

AI Accelerators 2

26 July 2025 09:51



What are Artificial Intelligence (AI) Accelerators?



An AI accelerator is **high-performance specialized hardware** that is **optimized for AI workloads** such as **neural networks, machine learning**, and other data-intensive or sensor-driven processes.

The three main types are

1. Central Processing Unit (CPU)
2. Graphics Processing Unit (GPU)
3. Field-Programmable Gate Arrays (FPGA)/Application-Specific Integrated Circuit (ASIC)



Where are they used



We can divide AI Accelerators into two groups (based on where we use them):

- Data centres
- Edge Devices

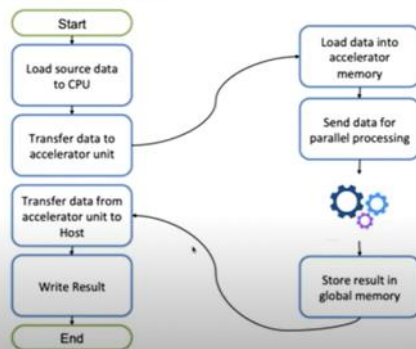
Data Centres



Edge



How do they work



→ AI accelerators should work in co-processing unit.

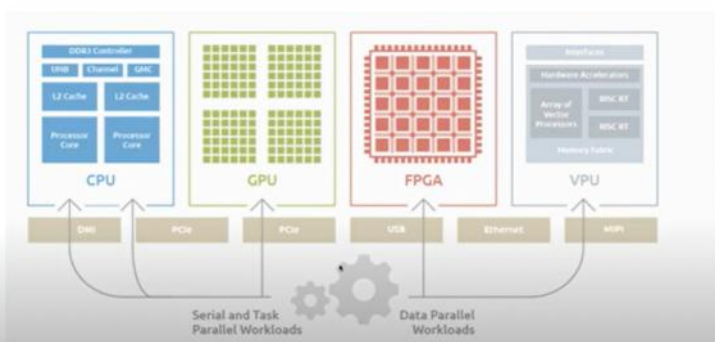
→ # this is how we develop program which could be run on AI accelerators.

→ host - CPU
device - GPU

Not everything is to be offloaded on GPU.



One View of AI Accelerators



Demo 1

- CPU related info
- 1. `lscpu` : displays CPU information
- 2. `cat /proc/cpuinfo`

```
ubuntu@ubuntu:~$ lscpu
Architecture: x86_64
CPU op-mode(s): 32-bit, 64-bit
Byte Order: Little Endian
CPU(s): 8
On-Line CPU(s): 8
Thread(s) per core: 2
Core(s) per socket: 4
Socket(s): 2
NUMA node(s): 2
Vendor ID: GenuineIntel
CPU family: 6
Model: 58
Model name: Intel(R) Core(TM) i7-4770 CPU @ 3.40GHz
Stepping: 3
CPU MHz: 3400.383
CPU max MHz: 3800.0000
CPU min MHz: 800.0000
BogoMIPS: 6553.61
Virtualization: VT-x
L1d cache: 32K
L1i cache: 32K
L2 cache: 256K
L3 cache: 8192K
NUMA node CPU(s): 0-7
```

- GPU related info
- 1. `nvidia-smi` : displays GPU information
- 2. `gpustat` : GPU information

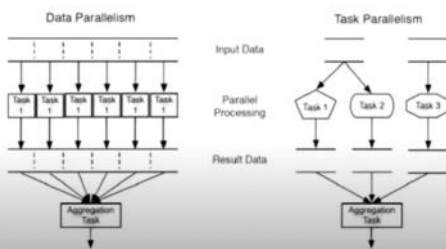
nvidia-smi									
Fri May 19 10:46:20 2022									
GPU	Name	Persistence	API	Bus-Id	Display Mode	GPU-Util	Memory-Usage	GPU-VRAM	Compute M.
0	GeForce 4200M	ON	430.31	000000:00:02.0	Off	0%	10000K / 4096K	0K	Notset
1	GeForce 4200M	ON	430.31	000000:00:02.1	Off	0%	10000K / 4096K	0K	Notset

Process:									
GPU	PID	UID	Type	Process name	GPU Memory Usage				
0	1616	root	G	/usr/lib/plymouth/plymouth	0MB				
0	1616	root	G	/usr/lib/plymouth/plymouth	0MB				
0	1616	root	G	/usr/lib/plymouth/plymouth	0MB				
0	1616	root	G	/usr/lib/plymouth/plymouth	0MB				

only compute intensive tasks which require Parallelism must be offloaded on GPUs

Task and Data Parallelism

Task parallelism vs Data parallelism



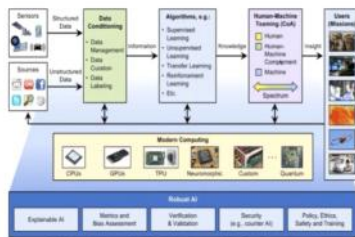
• TASK vs DATA Parallelism

Data Parallelism :- data is split & given to various tasks } data level decomposition

Task Parallelism :- small section of data split that particular data which is unique for each particular task. } data & task level decomposition

(not necessary task 'll use same data)
not tasks 'll do same work.

Second View of AI Accelerators



GPUs

diff. level of complexities available in these accelerators can be used as diff. levels of programming applications

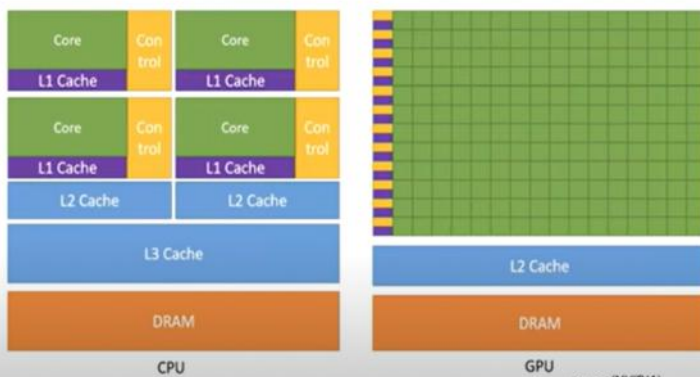




Figure 6. GA100 Full GPU with 128 SMs (A100 Tensor Core GPU has 108 SMs)

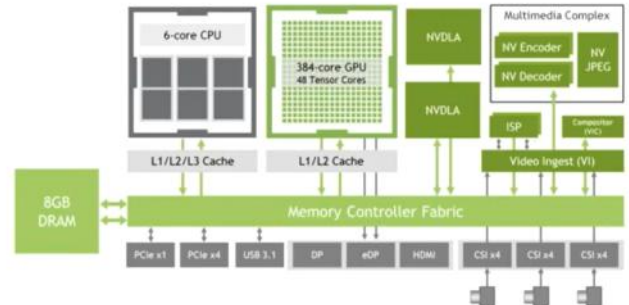


Figure 2. Block diagram of Jetson Xavier NX processor engines including high-speed I/O and memory fabric.

CPU vs GPU Processing

Sequential reduction : (((((((13+27)+15)+14)+33)+2)+24)+6)

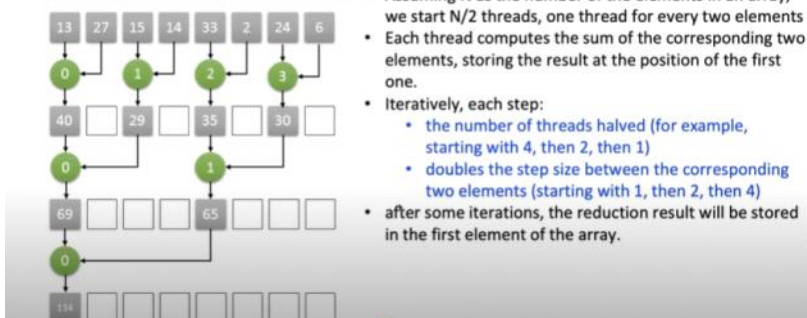


Parallel reduction : ((13+27)+(15+14))+((33+2)+(24+6))



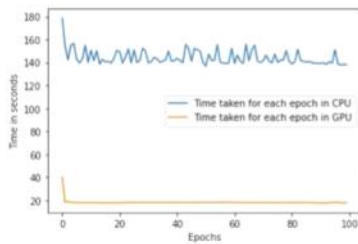
How does parallel reduction happens ?

Parallel Reduction with a GPU



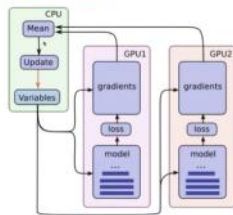


CPU vs. GPU Processing – Training time (Classification Problem)



Epochs	Time taken for Epochs in CPU (s)	Time taken for Epochs in GPU (s)
1	154.26	18.23
2	142.35	18.84
3	154.84	17.92
4	154.78	17.85
5	142.88	17.77
6	139.54	17.69
7	143.02	17.58
8	154.93	17.72
9	139.99	17.76
10	150.66	17.65
11	141.10	17.72
12	150.30	17.65
13	138.58	17.51
14	142.64	17.62
15	140.51	17.52
16	141.09	17.65
17	139.49	17.61
18	143.42	17.61
19	141.08	17.75
20	138.68	17.77
21	140.72	17.88
22	151.33	17.78
23	141.76	17.74
24	141.02	17.80
25	140.29	17.80
26	140.89	17.74
27	139.63	17.73
28	139.36	17.53
29	139.15	17.56
30	139.25	17.41
31	139.62	17.44
32	138.27	17.44
33	140.55	17.72
34	139.36	17.85
35	151.21	17.97
36	138.70	18.07
37	137.66	17.57
38	138.35	17.51
39	138.08	17.61

Train a convolutional neural network on multiple GPU with TensorFlow.

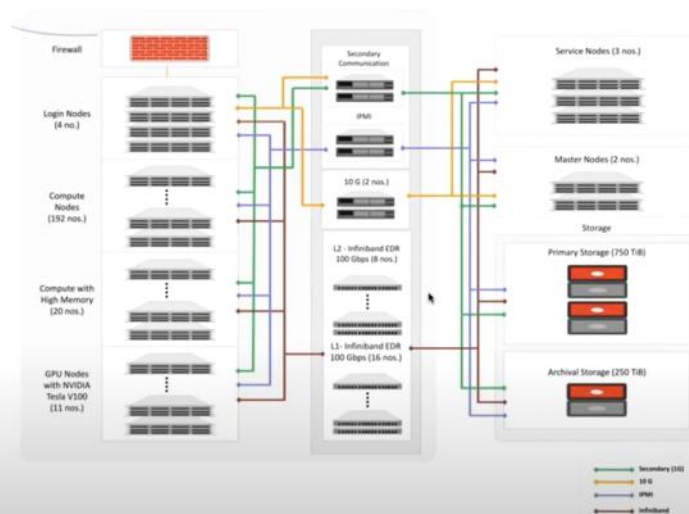


Not all tasks/applications are suitable to work good on gpu they have to be massively parallel computing requirements. CPU does certain specific portion which is not effectively done by gpu, as gpu is compute hungry type of processing units. You cannot depend on this data intensive and i/o intensive applications to be effectively run on gpu.

There can be single gpu/ multi gpu or gpu clusters type of setup



PARAM Shivay – Super Computer of 837 TFLOPS



CPU only Compute Nodes

- + 192 Nodes
- + 7680 Cores
- + Compute power of Peak 589 TFLOPS
- + Each Node with
 - + 2 X Intel Xeon Skylake 6148, 20 cores, 2.4 GHz, processors
 - + 192 GB memory
 - + 480 GB SSD

High Memory Compute Nodes

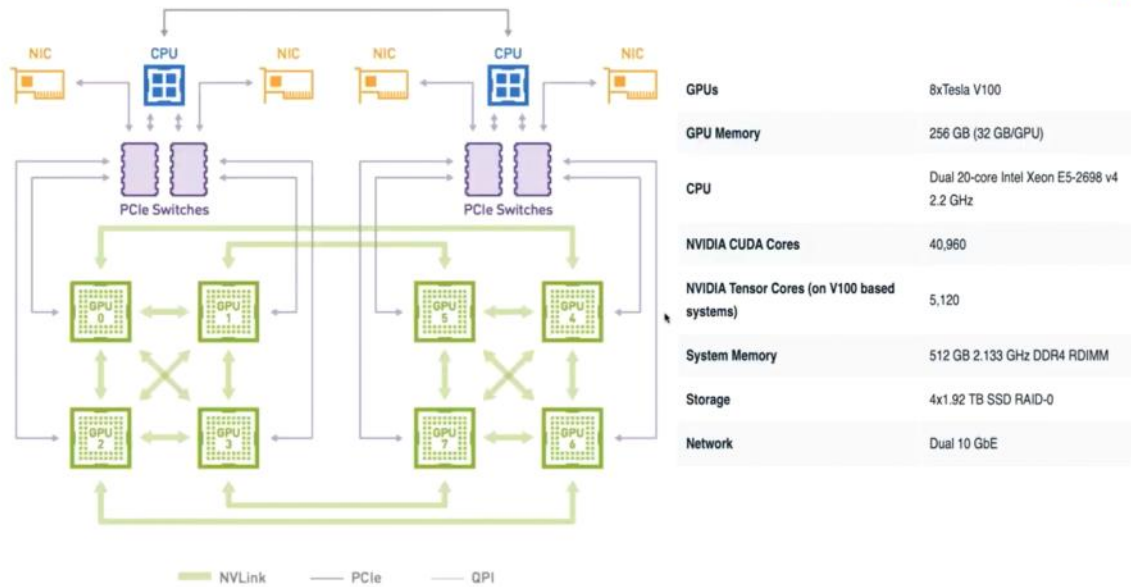
- + 20 Nodes
- + 800 Cores
- + Compute power of Peak 61 TFLOPS
- + Each Node with
 - + 2 X Intel Xeon Skylake 6148, 20 Cores, 2.4 GHz, Processors
 - + 768 GB Memory
 - + 480 GB SSD

GPU Compute Nodes

- + 11 Nodes
- + 440 CPU cores
- + Compute power of Peak 187 TFLOPS
- + Each Node with
 - + 2 X Intel Xeon Skylake 6148, 20 Cores, 2.4 GHz, Processors
 - + 192 GB Memory
 - + 480 GB SSD
 - + 2xNvidia V100 PCIe accelerator cards each with 5120 CUDA cores, 16 GB HBM2

Source: [https://www.paramshivay.in/](#) 26:54 / 46:25

CC BY-NC-SA YouTube



Dgx 1 server

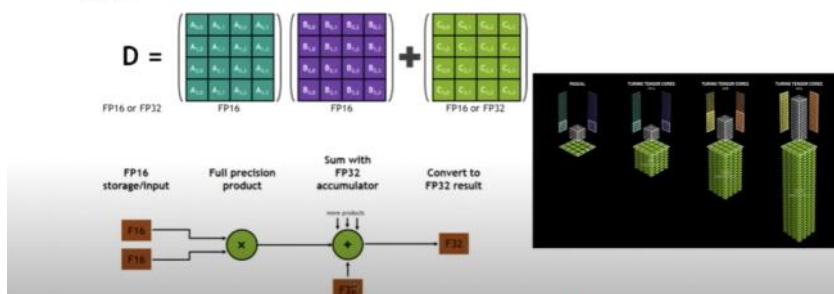
- Has Nvlink propriety bus link specially developed by nvidia



DGX - I - Tensor Cores



- The Tesla V100 GPU contains 640 Tensor Cores: 8 per SM (Streaming Multiprocessor).
- Each Tensor Core provides a 4x4 matrix processing array which performs the operation $D = A * B + C$, where A, B, C and D are 4x4 matrices. The matrix multiplies inputs A and B are FP16 matrices, while the accumulation matrices C and D may be FP16 or FP32 matrices.



Why is tensor cores very important



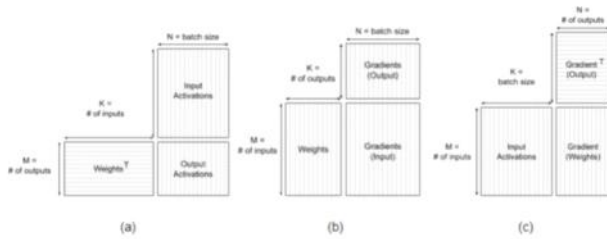
DGX - I - Tensor Cores



- One arithmetic operation that holds high importance is matrix multiplication.
- Multiplying two 4x4 matrices involves 64 multiplications and 48 additions.
- Convolution and Multiplication are the areas where the new cores shine.
- Typically, the notion is that CUDA cores are slower, but offer more significant precision. Whereas a Tensor cores are lightning fast, however lose some precision along the way.
- All Tensor core does is that it accelerates the speed of matrix multiplication.
- Tensor Cores are able to multiply two fp16 matrices 4x4 and add the multiplication product fp32 matrix (size: 4x4) to accumulator (that is also fp32 4x4 matrix).
- **General Matrix Multiplication (GEMM)**
- Instead of needing to use many CUDA cores and more clocks to accomplish the same task it can be done in a single clock cycle causing a dramatic speed up in machine learning applications.



DGX - I - Tensor Cores



Equivalent matrix multiplies for (a) forward propagation, (b) activation gradient calculation, and (c) weight gradient calculation of a fully-connected layer.

- Tensor Cores in CUDA Libraries : cuBLAS uses Tensor Cores to speed up GEMM computations (GEMM is the BLAS term for a matrix-matrix multiplication); cuDNN uses Tensor Cores to speed up both convolutions and recurrent neural networks (RNNs)
- CUDA – WMMA (Warp Matrix Multiply-Accumulate)



DGX - I - Tensor Cores



- DLSS -- Deep Learning Super Sampling
- Basic premise is simple: render a frame at low-ish resolution and when finished, increase the resolution of the end result so that it matches the native screen dimensions of the monitor (e.g. render at 1080p, then resize it to 1400p). That way you get the performance benefit of processing fewer pixels, but still get a nice looking image on the screen.



- Tensor core vs. Gpu cores performance

```
chpc@LAPTOP-TDCSOMDA:~/wmma_tensorcore_sample/project/project$ ./kernel
[*] Initializing Matrix...
[+] A: 4096 x 4096
[+] B: 4096 x 4096
[+] C: 4096 x 4096
[*] Computing D = A * B + C on GPU without Tensor Cores...
[+] GPU(without Tensor Cores) Elapsed Time: 1409.558716 ms
[+] TFLOPS: 0.10
[*] Computing D = A * B + C on GPU with Tensor Cores...
[+] GPU(with Tensor Cores) Elapsed Time: 1386.487549 ms
[+] TFLOPS: 0.10
```