

Deep Learning

A course about theory & practice



Hardware for Deep Learning

Marco Piastra

GPU vs. CPU

- **The GPU resides on a separate board**

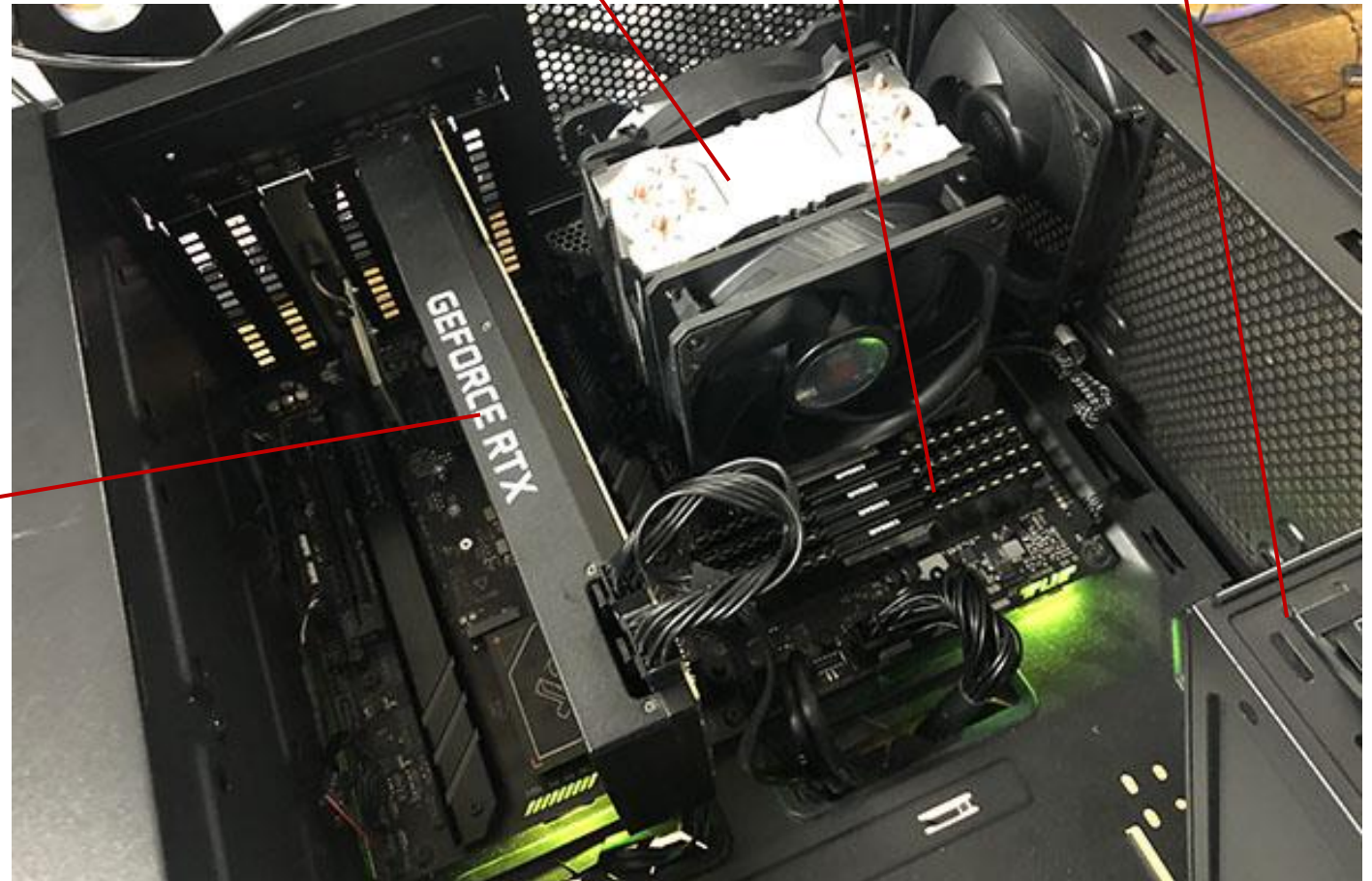
An almost independent computer

CPU, with ventilation

RAM chips

Power Supplier

GPU Board, with its own DRAM



[image from <https://www.researchgate.net/publication/322525660>]

GPU vs. CPU

▪ Different hardware architectures

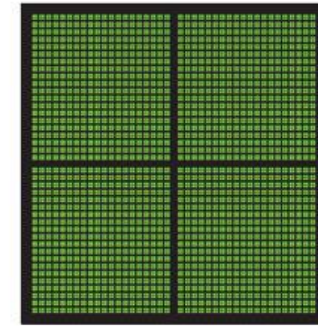
Different computing paradigms

A trade-off between

- fully independent cores (CPU)
- interdependent cores (GPU)
with some (limited) degrees of independence



CPU
MULTIPLE CORES



GPU
THOUSANDS OF CORES



[images from <http://www.nvidia.com/docs/>]

GPU vs. CPU

- **Different hardware architectures**

Different computing paradigms

	Cores	Clock Speed	Memory	Price	Speed
CPU (Intel Core i7-7700k)	10	4.3 GHz	System RAM	\$385	~640 GFLOPs FP32
GPU (NVIDIA RTX 3090)	10496	1.6 GHz	24 GB GDDR6X	\$1499	~35.6 TFLOPs FP32
GPU (Data Center) NVIDIA A100	6912 CUDA, 432 Tensor	1.5 GHz	40/80 GB HBM2	\$3/hr (GCP)	~9.7 TFLOPs FP64 ~20 TFLOPs FP32 ~312 TFLOPs FP16
TPU Google Cloud TPUv3	2 Matrix Units (MXUs) per core, 4 cores	?	128 GB HBM	\$8/hr (GCP)	~420 TFLOPs (non-standard FP)

[image http://cs231n.stanford.edu/slides/2021/lecture_6.pdf]

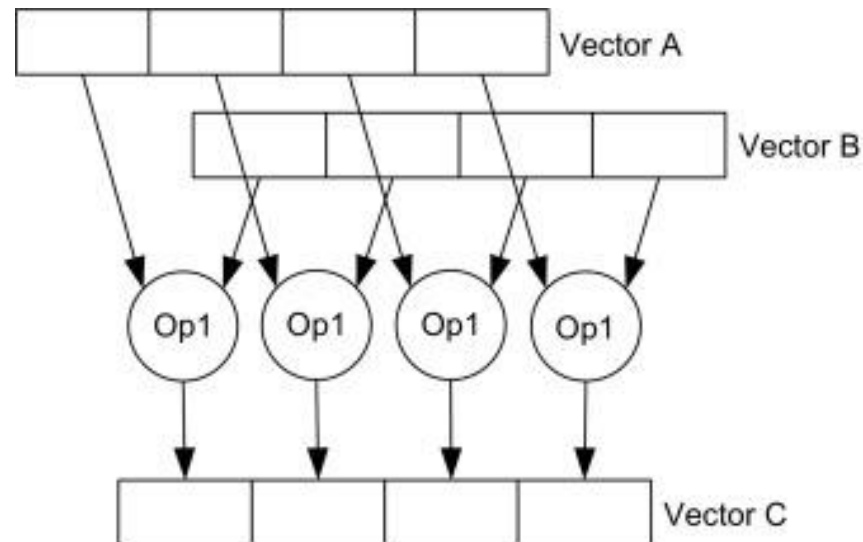
SIMT Parallelism

- **Single Instruction, Multiple Data (SIMD)**

Execution is parallel

All cores are executing the same instruction, in sync

Each core works on specific data



[images from <https://www.sciencedirect.com/topics/computer-science/single-instruction-multiple-data>]

SIMT Parallelism

▪ *Single Instruction, Multiple Threads* (SIMT)

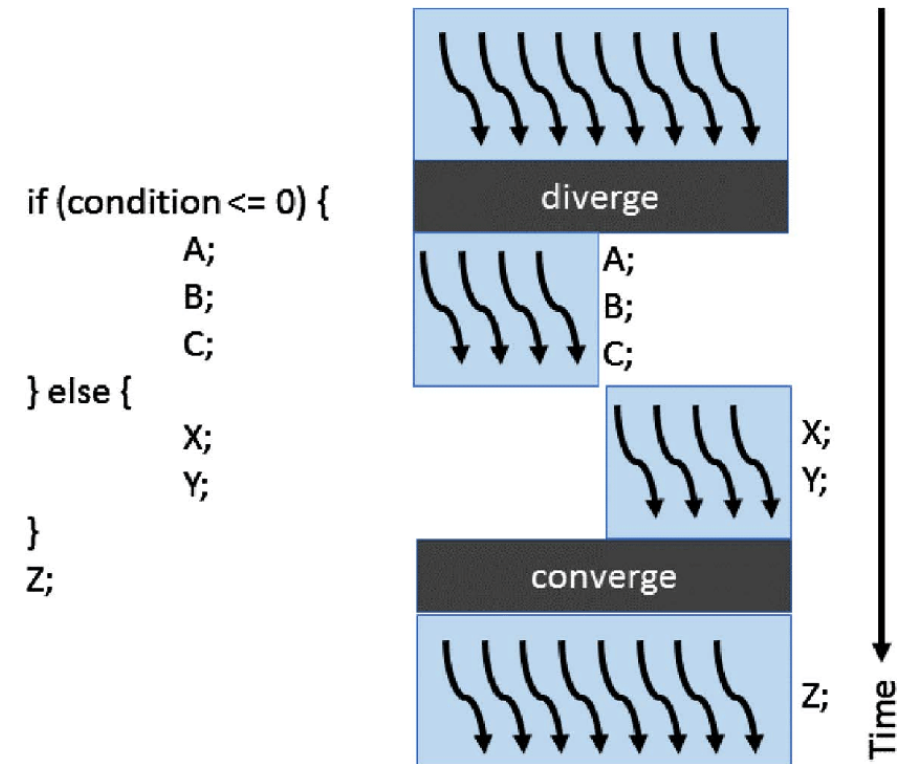
Execution is parallel

All ***active*** cores are executing the same instruction, in sync

Each core works on specific data

The control system activates and deactivates cores on each execution branch

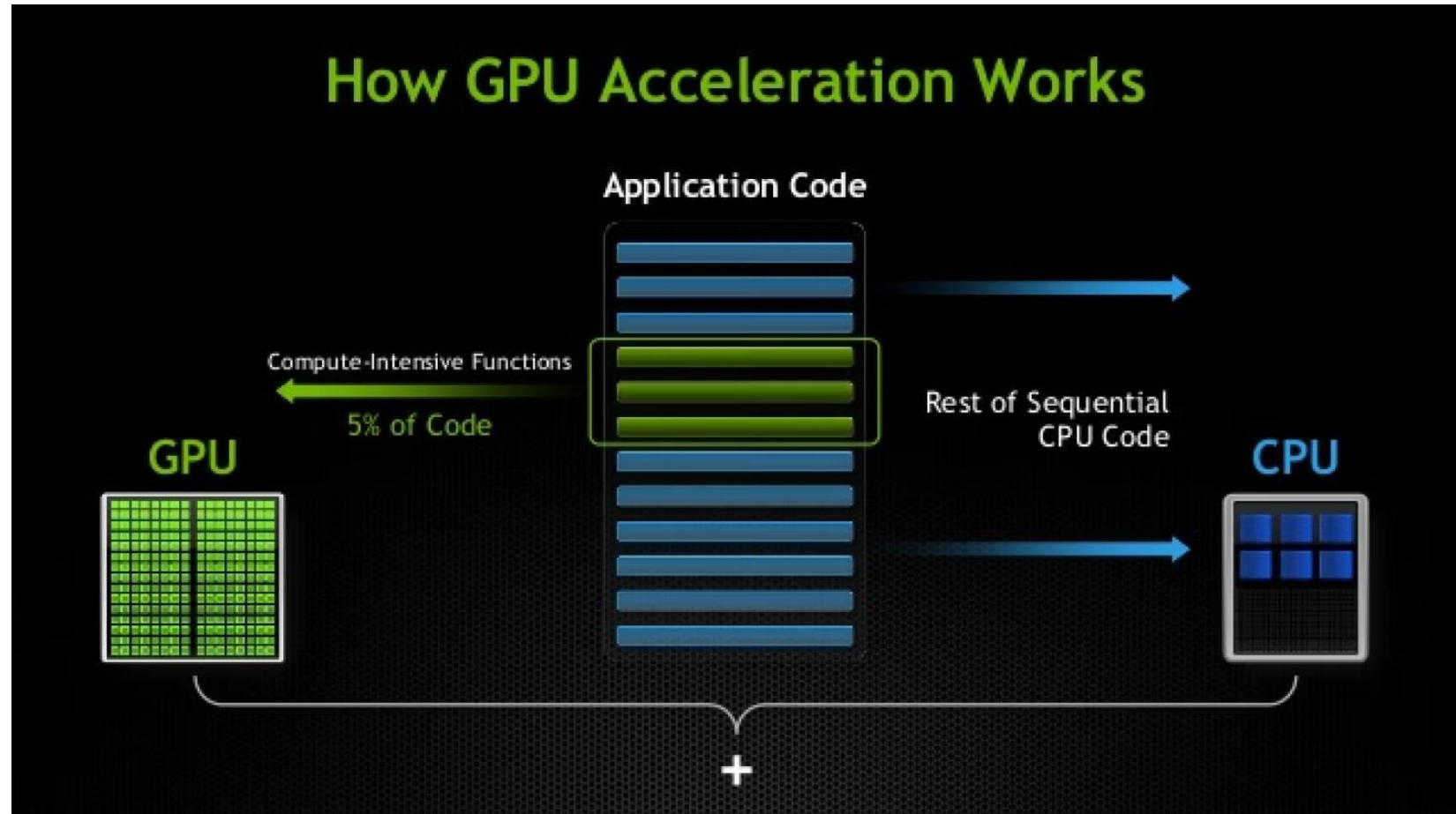
Moral: any computation might be performed, but divergent threads will be sequentialized



[images from <https://www.sciencedirect.com/topics/computer-science/single-instruction-multiple-data>]

Selective parallelization

Not all parts of a program are worth executing in parallel...



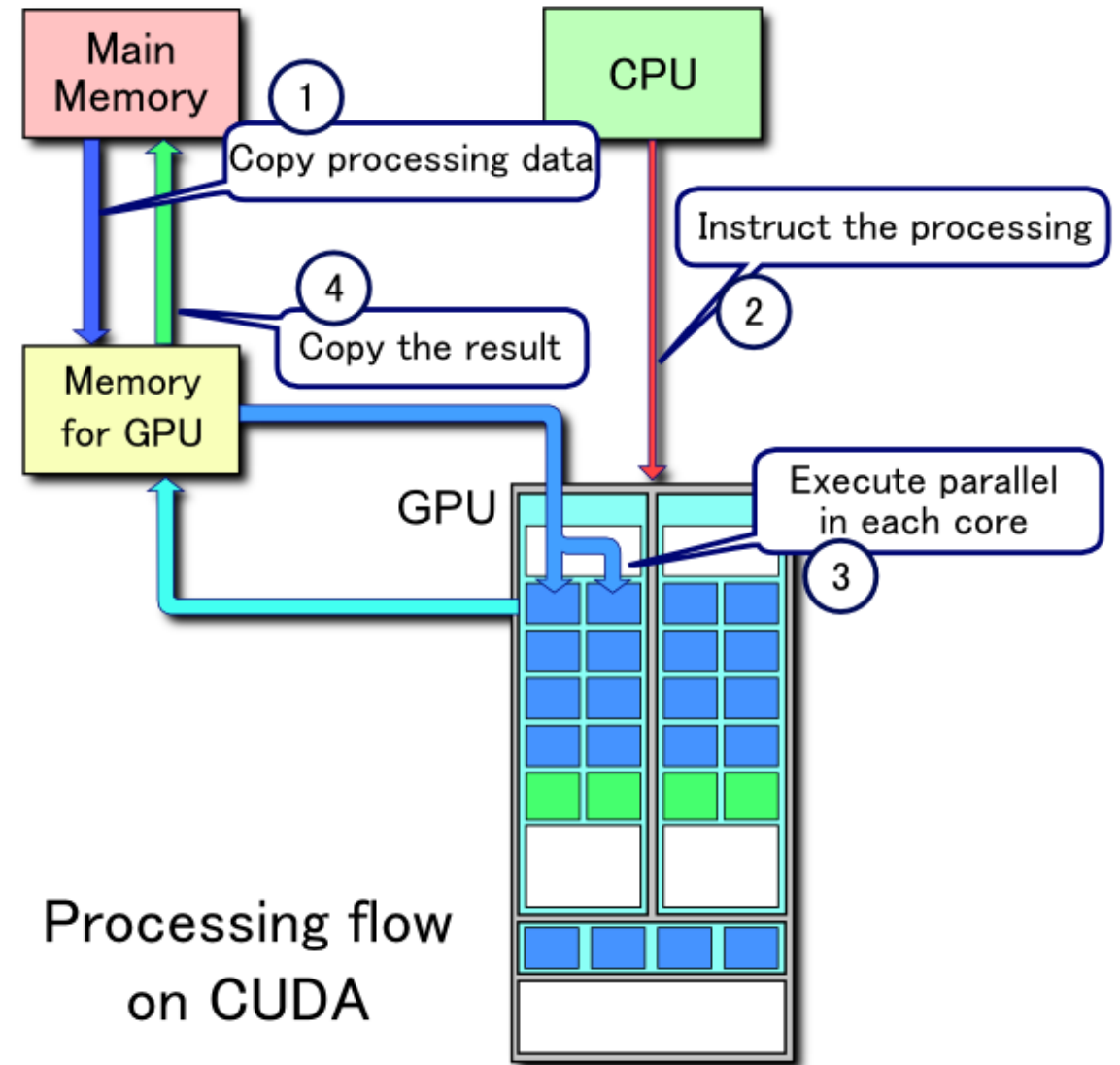
[images from <http://www.nvidia.com/docs/>]

GPU Processing Cycle

■ CPU > Memory Transfer > GPU and back

The program on the CPU drives the execution:

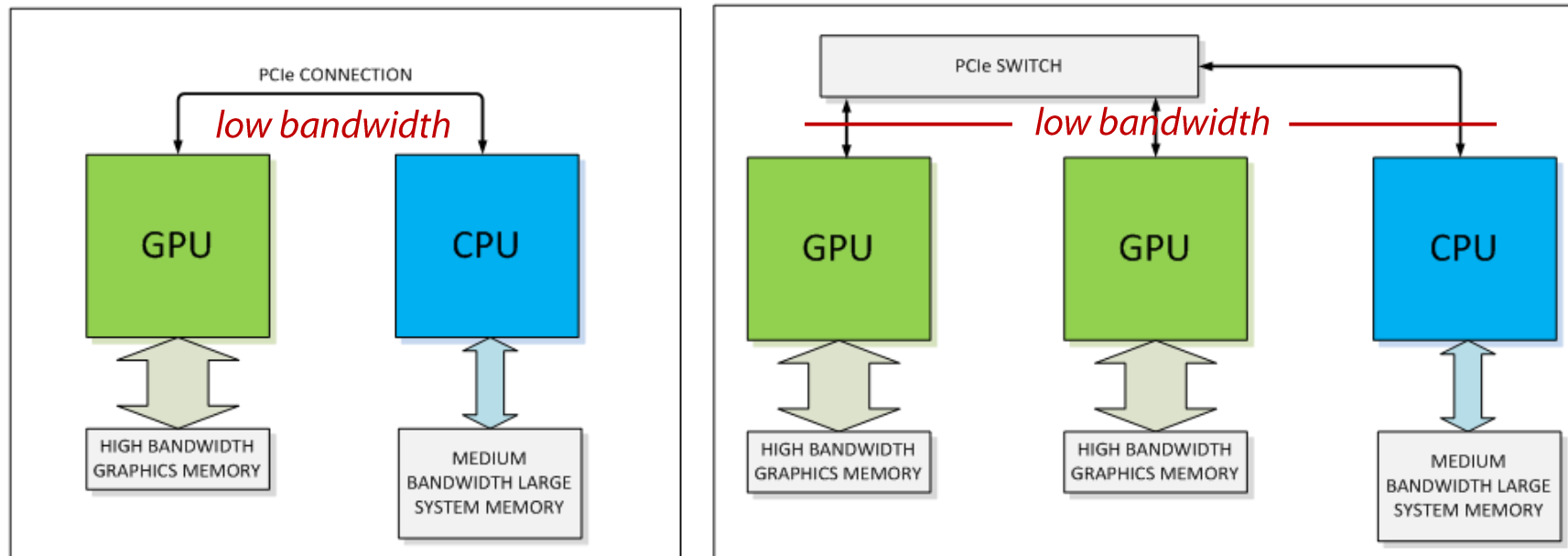
1. All data (program + actual data) are transferred from main memory to GPU DRAM
2. The GPU kernel is launched
3. The Kernel is executed in parallel onto GPU cores, using GPU DRAM
4. Results are copied back from GPU DRAM to main memory



[image from [https://commons.wikimedia.org/wiki/File:CUDA_processing_flow_\(En\).PNG](https://commons.wikimedia.org/wiki/File:CUDA_processing_flow_(En).PNG)]

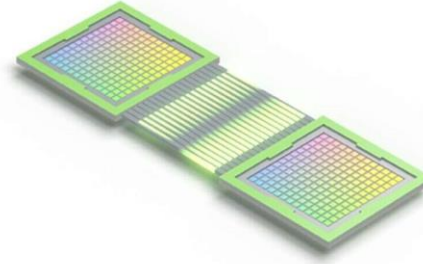
PyTorch and GPUs

- PyTorch computations are optimized to be run on **GPUs**
For the programmer, these implementation details are (mostly) **transparent**
TF can also run on the CPU only, but with lower performance.
- PyTorch automatically manages **memory transfers** to/from GPUs
Memory transfers are very costly, due to low bandwidth PCIe

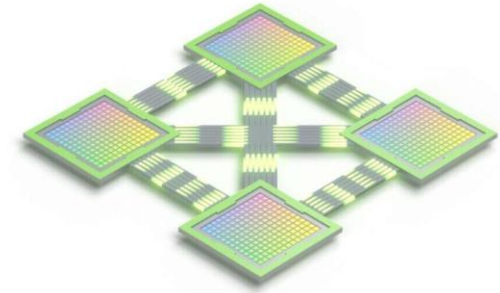


High Speed Interconnect

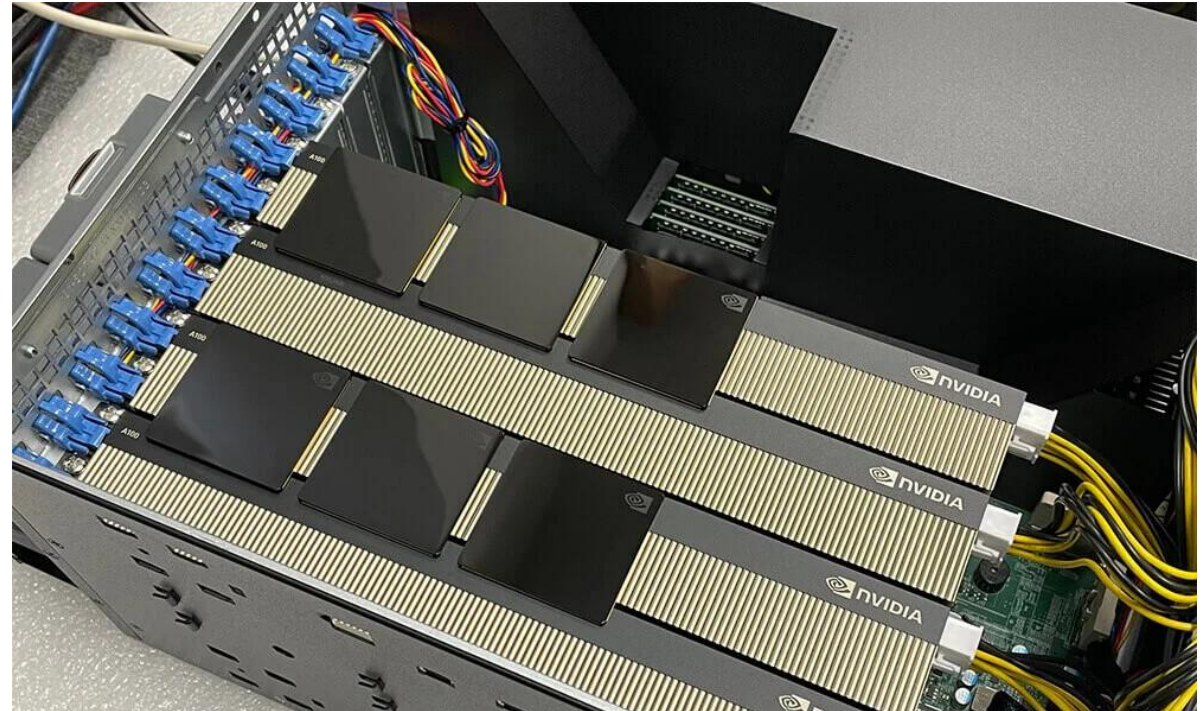
- Available for large GPUs (data center)
 - Dedicate direct link between GPUs:
 - High bandwidth, for faster data communication
 - Low latency
 - Scalability
 - Energy efficiency



NVIDIA H100 PCIe with NVLink GPU-to-GPU connection



NVIDIA H100 with NVLink GPU-to-GPU connections

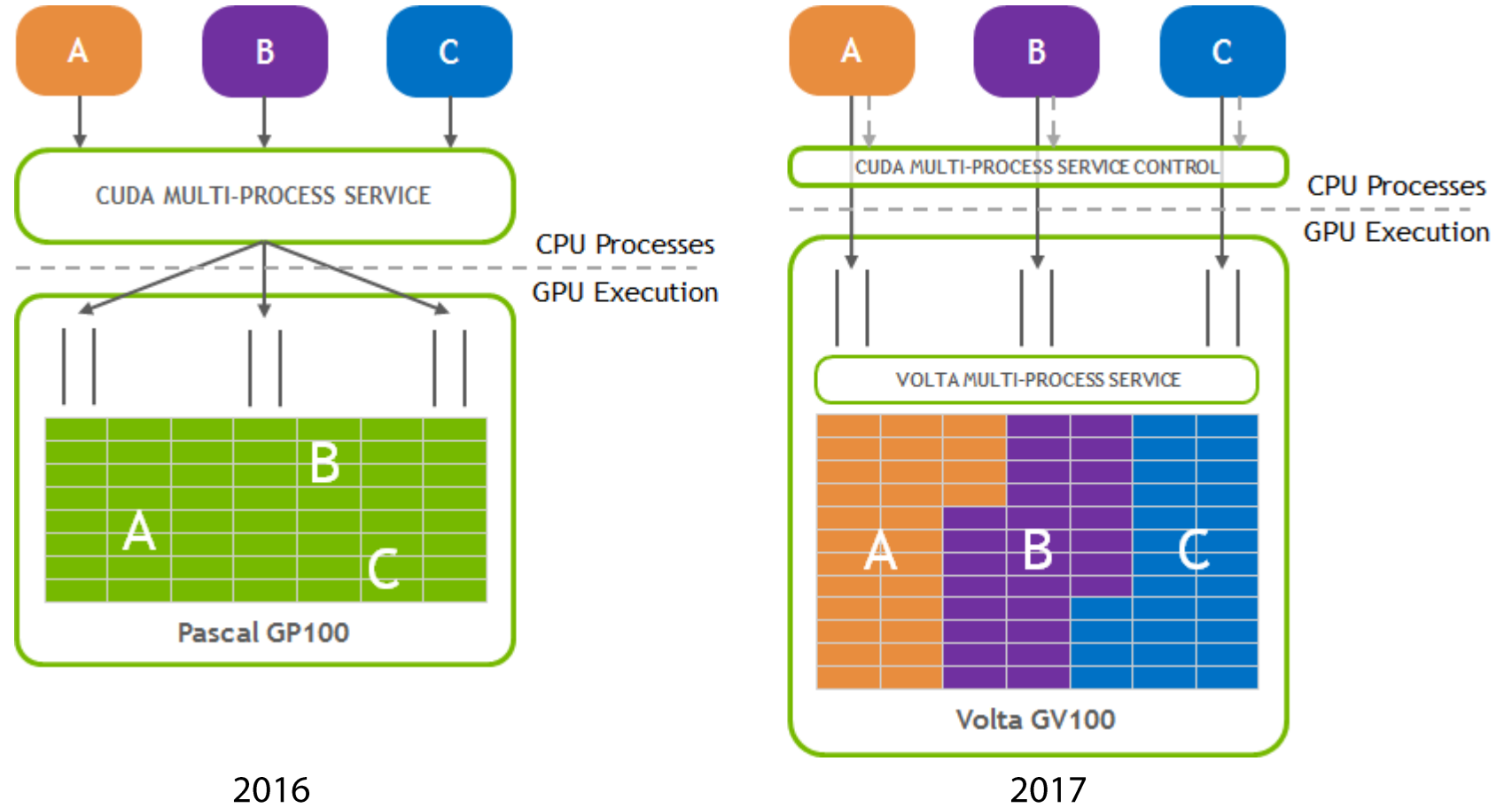


[image from <https://www.cudocompute.com/blog/a-beginners-guide-to-nvidia-gpus>]

GPU Multiprocessing

Until recently, GPUs could only serve one process at time

Now they can be partitioned among several processes

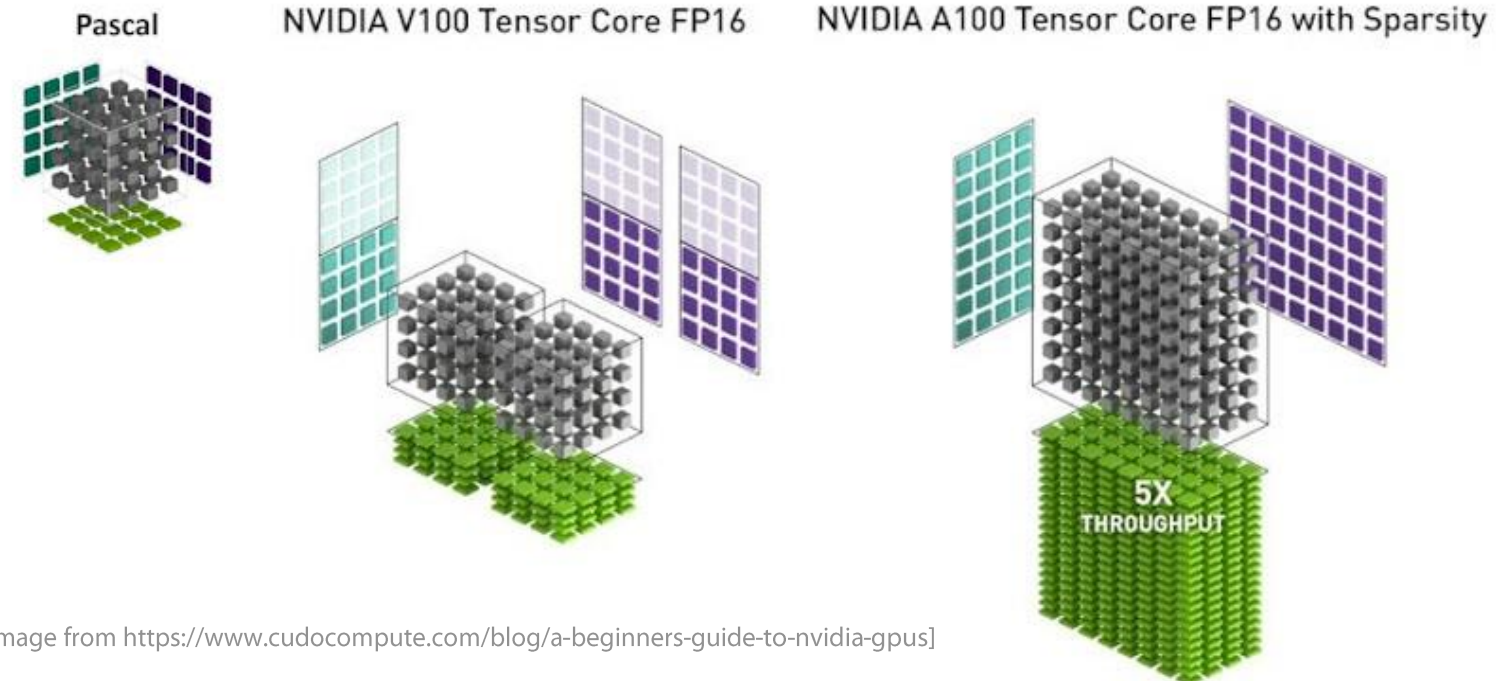


[images from <https://docs.nvidia.com/deploy/mps/>]

Tensor Cores

Specialized processing units to accelerate tensor algebra operations

- **Matrix Multiply-Accumulate (MMA) units**
Each MMA unit can perform a 4x4 matrix multiply-accumulate operation in a single clock cycle
- **Warp schedulers**
MMA units are kept busy and the data flow is optimized
- **High-speed registers and shared memory**
For storing and sharing intermediate among threads



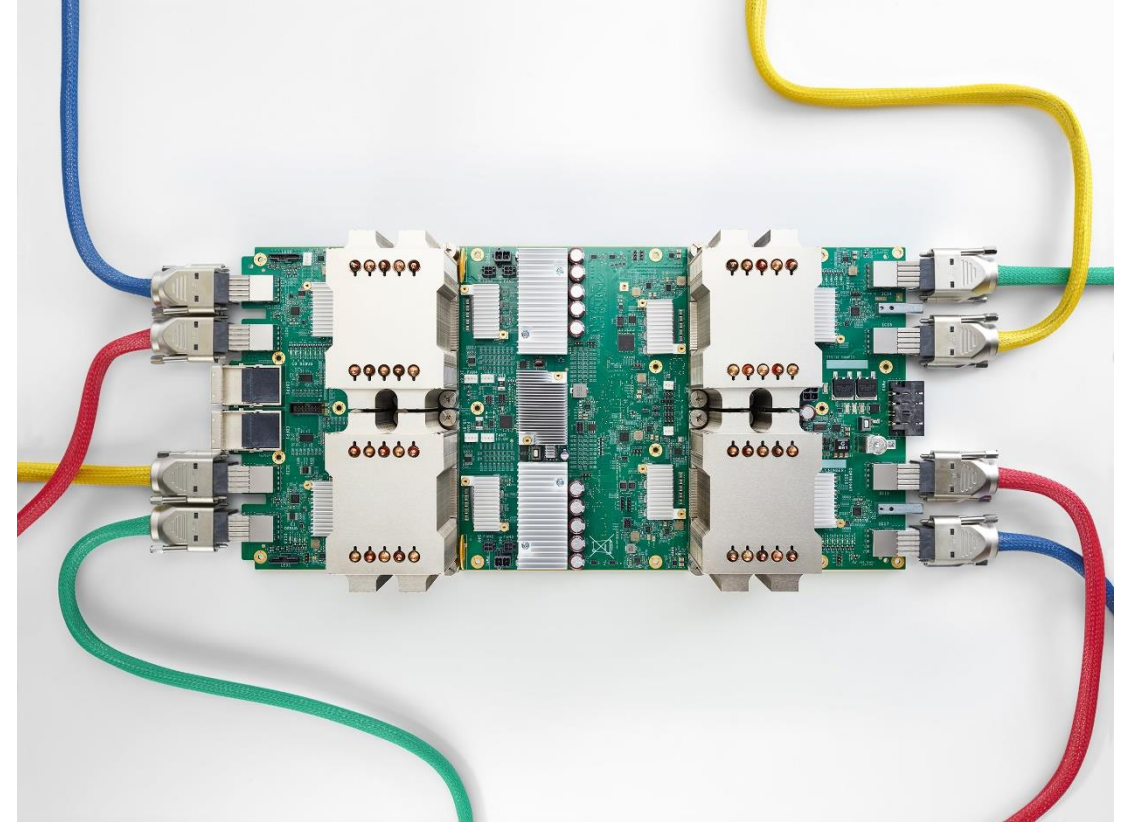
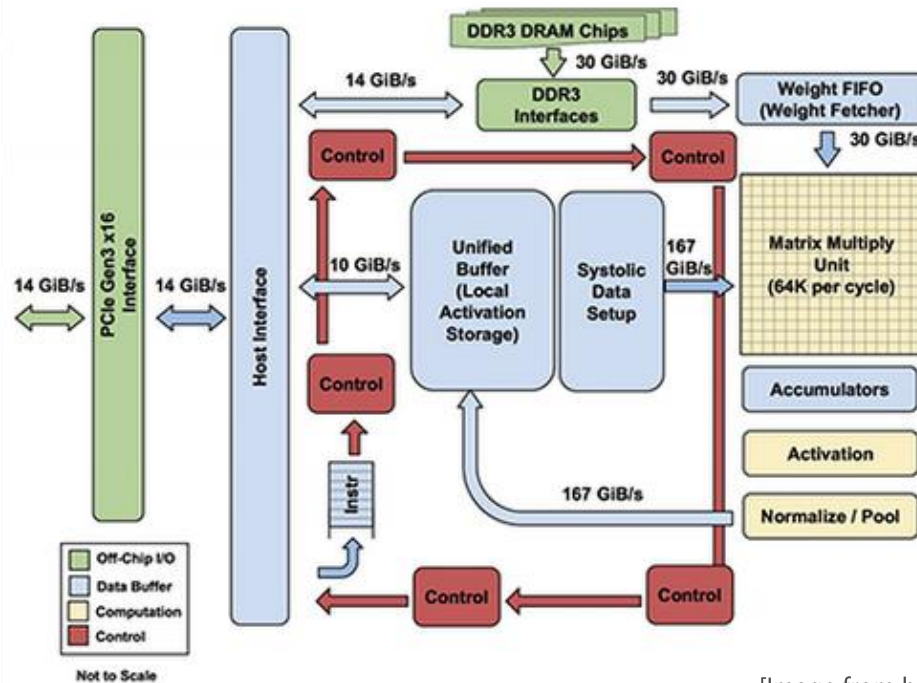
In-Cloud TPUs

- Tensor Processing Units (TPUs)

They are ASICs (Application-Specific Integrated Circuits) and are not on sale

As computation resources, they are only available in cloud (at Google)

TPUs are mounted on separate boards, much like GPUs



[Image from <https://cloud.google.com/blog/products/ai-machine-learning/an-in-depth-look-at-googles-first-tensor-processing-unit-tpu>]

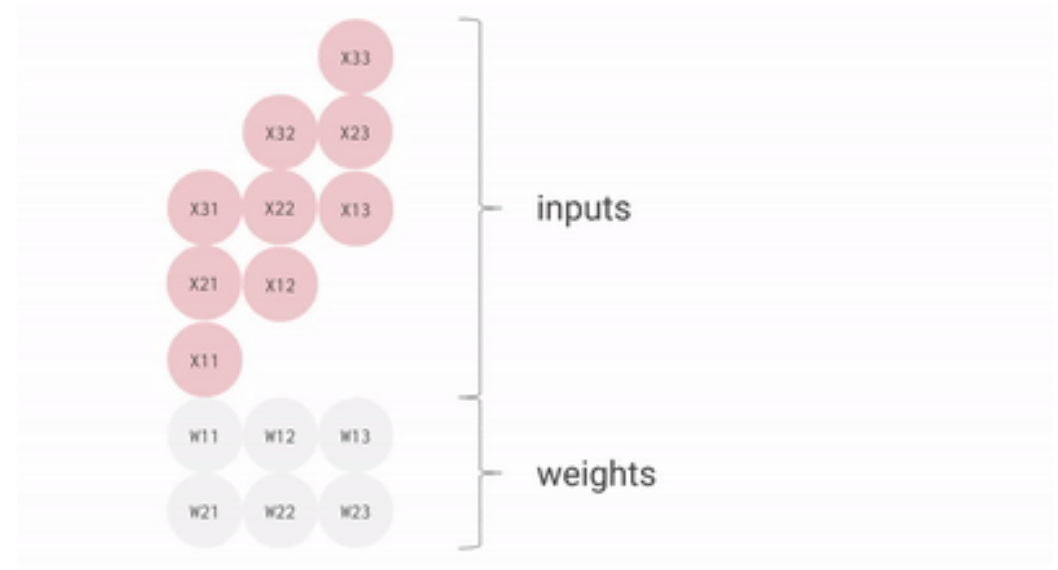
Systolic Parallel Processing

■ Data flow through cores

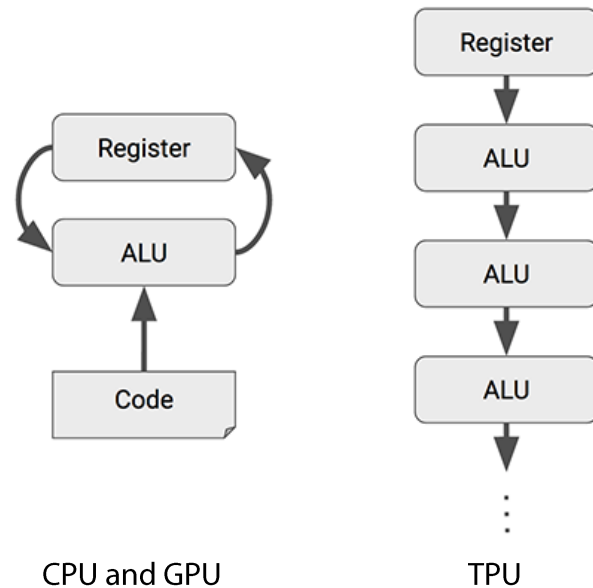
TPU architecture is optimized natively for tensor processing and not for graphics

Arithmetic Logic Units are organized in a pipeline

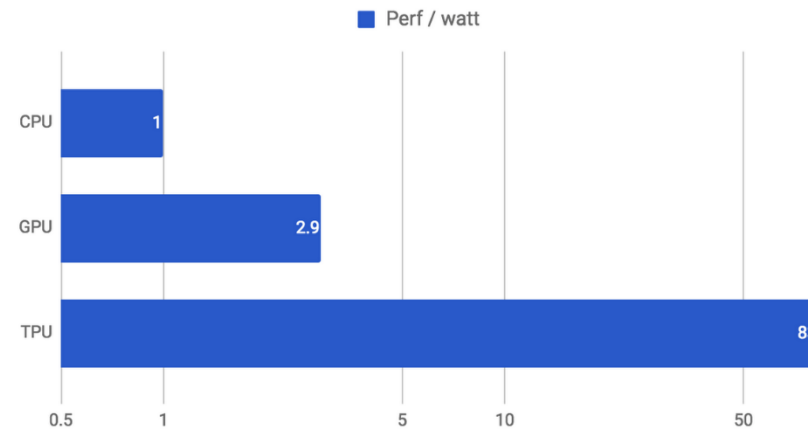
Tensor data are made to 'flow through' the pipeline



Matrix multiplication



TPUs can be much more efficient for tensor computations



[Image from <https://cloud.google.com/blog/products/ai-machine-learning/an-in-depth-look-at-googles-first-tensor-processing-unit-tpu>]