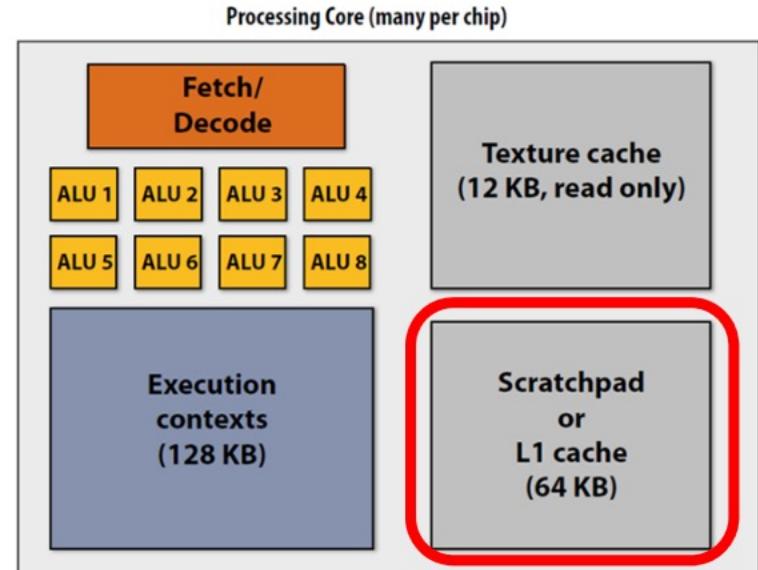
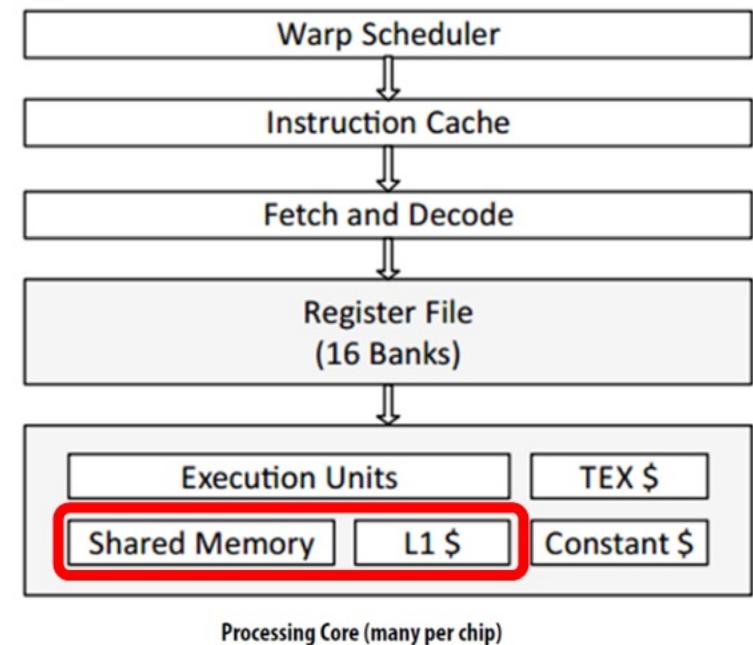
# NVIDIA GPU memory hierarchy

- Streaming multiprocessors registers
- Very fast (**every 1 clock cycle**)
- Only usable by 1 thread



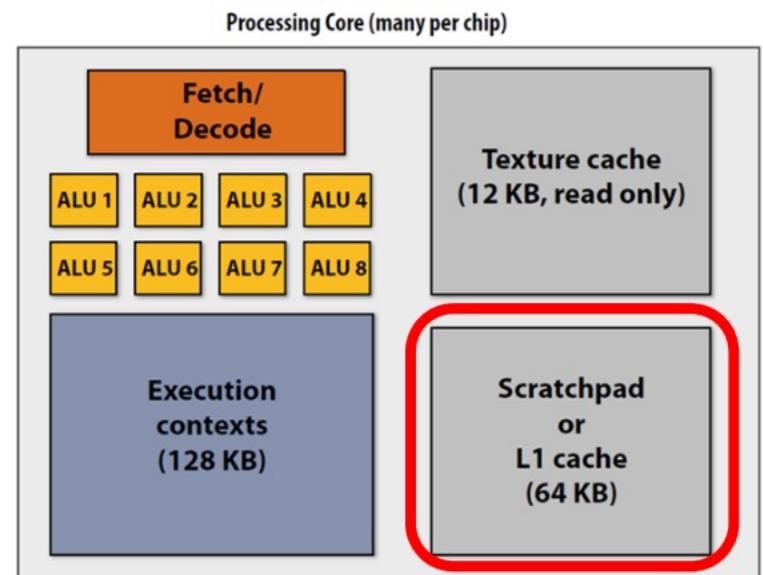
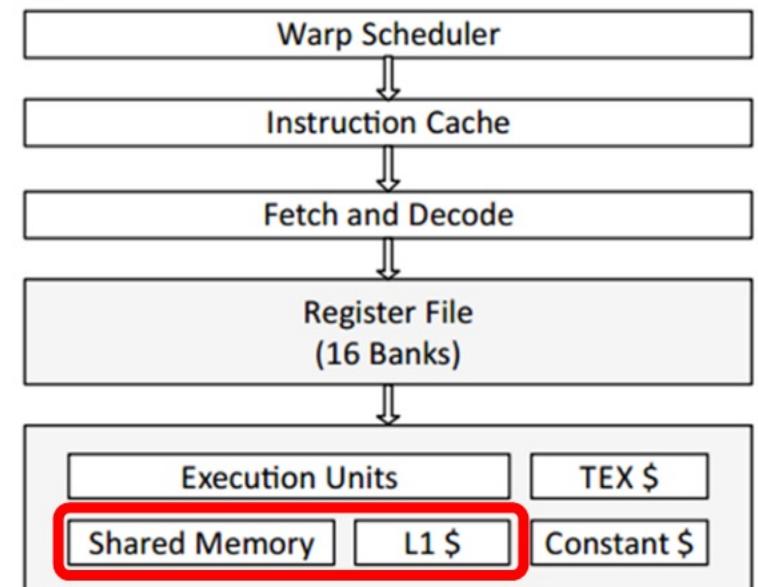
# NVIDIA GPU memory hierarchy

- Shared/L1 Memory
- Organized into banks that can be accessed simultaneously



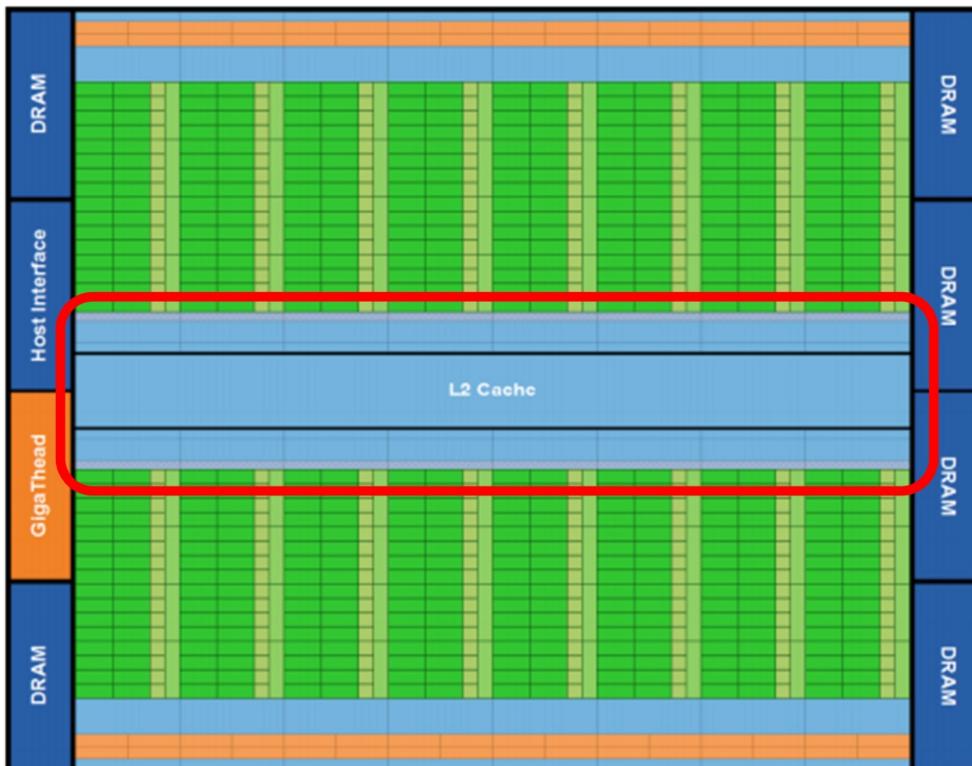
# NVIDIA GPU memory hierarchy

- Shared/L1 Memory
- Organized into banks that can be accessed simultaneously
- Located in an SM



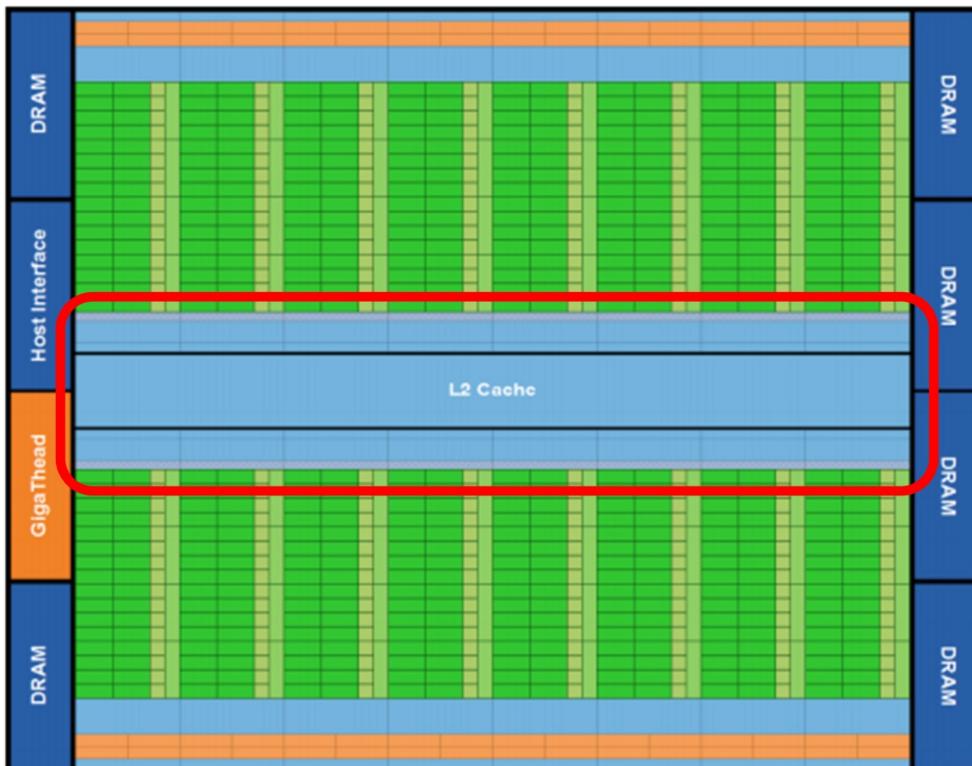
# NVIDIA GPU memory hierarchy

- L2 cache
- Shared among SMs



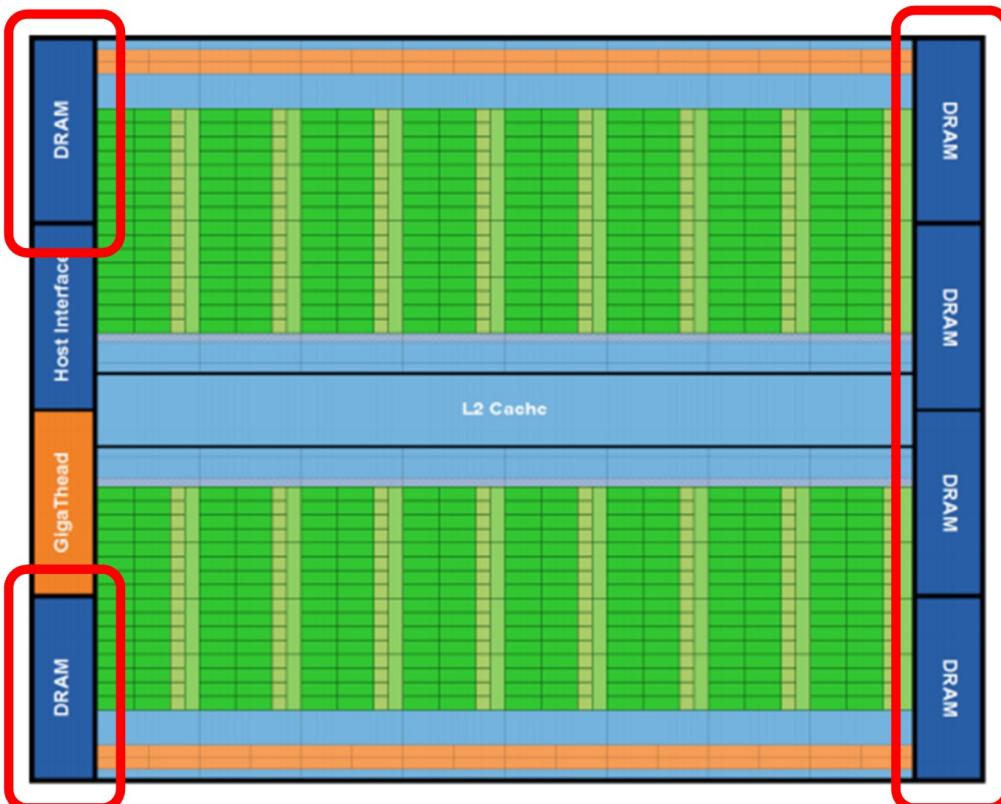
# NVIDIA GPU memory hierarchy

- L2 cache
- Shared among SMs
- Slower bandwidth than L1



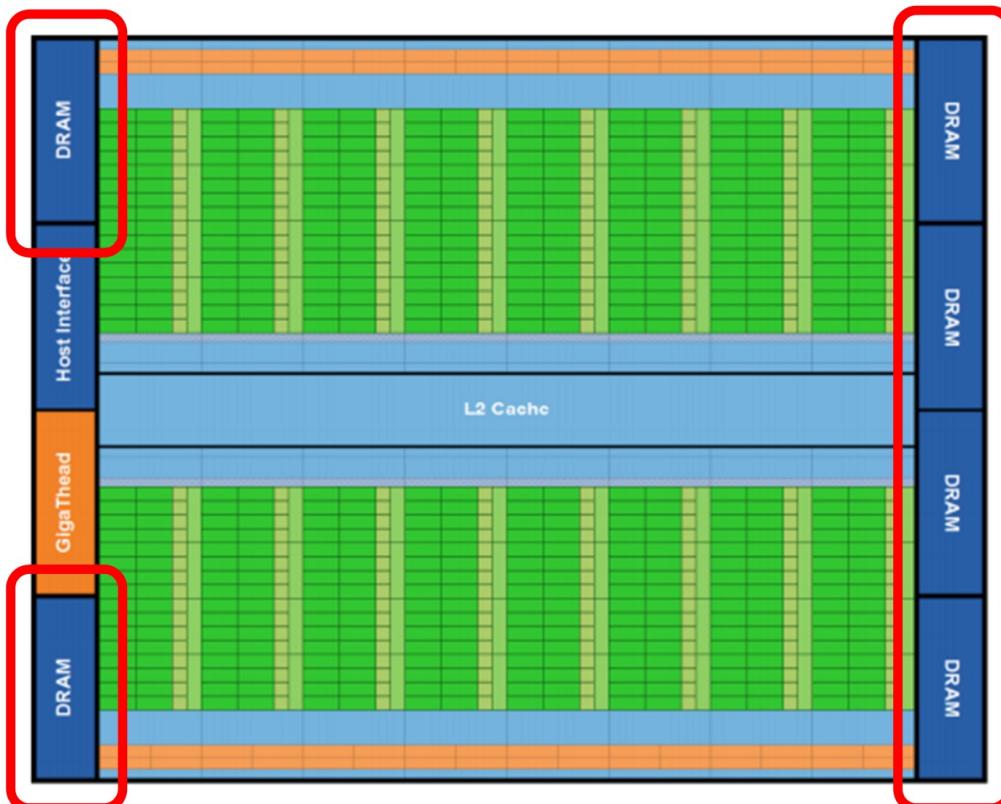
# NVIDIA GPU memory hierarchy

- DRAM
- Much slower than L2 (hundreds of clock cycles)



# NVIDIA GPU memory hierarchy

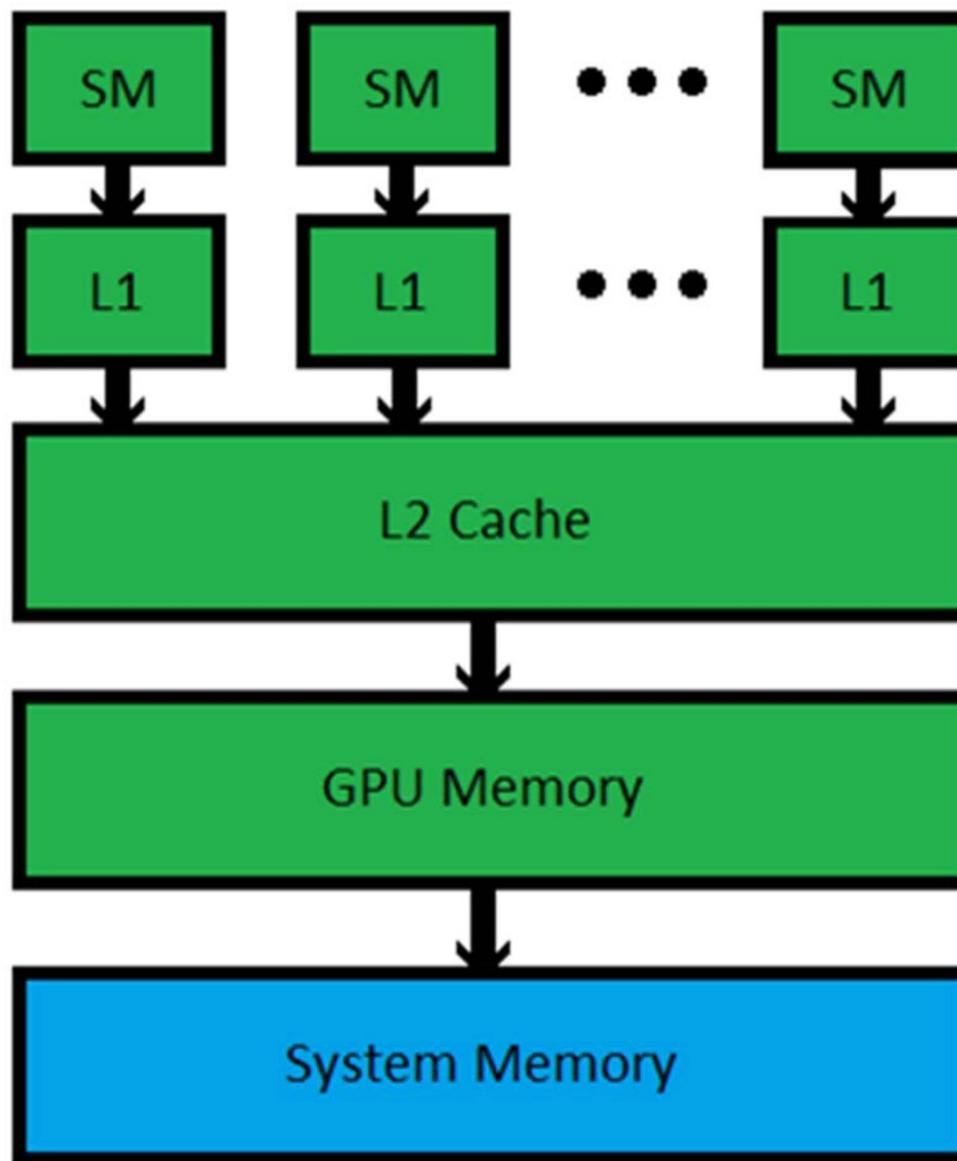
- DRAM
- Much slower than L2 (hundreds of clock cycles)
- Can be accessed by both CPU and GPU



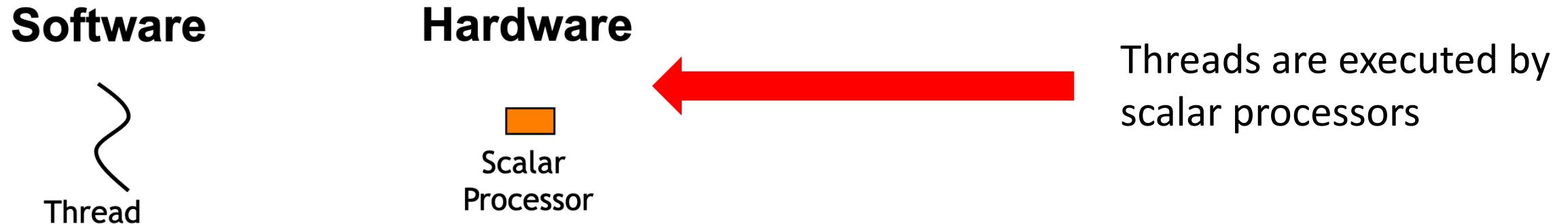
# NVIDIA GPU memory hierarchy

- PCIe (10-100 GB/s)

# NVIDIA GPU memory hierarchy



# GPU software programming



# GPU software programming

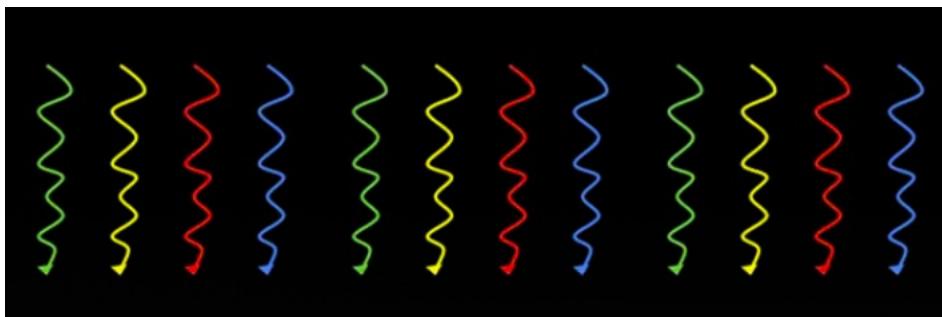
## Software



## Hardware



Threads are executed by scalar processors



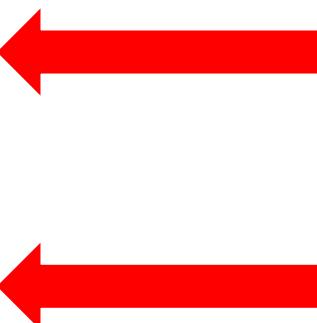
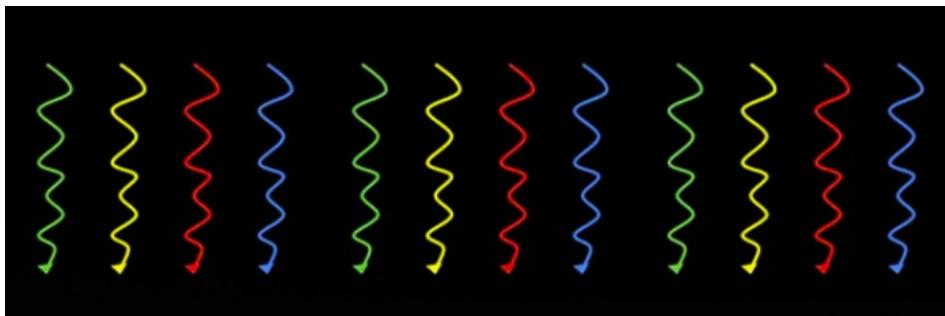
GPUs can handle thousands of concurrent threads

# GPU software programming

## Software



## Hardware



Threads are executed by scalar processors

GPUs can handle thousands of concurrent threads

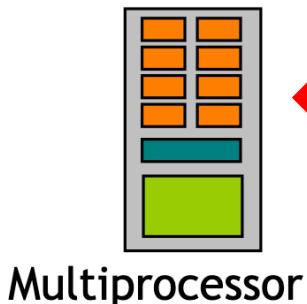
With CUDA, more threads than the GPU can handle concurrently can be launched.

# GPU software programming

## Software



## Hardware



Threads are executed by scalar processors

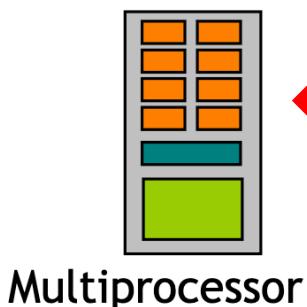
Threads blocks are executed on multiprocessors

# GPU software programming

## Software



## Hardware



Threads are executed by scalar processors

They do not migrate

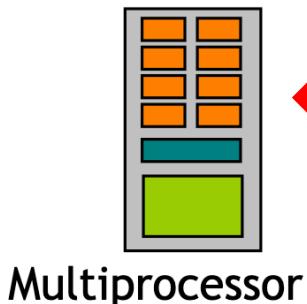


# GPU software programming

## Software



## Hardware



Threads are executed by scalar processors

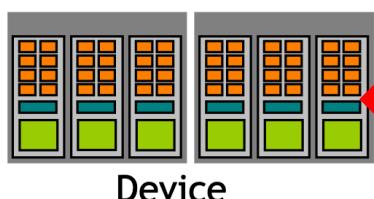
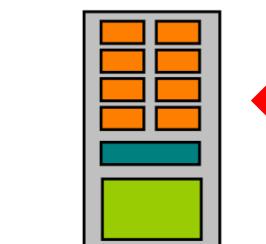
Concurrent thread block can reside on 1 multiprocessor

# GPU software programming

## Software



## Hardware

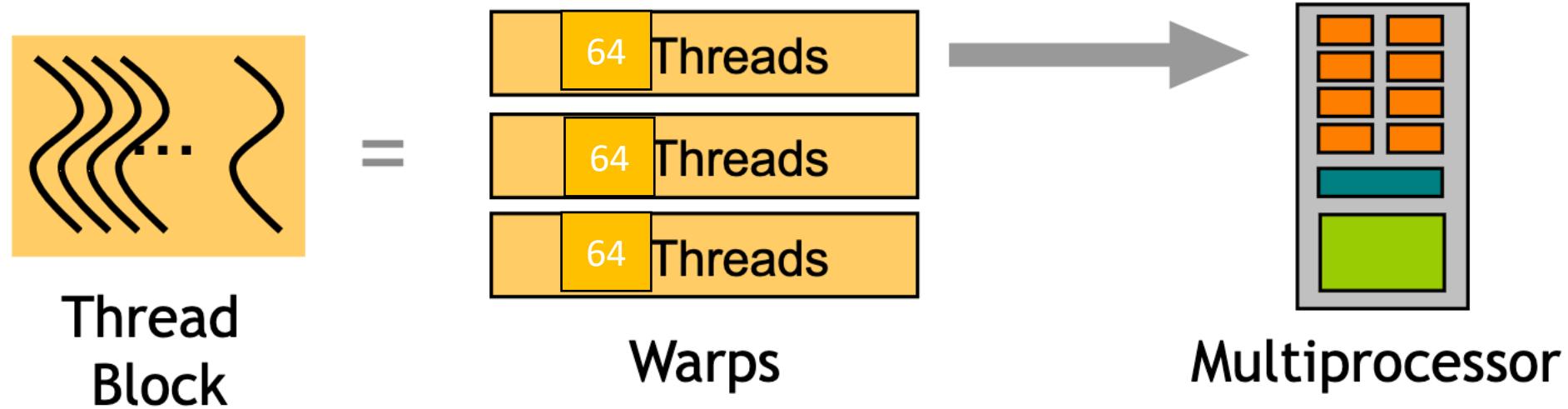


Threads are executed by scalar processors

Concurrent thread block can reside on 1 multiprocessor

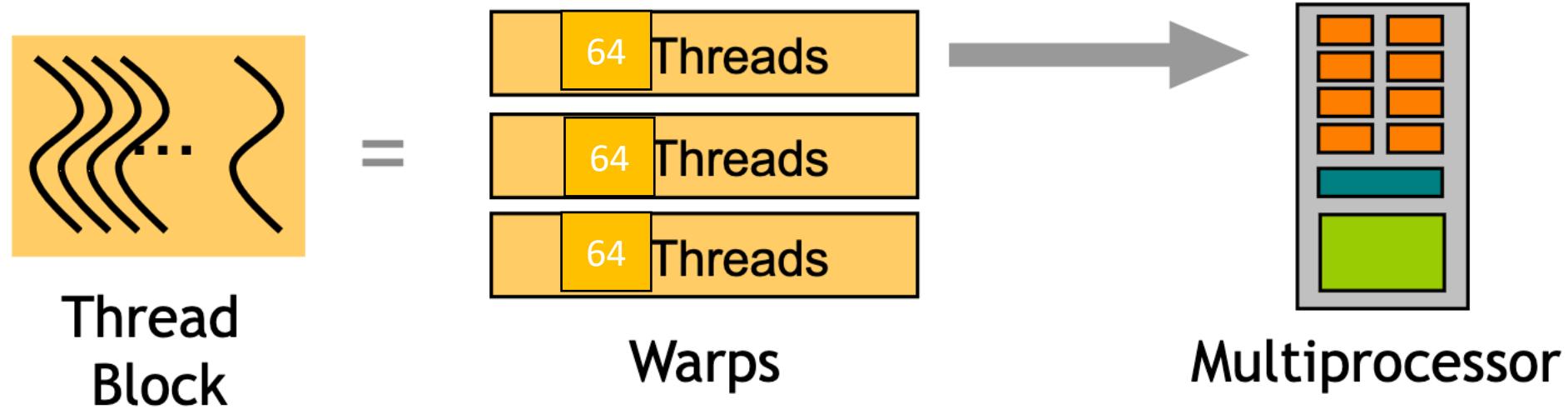
A kernel is launched as a grid of thread blocks

# Warps



- \* A thread block consists of 64 thread warps
- \* A warp is executed in parallel on a multiprocessor

# Warps

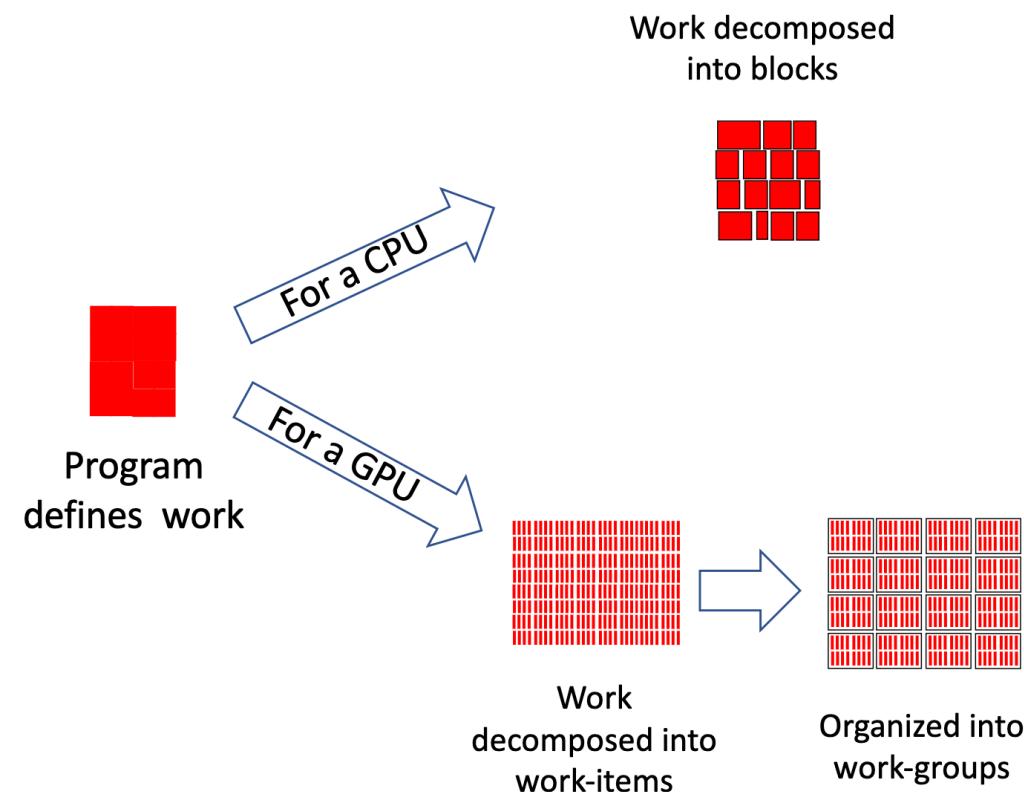


- \* A thread block consists of 64 thread warps

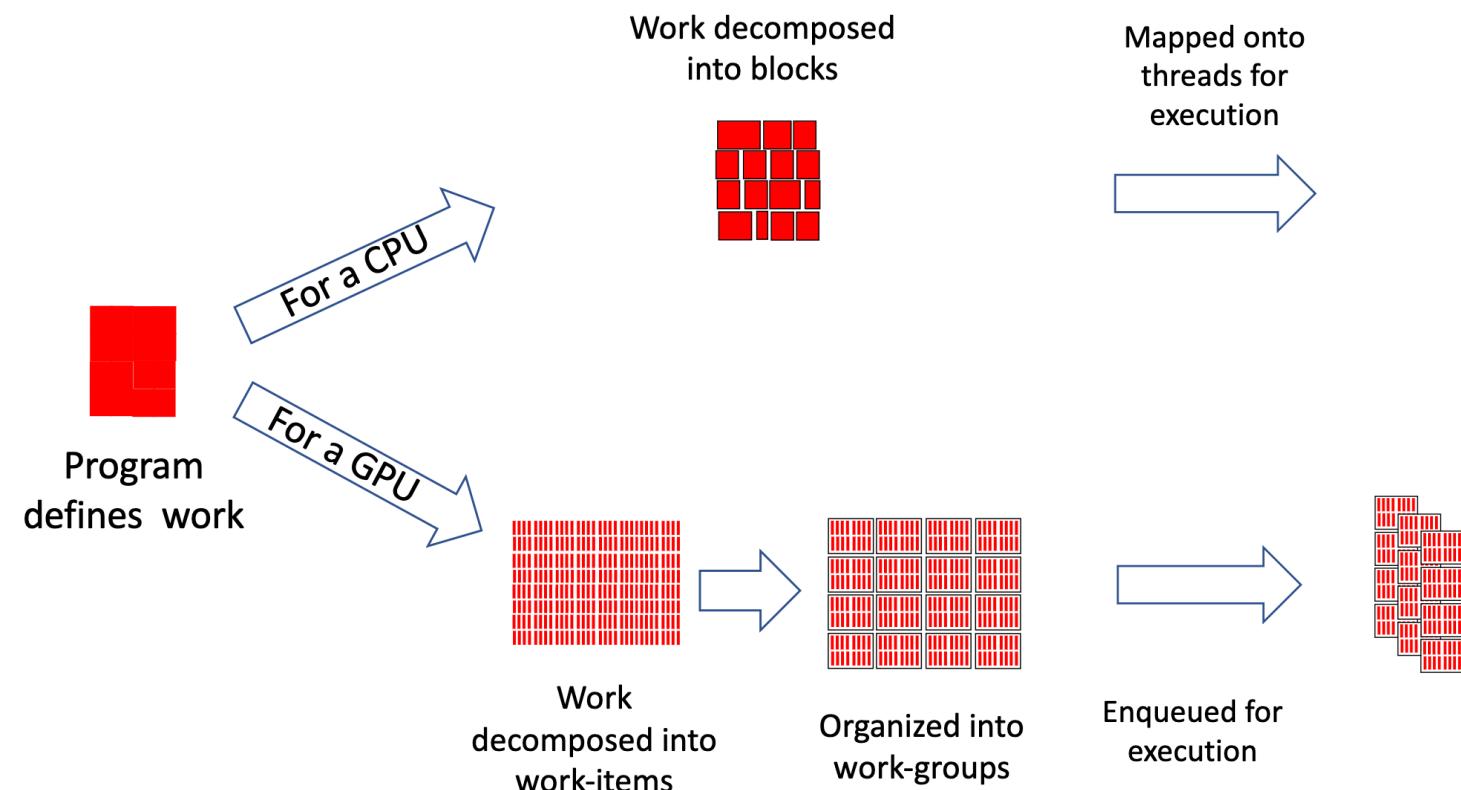
- \* A warp is executed in parallel on a multiprocessor

- \* Thread within the same warp execute the same operation at a given clock cycle.

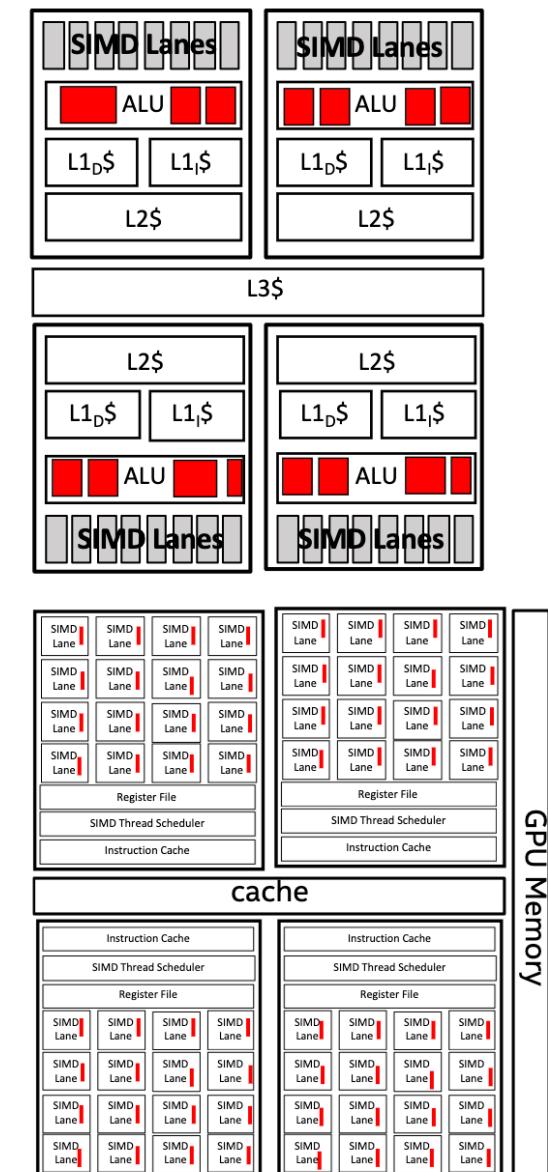
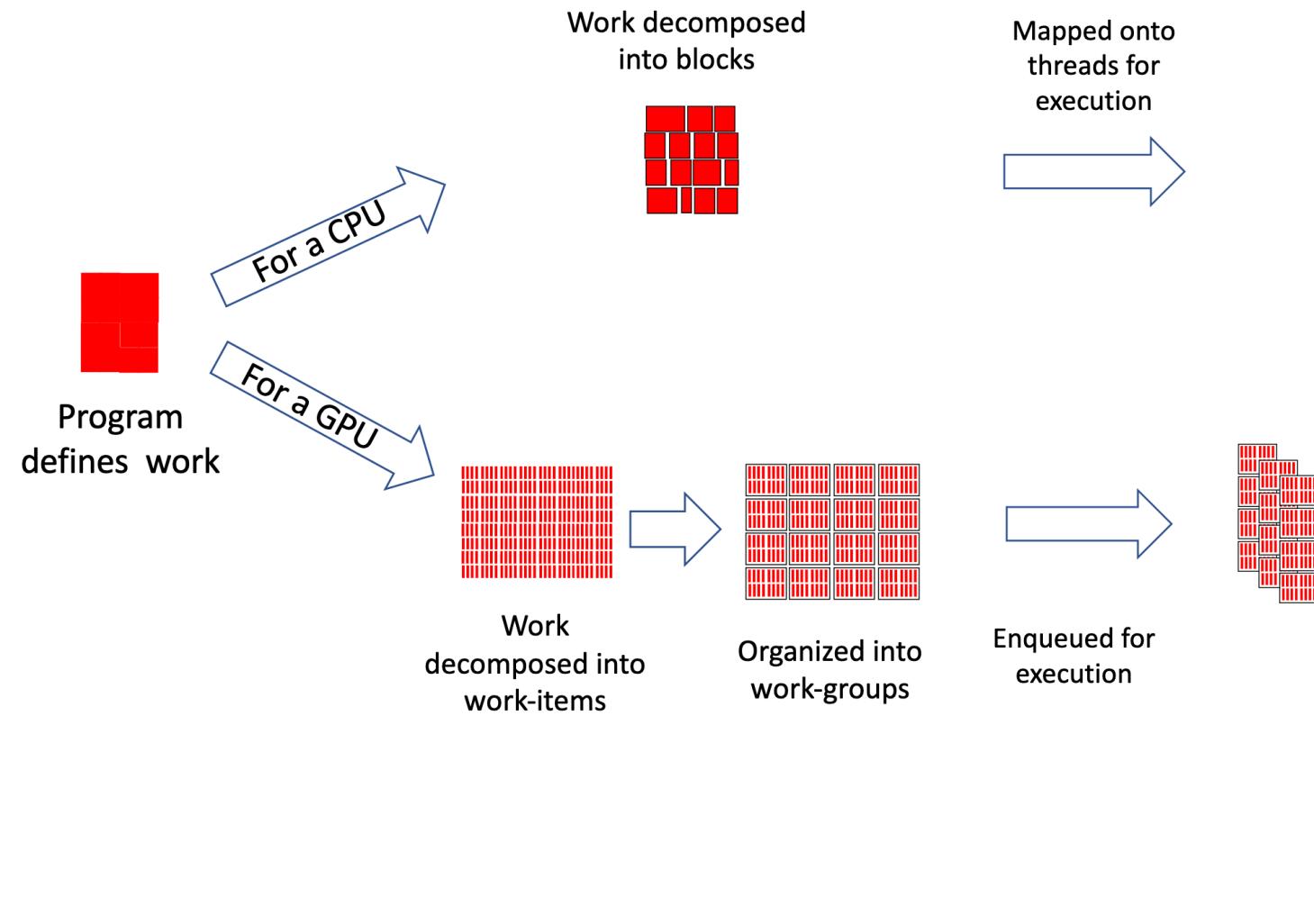
# Summary CPU vs GPU computation



# Summary CPU vs GPU computation



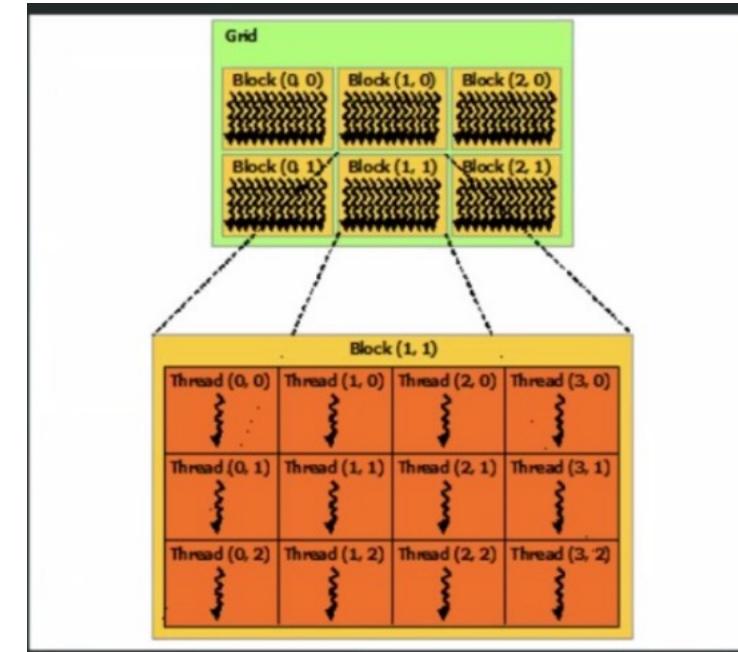
# Summary CPU vs GPU computation



One work-group per compute-unit executing

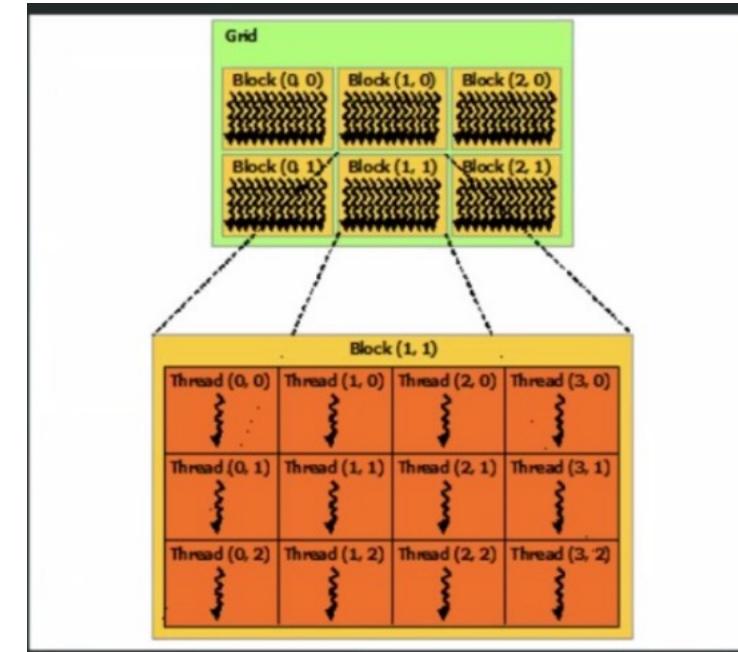
# CUDA terminology

- \* A kernel is the operation you want to run on your data which will be shipped to the GPU



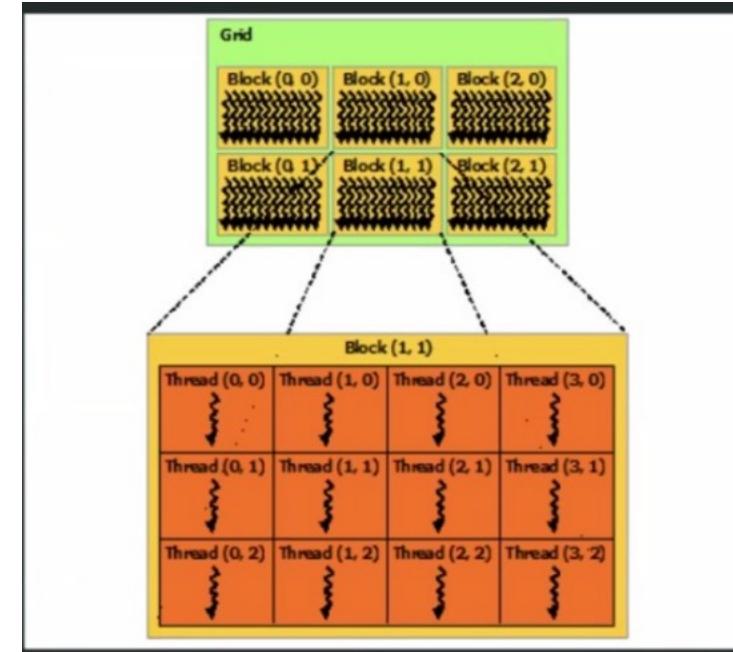
# CUDA terminology

- \* A kernel is the operation you want to run on your data which will be shipped to the GPU
- \* Kernels are launched in grid of blocks



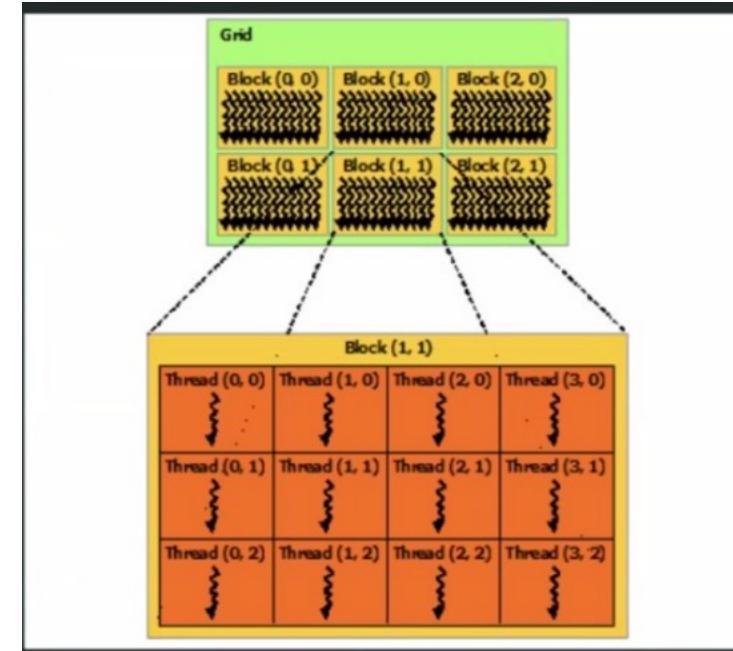
# CUDA terminology

- \* A kernel is the operation you want to run on your data which will be shipped to the GPU
- \* Kernels are launched in grid of blocks
- \* Grids and blocks can be 1D, 2D or 3D



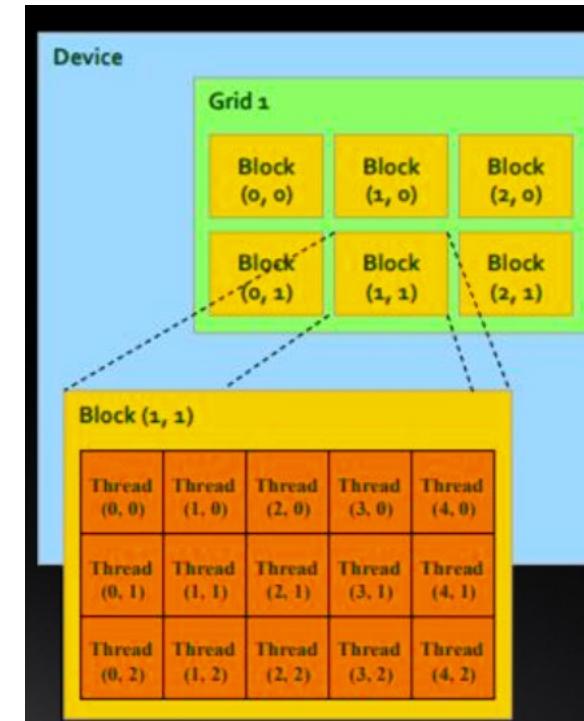
# CUDA terminology

- \* A kernel is the operation you want to run on your data which will be shipped to the GPU
- \* Kernels are launched in grid of blocks
- \* Grids and blocks can be 1D, 2D or 3D



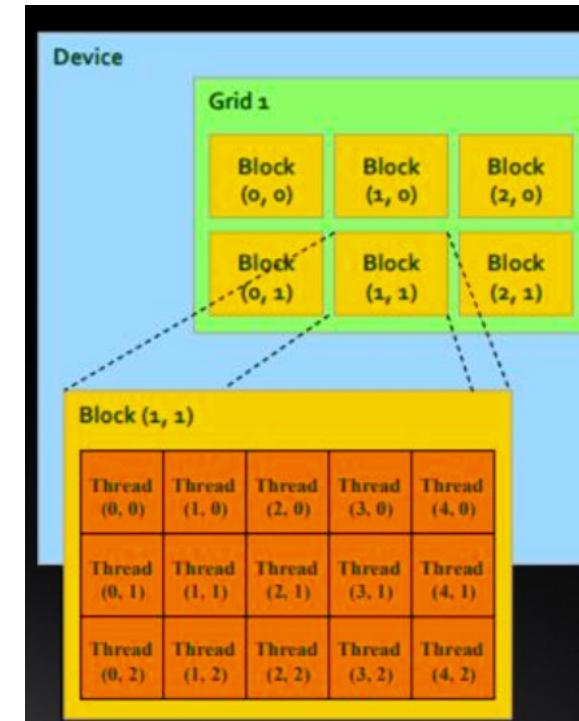
# Mapping Software to Hardware

- \* Threads have 3D IDs unique within block



# Mapping Software to Hardware

- \* Threads have 3D IDs unique within block
- \* Thread blocks have 2D ids unique within grid

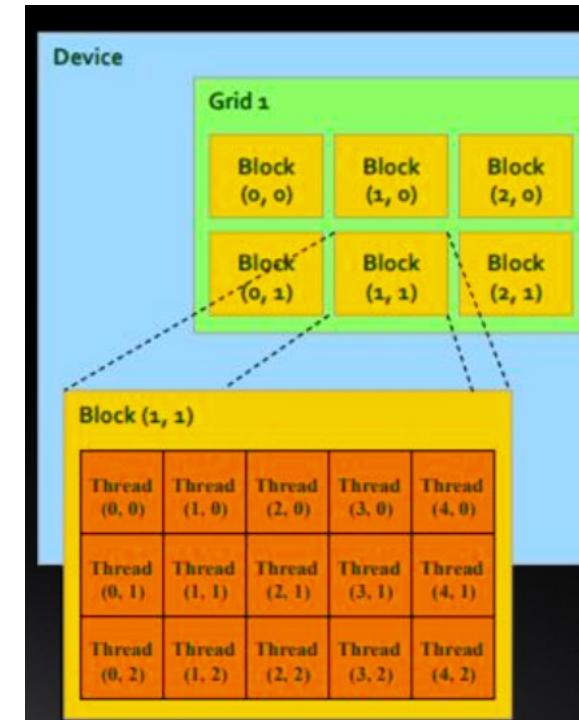


# Mapping Software to Hardware

- \* Threads have 3D IDs unique within block

- \* Thread blocks have 2D ids unique within grid

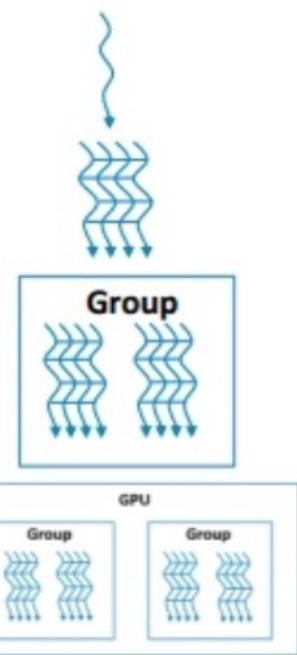
- \* Built-in variables”  
**threadIdx, blockIdx**  
**blockDim, gridDim**



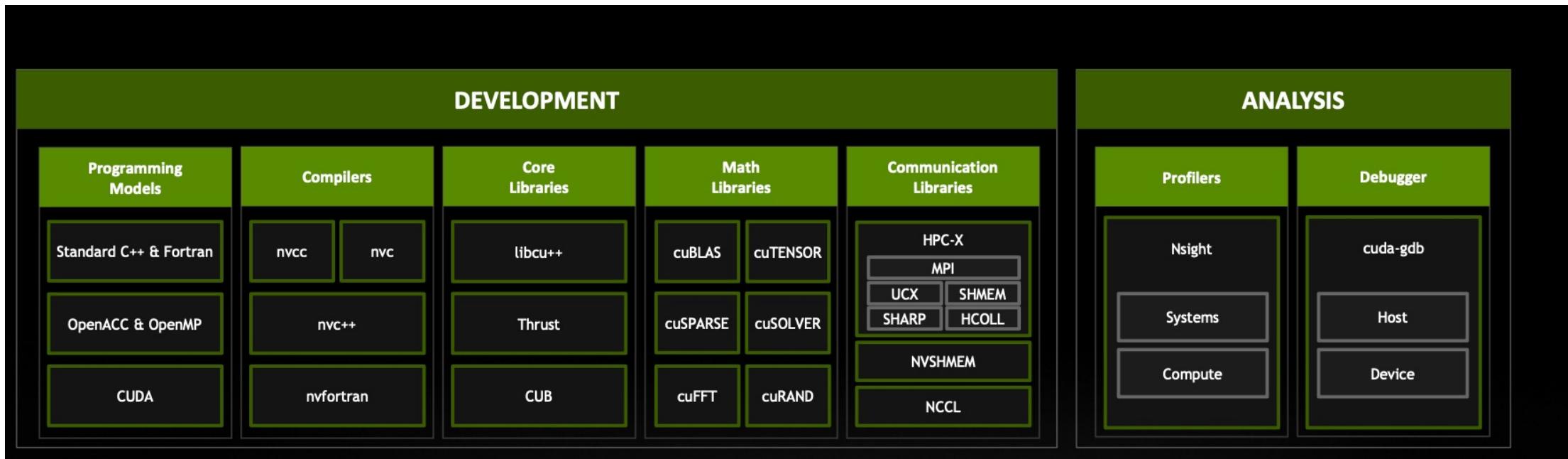
# GPU terminology varies a lot!!!

## Terminology Headaches #6-9

| CUDA/Nvidia | OpenCL/AMD | Henn&Patt                        |
|-------------|------------|----------------------------------|
| Thread      | Work-item  | Sequence of SIMD Lane Operations |
| Warp        | Wavefront  | Thread of SIMD Instructions      |
| Block       | Workgroup  | Body of vectorized loop          |
| Grid        | NDRange    | Vectorized loop                  |



# NVIDIA HPC SDK



# GPUs in Python

- NumPy and SciPy do not utilize GPUs out of the box

# GPUs in Python

- NumPy and SciPy do not utilize GPUs out of the box
- There are many Python GPU frameworks out there:
  - “drop in” replacements for numpy, scipy, pandas, scikit-learn, etc
    - **CuPy, RAPIDS**

# GPUs in Python

- NumPy and SciPy do not utilize GPUs out of the box
- There are many Python GPU frameworks out there:
  - “drop in” replacements for numpy, scipy, pandas, scikit-learn, etc
    - **CuPy, RAPIDS**
  - “machine learning” libraries that also support general GPU computing
    - **PyTorch, TensorFlow, JAX**

# GPUs in Python

- NumPy and SciPy do not utilize GPUs out of the box
- There are many Python GPU frameworks out there:
  - “drop in” replacements for numpy, scipy, pandas, scikit-learn, etc
    - **CuPy, RAPIDS**
  - “machine learning” libraries that also support general GPU computing
    - **PyTorch, TensorFlow, JAX**
  - “I want to write my own GPU kernels”
    - **Numba, PyOpenCL, PyCUDA, CUDA Python**
  - multi-node / distributed memory:
    - **mpi4py+X, dask, cuNumeric**

# GPUs in Python

- NumPy and SciPy do not utilize GPUs out of the box
- There are many Python GPU frameworks out there:
  - “drop in” replacements for numpy, scipy, pandas, scikit-learn, etc
    - **CuPy, RAPIDS**
  - “machine learning” libraries that also support general GPU computing
    - **PyTorch, TensorFlow, JAX**
  - “I want to write my own GPU kernels”
    - **Numba, PyOpenCL, PyCUDA, CUDA Python**
  - multi-node / distributed memory:
    - **mpi4py+X, dask, cuNumeric**



# Next time

- We will learn about GPU programming with Python (Numba, Cupy)
- Learn basic kernel programming and operations
- Things to be aware of (how to locate the threadID, determine the correct number of threads to launch etc ...)