

Which GPU should I use?

Pawel Pomorski

SHARCNET/Compute Ontario/Alliance

Current GPUs in Alliance systems

GPU cards available

- **P100** Pascal (2016)
- **V100** Volta (2017)
- **T4** Turing (2018)
- **A100** Ampere (2020)

Alliance main clusters

- Graham - 320 **P100** , 144 **T4**, 72 **V100**, 8 **A100**
- Cedar - 584 **P100**, 768 **V100**
- Beluga - 688 **V100**
- Narval - 636 **A100**

Other Alliance clusters

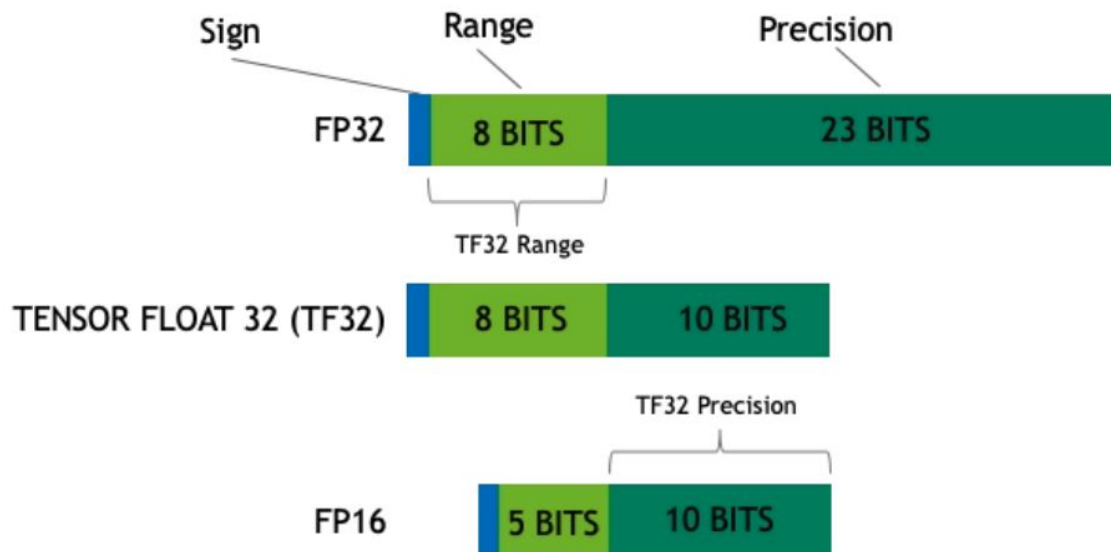
- Mist - 216 **V100** (IBM Power9 CPUs)
- Rouge - 160 **MI-50** (AMD)

NVIDIA GPU performance trends (TFLOPS)

| GPU | FP64 | FP64 TC | FP32 | TF32 TC | FP16/32 TC | INT8 TC |
|------|------|------------|------|---------------|---------------|----------------|
| P100 | 4.7 | - | 9.3 | - | - | - |
| V100 | 7.8 | - | 15.7 | - | 125 | - |
| T4 | 0.25 | - | 8.1 | - | 65 | 130 |
| A100 | 9.7 | 19.5 | 19.5 | 156 312 sp | 312 624 sp | 624 1048 sp |

FP64 - 64 bit floating point, TC - Tensor Core , sp - with sparsity

Reduce precision to gain higher performance



Tensor cores

- Units dedicated to matrix multiplication
- First appeared on Volta doing FP32/FP16 multiply
- Turin and Ampere implement more operations

$$\mathbf{D} = \begin{pmatrix} \begin{matrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{matrix} & \begin{matrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{matrix} & + & \begin{matrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{matrix} \end{pmatrix}$$

FP16 or FP32 FP16 FP16 FP16 or FP32

$$\mathbf{D} = \mathbf{AB} + \mathbf{C}$$

Upcoming GPU hardware in 2023

- NVIDIA
 - Hopper H100 GPU (up to 700 W, available Q3 2022)
 - Grace CPU (ARM)
- AMD:
 - CDNA2 architecture GPUs - MI250, MI250X (November, 2021)
 - MI300 in development, performance 2x of MI250, heterogeneous CPU+GPU
- Intel:
 - Ponte Vecchio GPU expected in 2022, will be in DOE Aurora exascale system
 - Laptop gaming GPU released March 30 - Intel Arc 3, will be followed by Arc 5 and Arc 7.

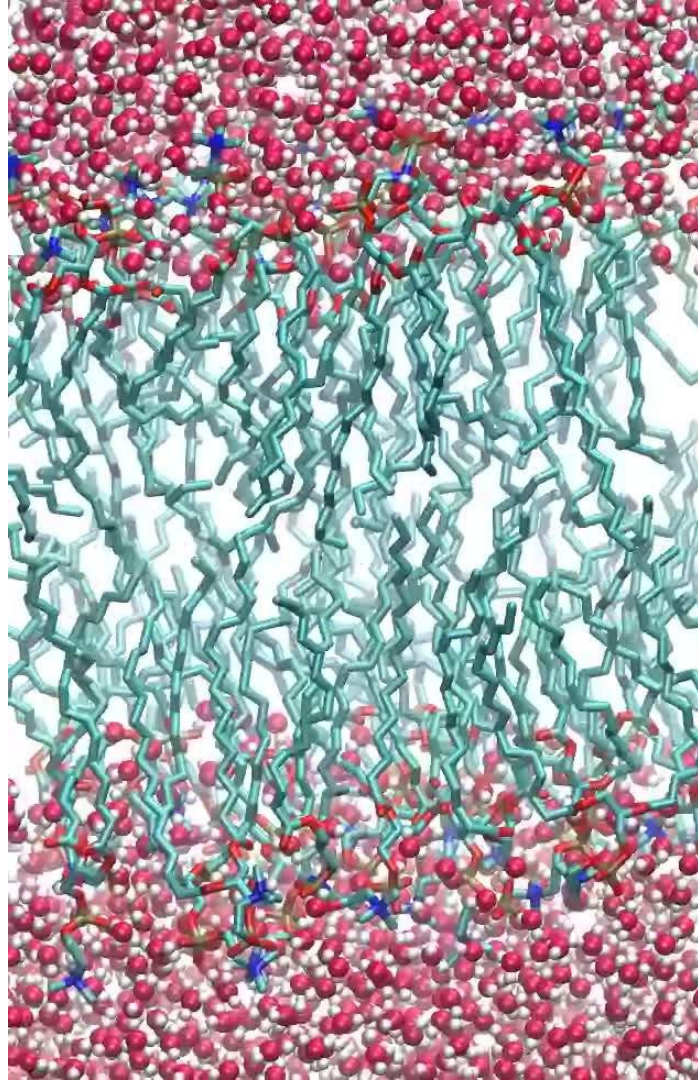
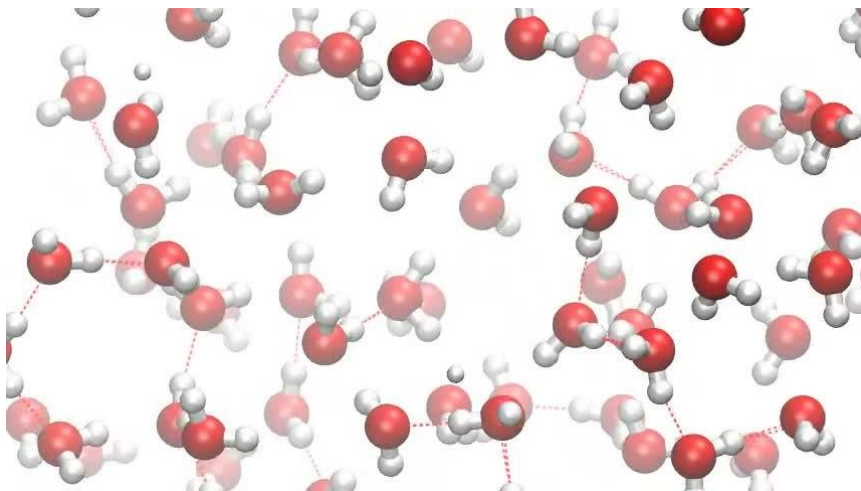
Future systems: NVIDIA versus AMD

| GPU | FP64 | FP64 TC/MC | FP32 | FP32 MC | TF32 TC | FP16 TC/MC | FP8 TC | INT8 TC/MC |
|-------------|------|---------------|------|------------|----------------|-----------------|-----------------|-----------------|
| MI200 | 45.3 | 90.5 | 45.3 | 90.5 | - | 362.1 | - | 362.1 |
| MI250x | 47.9 | 95.7 | 47.9 | 95.7 | - | 383 | - | 383 |
| A100 | 9.7 | 19.5 | 19.5 | - | 156 312 sp | 312 624 sp | - | 624 1048 sp |
| H100 SXM | 30 | 60 | 60 | - | 500 1000 sp | 1000 2000 sp | 2000 4000 sp | 2000 4000 sp |

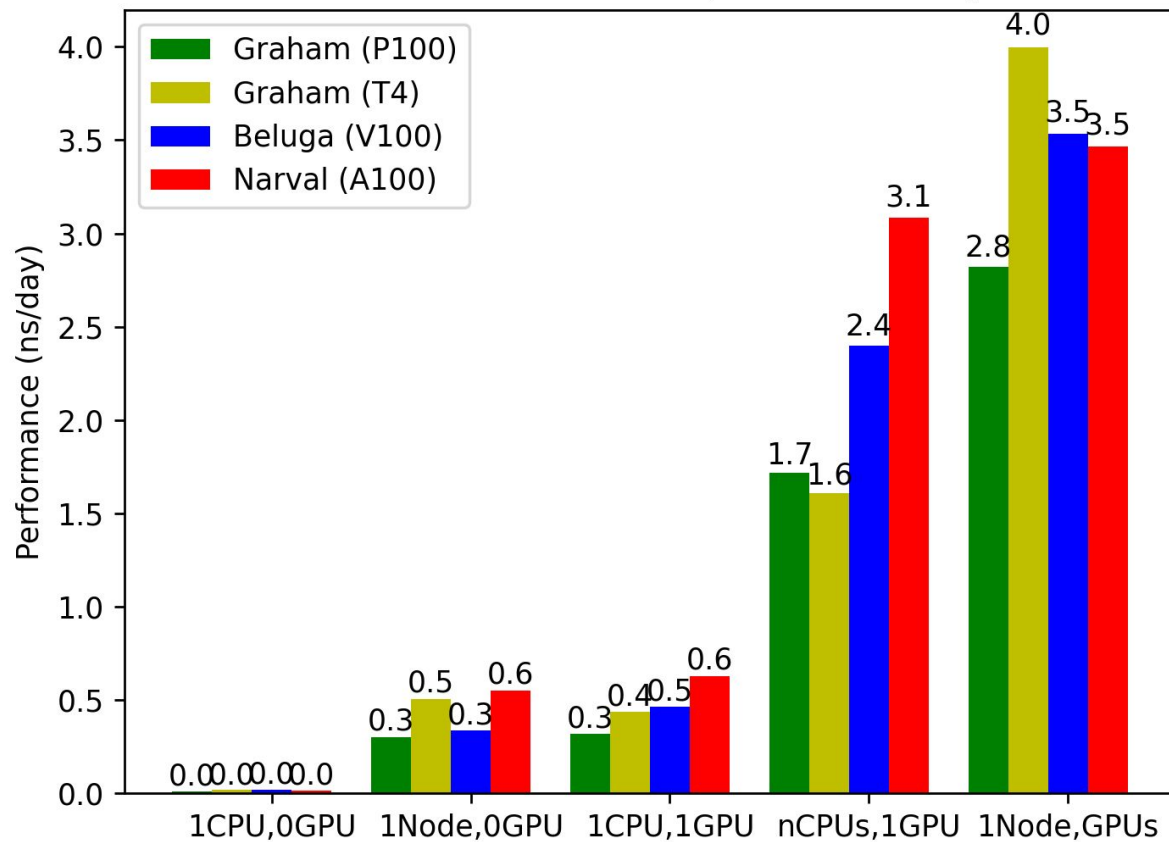
FP64 - 64 bit floating point, TC - Tensor Core , MC - Matrix Core, sp - with sparsity

NAMD

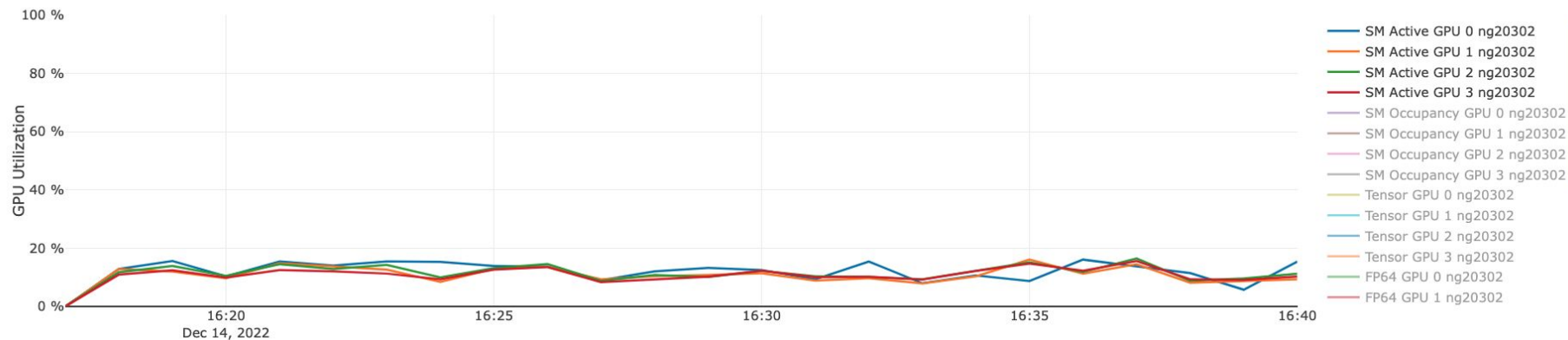
- Standard Molecular Dynamics (MD code)
 - Well-accelerated on GPUs, also uses CPUs
 - MD relies single precision (FP32)
 - Using more parallel resources efficiently requires larger system (more atoms)



NAMD 2.14 STMV Benchmark performance by Cluster



Why narval does not gain performance with 4 GPUs?



Plot from <https://portail.narval.calculquebec.ca>

NAMD 3.0alpha

In this version work moved to the GPU, so can run efficiently with just 1 CPU core per GPU. Enabled with:

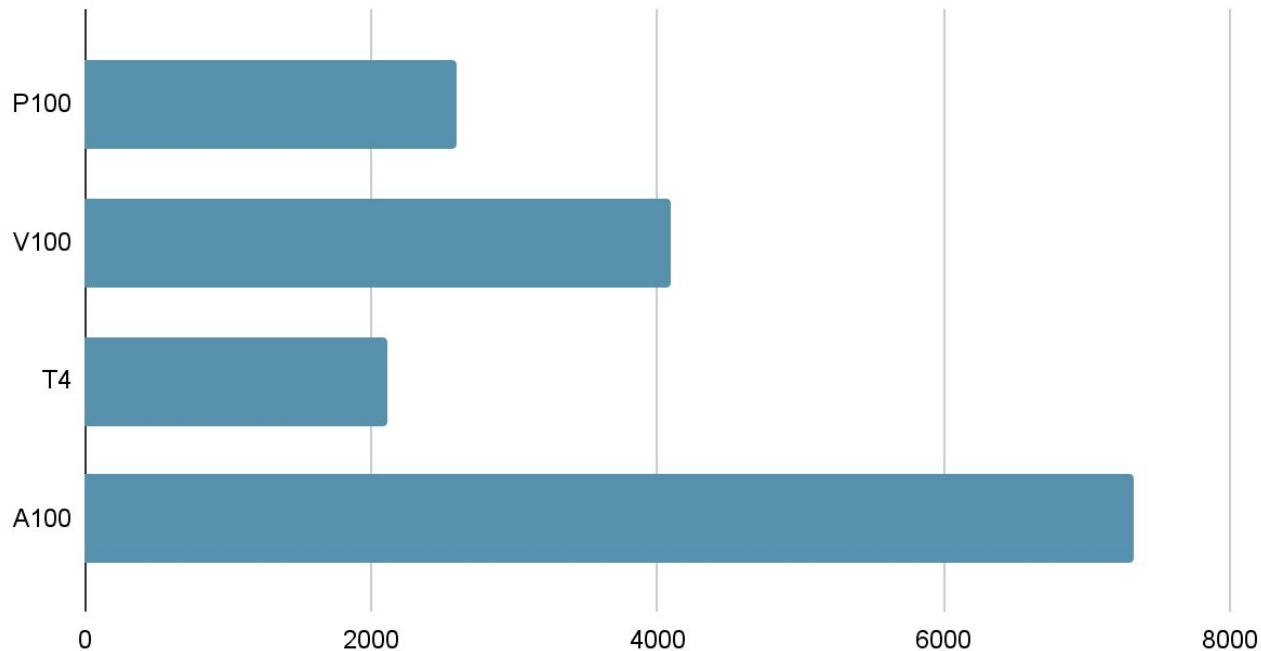
CUDASOAintegrate on

1 GPU with 1 core - 7.6 ns/day (compared to 3.1 ns/day for NAMD 2.14)

4 GPU with 4 core - 26 ns/day (compared to 3.5 ns/day for NAMD 2.14)

Tensorflow performance

keras_cifar_benchmark.Resnet56KerasBenchmarkSynth



MIG - subdividing Nvidia GPUs

GPUs have O(64) actual cores (called SM, EU, etc)

MIG allows isolation of SM+MC combinations: up to 7 pieces

MIG configuration can be changed whenever the whole GPU is idle

Some talk about subdivision from AMD or Intel

MIG performance

