

NAIRR Pilot Compute Calculations

We made the following calculations to estimate the compute available through the NAIRR pilot, which encompasses six Federal HPC resources. This includes calculations to (a.) estimate the FLOPS of CPUs, (b.) convert GPU FLOPS to FP16, and (c.) aggregate GPU, CPU, and ASIC FLOPS across nodes. Information on the hardware for each NAIRR pilot resource was sourced primarily from the NAIRR website,¹ information used to convert all FLOPS to FP16 was sourced primarily from TechPowerUp,² and information on CPU specifications was sourced from the developers' websites (i.e., Intel and AMD's websites). Note that FP16 FLOPS for some resources were directly provided by sources, and in these cases we did not perform calculations separately.

Delta Resource³

GPUs:

1. NVIDIA A100 40GB HBM2e⁴ [tensor FP16 without sparsity] = 312TF
 - a. [400 NVIDIA A100 40GB HBM2] $400 * 312\text{TF} = 124800\text{TF}$
2. NVIDIA A40 48GB GDDR6 PCIe 4.0⁵ = 37.42TF
 - a. [400 NVIDIA A40 48GB GDDR6 PCIe 4.0] $400 * 37.42 = 14968\text{TF}$
3. NVIDIA A100 GPUs⁶ [tensor FP16 without sparsity = 312TF]
 - a. [40 NVIDIA A100 GPUs (assuming 40GB HBM2)] $40 * 312\text{TF} = 12480\text{TF}$
4. AMD MI100⁷ = 184.6TF
 - a. [8 AMD MI100] $184.6 * 8 = 1476.8\text{TF}$
5. Total GPU FLOPS: $124800\text{TF} + 14968\text{TF} + 12480\text{TF} + 1476.8\text{TF} = 153724.8\text{TF} = 153.7248\text{PF}$

¹ <https://nairrpilot.org/allocations>

² <https://www.techpowerup.com/gpu-specs/>

³ <https://www.ncsa.illinois.edu/research/project-highlights/delta/>

⁴

<https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a100-datasheet-nvidia-us-2188504-web.pdf>

⁵

[https://www.techpowerup.com/gpu-specs/a40-pcie.c3700#:~:text=NVIDIA%20has%20paired%2048%20GB,MHz%20\(14.5%20Gbps%20effective\)](https://www.techpowerup.com/gpu-specs/a40-pcie.c3700#:~:text=NVIDIA%20has%20paired%2048%20GB,MHz%20(14.5%20Gbps%20effective))

⁶

<https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a100-datasheet-nvidia-us-2188504-web.pdf>

⁷ <https://www.techpowerup.com/gpu-specs/radeon-instinct-mi100.c3496>

CPUs: 200 AMD EPYC 7763 ("Milan") CPUs⁸ with 64-cores/socket (64-cores/node) at 2.55GHz and 256GB of DDR4-3200 RAM

6. $\text{AVX 256 [bits]} / 16 [\text{FP16}] * 2 [\text{multiple and add; 2 operations}] * 2 [\text{two AVX 256 blocks per core}] * 64 [64 \text{ cores per CPU}] * 2.55 \text{ Ghz} = 10444.8 \text{ GF/node}$
 - a. $10444.8 \text{ GF/node} * 200 [\text{total CPUs}] = 2088960 \text{ GF} = 2.08896 \text{ PF}$

Total FLOPS:

7. $153.7248 \text{ PF} [\text{from GPUs}] + 2.08896 \text{ PF} [\text{from CPUs}] = 155.81376 \text{ PF}$

Frontera Resource⁹

GPUs:

1. The 90 nodes with NVIDIA Quadro RTX 5000 GPUs can do 4PF in FP32 in total.¹⁰
Based on conversions using the FP32/FP16 numbers here,¹¹ this would be 8.028PF.
 - a. Note that we used non-tensor FP16 FLOPS because the Quadro RTX 5000 specs sheet does not indicate the precision of the tensor performance.¹²

CPUs:

2. Intel® Xeon® Platinum 8280 Processors¹³
 - a. $\text{AVX 512 [bits]} / 16 [\text{FP16}] * 2 [\text{multiple and add; 2 operations}] * 2 [2 \text{ AVX 512 blocks per core}] * 28 [28 \text{ cores per CPU}] * 2.7 \text{ Ghz} = 9676.8 \text{ GF}$
 - i. $9676.8 \text{ GF} * 2 [2 \text{ CPUs per node}]^{14} = 19353.6 \text{ GF/node}$
 - ii. $19353.6 \text{ GF/node} * 8368 [\text{total nodes}] = 161950924.8 \text{ GF} = 161.95 \text{ PF}$
3. Intel Xeon CPU E5-2620 v4 Processors¹⁵
 - a. $256 [\text{bits}] / 16 [\text{FP16}] * 2 [\text{multiple and add; 2 operations}] * 2 [\text{blocks per core}] * 8 [8 \text{ cores per CPU}] * 2.1 \text{ GHz} = 1075.2 \text{ GF/CPU}$
 - b. $1075.2 \text{ GF/CPU} * 360 [360 \text{ CPUs}] = 387072 \text{ GF} = 387.072 \text{ TF} = 0.387 \text{ PF}$

⁸ <https://www.amd.com/en/products/cpu/amd-epyc-7763>

⁹ <https://tacc.utexas.edu/systems/frontera/>

¹⁰ <https://tacc.utexas.edu/systems/frontera/>

¹¹ <https://www.techpowerup.com/gpu-specs/quadro-rtx-5000.c3308>

¹²

<https://www.nvidia.com/content/dam/en-zz/Solutions/design-visualization/quadro-product-literature/quadro-rtx-5000-data-sheet-us-nvidia-704120-r4-web.pdf>

¹³

<https://ark.intel.com/content/www/us/en/ark/products/192478/intel-xeon-platinum-8280-processor-38-5m-cache-2-70-ghz.html>

¹⁴ <https://tacc.utexas.edu/systems/frontera/>

¹⁵

<https://www.intel.com/content/www/us/en/products/sku/92986/intel-xeon-processor-e52620-v4-20m-cache-2-10-ghz/specifications.html>

4. CPUs total
 - a. $161.95\text{PF} + 0.387\text{PF} = 162.337\text{PF}$

Total FLOPS:

5. 8.028PF [from GPUs] + 162.337PF [from CPUs] = 170.365PF

Lonestar6 Resource¹⁶

CPU: AMD "Milan" EPYC 7763 processors

1. $\text{AVX } 256 \text{ [256 bits]} / 16 \text{ [FP16]} * 2 \text{ [multiple and add; 2 operations]} * 2 \text{ [2 AVX 256 blocks per core]} * 64 \text{ [64 cores per CPU]} * 2.55 \text{ Ghz} = 10444.8\text{GF}$
 - a. $10444.8\text{GF} * 2 \text{ [2 CPUs per node]} = 20889.6 \text{ GF/node}$
 - b. $20889.6 \text{ GF/node} * 530 \text{ [total nodes]} = 11071488\text{GF} = 11.071488\text{PF}$

GPU:

1. $255 \text{ A100s} * 312\text{TF} \text{ [312TF per A100 without sparsity]}^{17} = 79560\text{TF} = 79.56\text{PF}$
2. $8 \text{ H100s} * 756\text{TF} \text{ [756TF per H100 without sparsity]}^{18} = 6048\text{TF} = 6.048\text{PF}$
3. Total: $79.56\text{PF} + 6.048\text{PF} = 85.608\text{PF}$

Total FLOPS:

2. 11.071488PF [from CPUs] + 85.608PF [from GPUs] = 96.68PF

Neocortex Resource¹⁹

ASICs: Two Cerebras CS-2s²⁰

1. $5780\text{TF} \text{ [per CS-2]} * 2 \text{ [2 CS-2s]} = 11.560\text{PF}^{21}$

CPUs: 32 Intel Xeon Platinum 8280L CPUs²²

2. $\text{AVX } 512 \text{ [bits]} / 16 \text{ [FP16]} * 2 \text{ [multiple and add; 2 operations]} * 2 \text{ [2 AVX 512 blocks per core]} * 28 \text{ [28 cores per CPU]} * 2.7 \text{ Ghz} = 9676.8\text{GF}$

¹⁶ <https://docs.tacc.utexas.edu/hpc/lonestar6/#intro>

¹⁷

<https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a100-datasheet-nvidia-us-2188504-web.pdf>

¹⁸ <https://resources.nvidia.com/en-us-tensor-core/nvidia-tensor-core-gpu-datasheet>

¹⁹ <https://www.cmu.edu/psc/aibd/neocortex/>

²⁰ <https://f.hubspotusercontent30.net/hubfs/8968533/CS-2%20Data%20Sheet.pdf>

²¹ <https://www.alcf.anl.gov/alcf-ai-testbed>

²²

<https://www.intel.com/content/www/us/en/products/sku/192472/intel-xeon-platinum-8280l-processor-38-5m-cache-2-70-ghz/specifications.html>

3. $9676.8\text{GF} * 32 \text{ [32 total CPUs]} = 309657.6\text{GF} = 0.30966\text{PF}$

Total FLOPS:

4. $11.560\text{PF} \text{ [from two Cerebras CS-2s]} + 0.30966\text{PF} \text{ [from CPUs]} = 11.87\text{PF}$

H100 Equivalent to Total NAIRR Compute

1. Total NAIRR compute = $3770.04\text{PF} = 3770040\text{TF}$
 - a. Note that some of the FP16 FLOPS numbers in the total NAIRR compute are not included in the above calculations because they are provided directly by sources.
See Table 1 in the blog post for more details.
2. $3770040\text{TF} / 756\text{TF} \text{ [756TF per H100 in tensor FP16 without sparsity]} = 4986.8$
H100s