

DS 7347

# High-Performance Computing (HPC) and Data Science

## Session 4

---

Robert Kalesky

Adjunct Professor of Data Science

HPC Research Scientist

May 5, 2022

Research and Data Sciences Services

Office of Information Technology

Center for Research Computing

Southern Methodist University



Session Question

Accelerators

Networks and Storage

Group Discussion

Readings and Assignments

## Session Question

---

## Session Question



What HPC components can perform application-specific computations?

## Accelerators

---



- 4992 NVIDIA CUDA cores with a dual-GPU design
- 2.91 TFLOPS double-precision
- 8.73 TFLOPS single-precision
- 24 GB of GDDR5 memory
- 480 GB/s aggregate memory bandwidth
- ECC protection for increased reliability

# NVIDIA Tesla P100 with PCIe



- 3584 NVIDIA CUDA cores
- 4.7 TFLOPS double-precision
- 9.3 TFLOPS single-precision
- 18.7 TFLOPS half-precision
- 16 GB HBM2
- 732 GB/s aggregate memory bandwidth
- 32 GB/s PCIe Gen3





- 3584 NVIDIA CUDA cores
- 5.3 TFLOPS double-precision
- 10.6 TFLOPS single-precision
- 21.1 TFLOPS half-precision
- 16 GB HBM2
- 732 GB/s aggregate memory bandwidth
- 160 GB/s NVLink

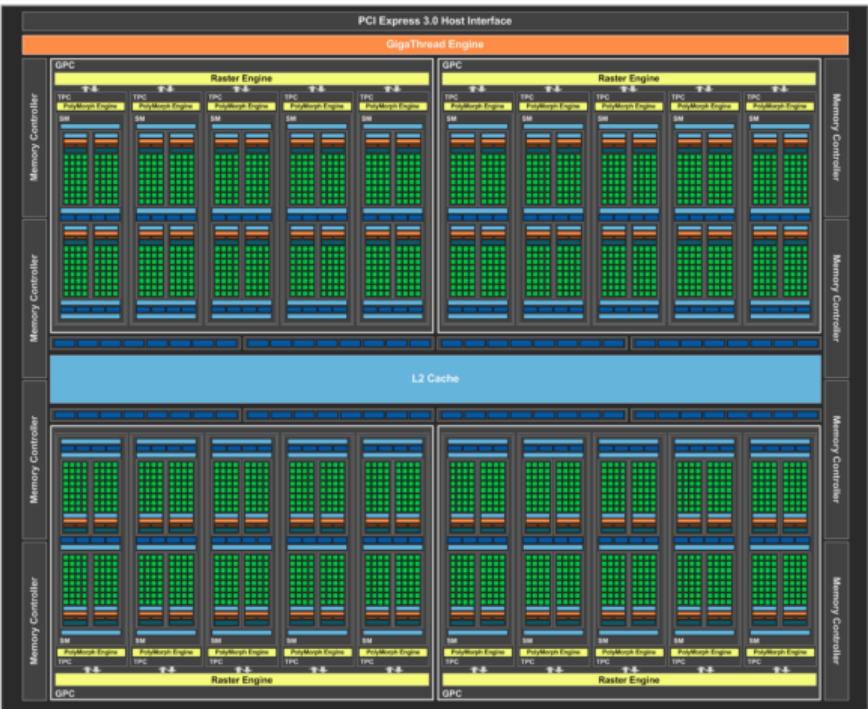


- 5120 NVIDIA CUDA cores
- 640 NVIDIA Tensor Cores
- 7.8 TFLOPS double-precision
- 15.7 TFLOPS single-precision
- 125 TFLOPS half-precision
- 32 GB HBM2
- 900 GB/s aggregate memory bandwidth
- 160 GB/s NVLink

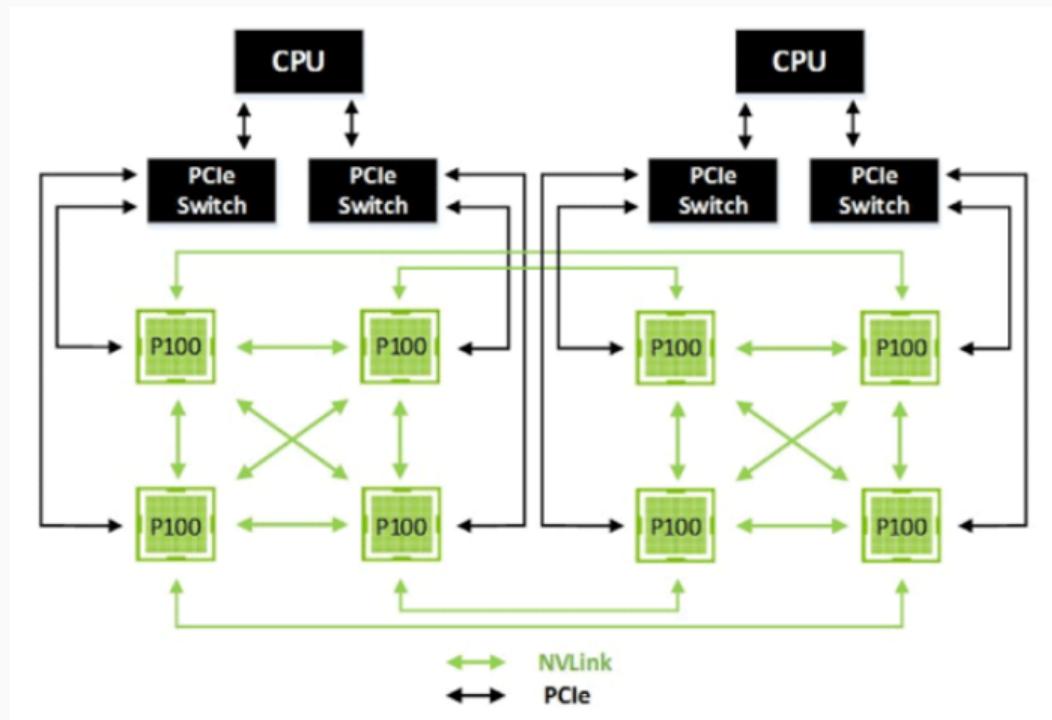


- 6912 NVIDIA CUDA cores
- 432 NVIDIA Tensor Cores (Third Generation)
- 9.7 TFLOPS double-precision
- 19.5 TFLOPS double-precision (Tensor Cores)
- 19.5 TFLOPS single-precision
- Tensor Float (TF32): 312 dense, 624 sparse
- BFLOAT16: 312 dense, 624 sparse TFLOPS (Tensor Cores)
- INT8: 624, 1248 TOPS (Tensor Cores)
- 80 GB HBM2e
- 2039 GB/s aggregate memory bandwidth
- 600 GB/s NVLink

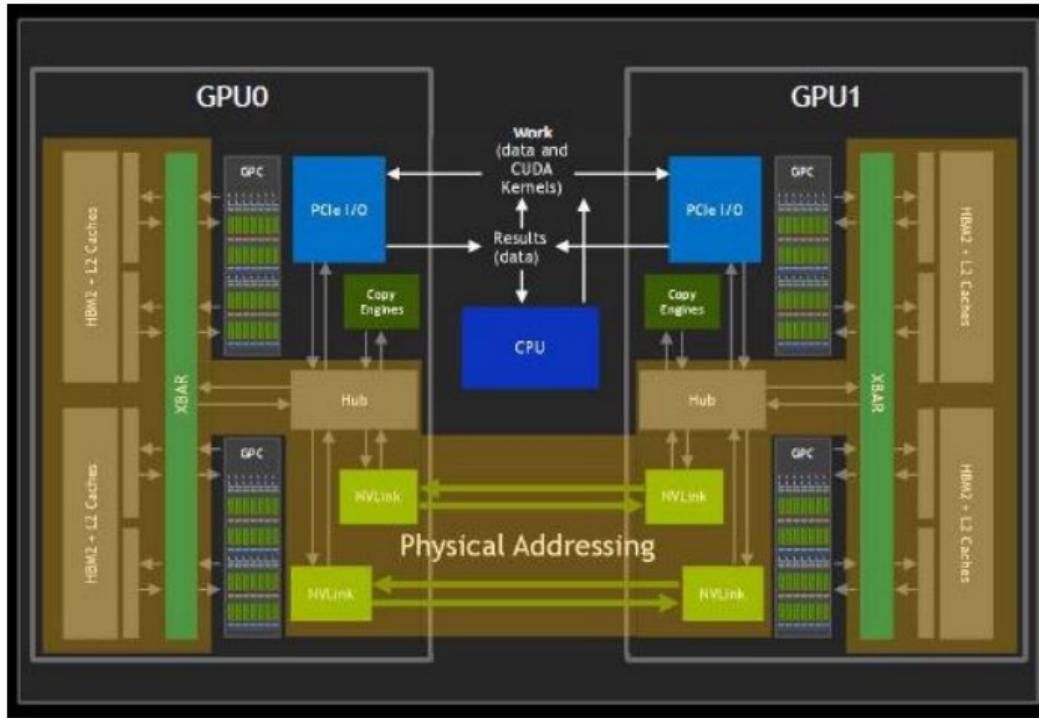
# NVIDIA Tesla Block Diagram



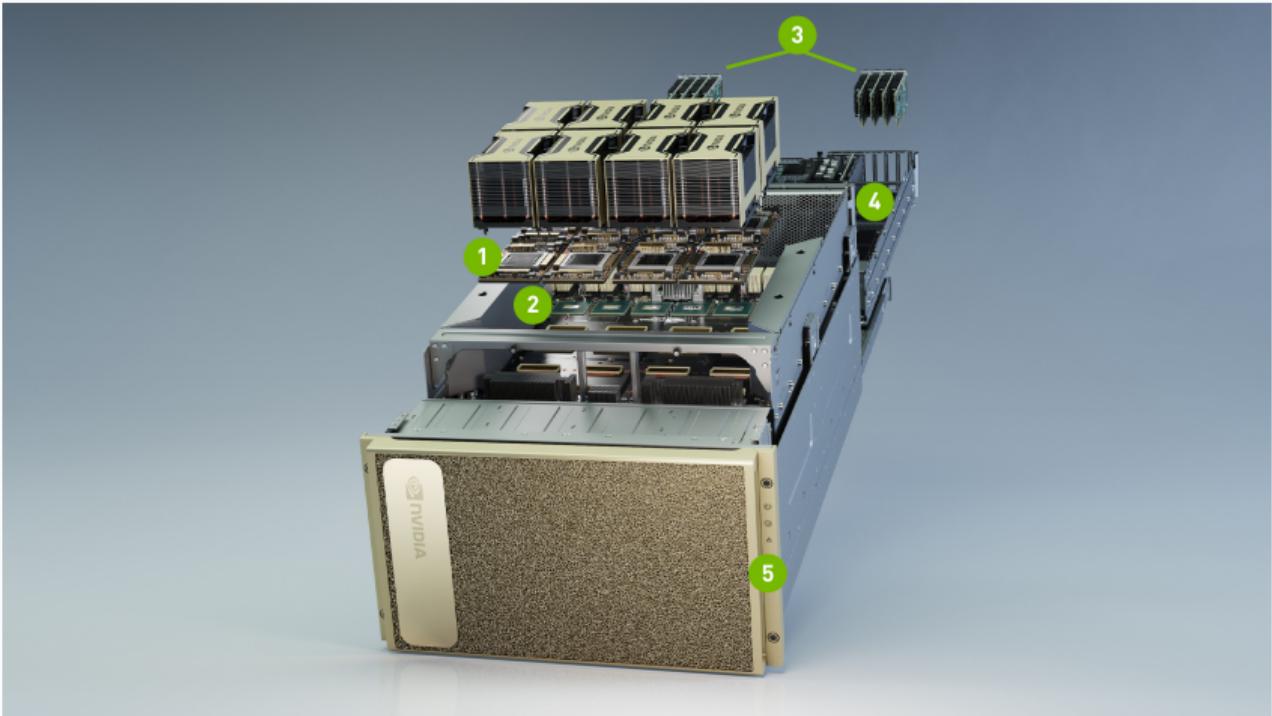
# NVIDIA NVLink and PCIe



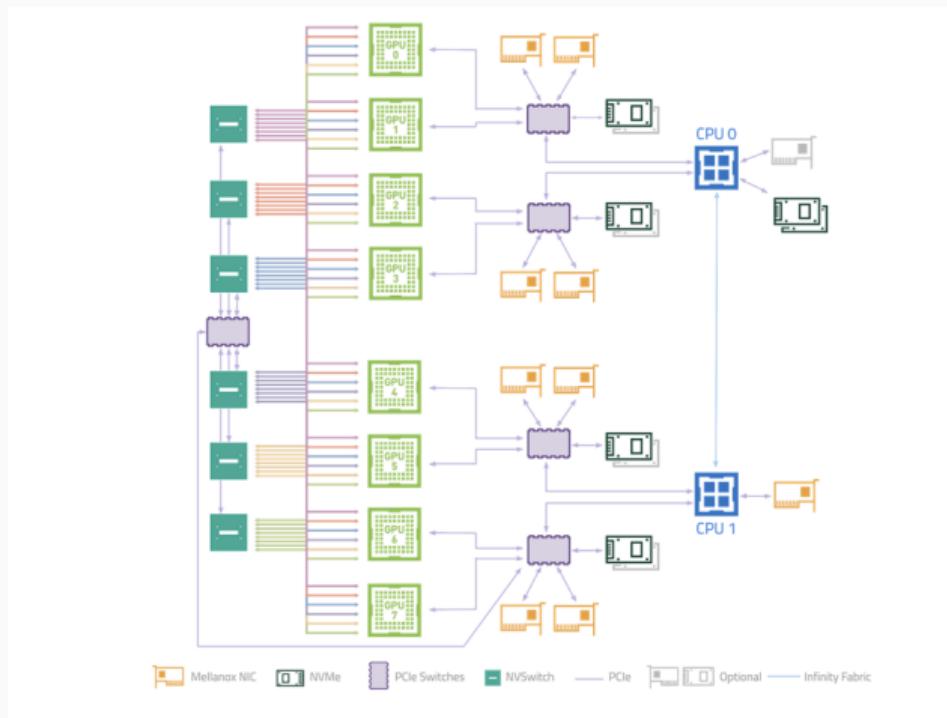
# NVIDIA NVSwitch Shared Memory



# NVIDIA SuperPOD DGX Node



# NVIDIA SuperPOD DGX Node



## Networks and Storage

---



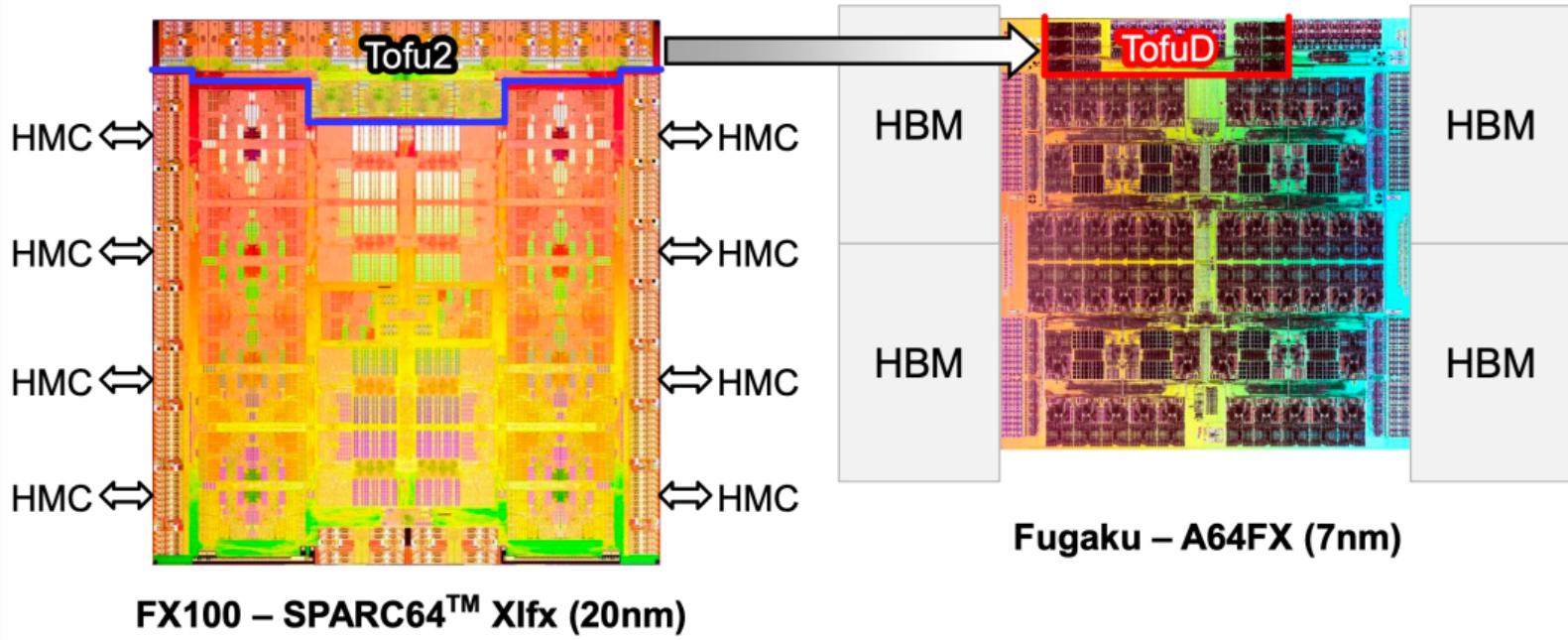
- Management Networks (2): 1 Gb/s Ethernet
- Network File System (NFS) Network: 10 Gb/s Ethernet
- High-Performance Network: > 50 Gb/s InfiniBand

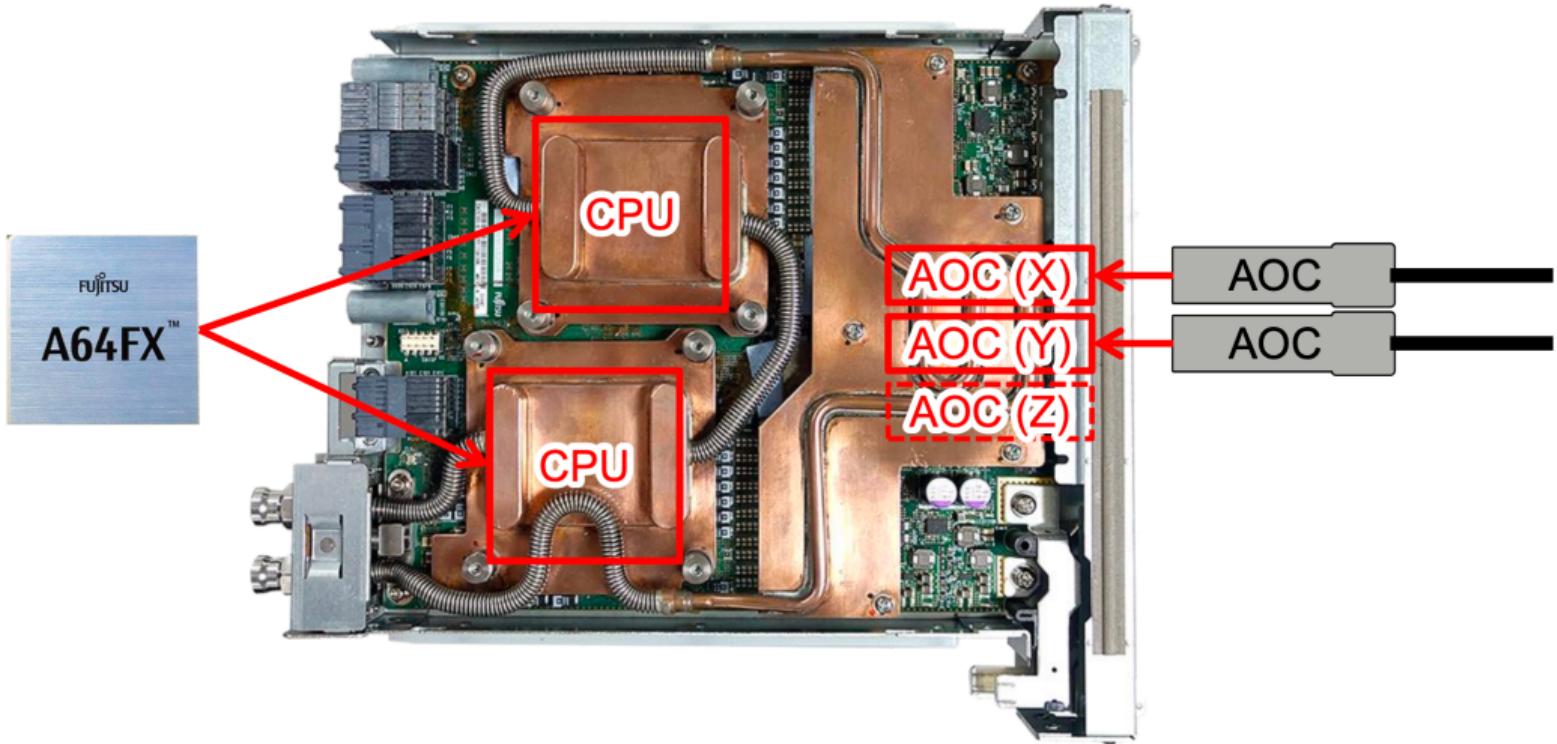


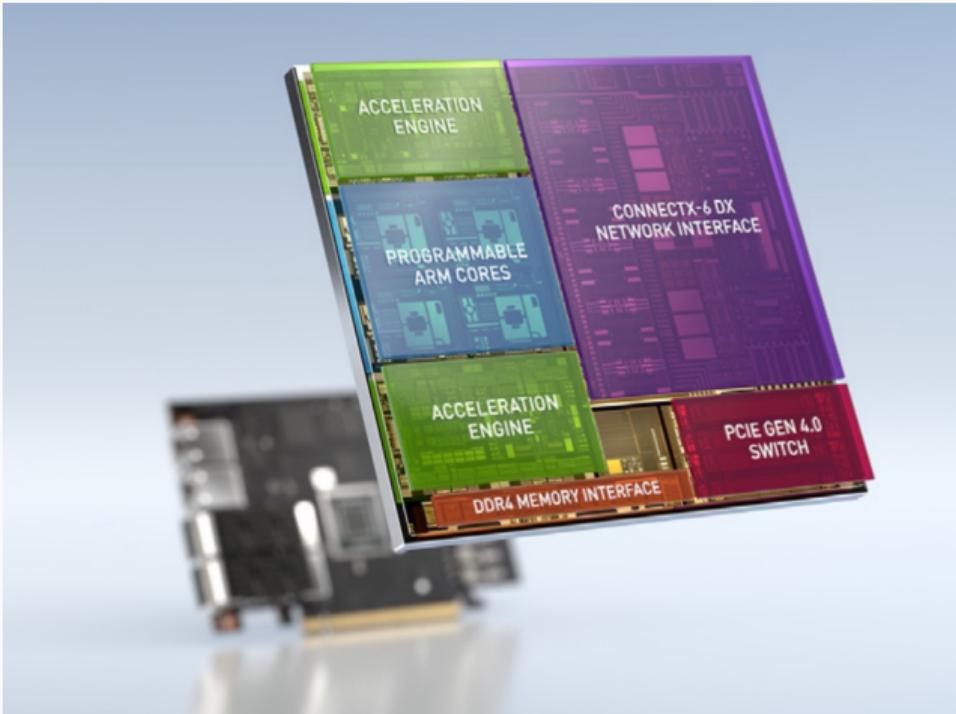
## Connected with Four Links

- QDR: 32 Gb/s
- FDR: 54 Gb/s
- EDR: 100 Gb/s
- HDR: 200 Gb/s
- NDR: 400 Gb/s

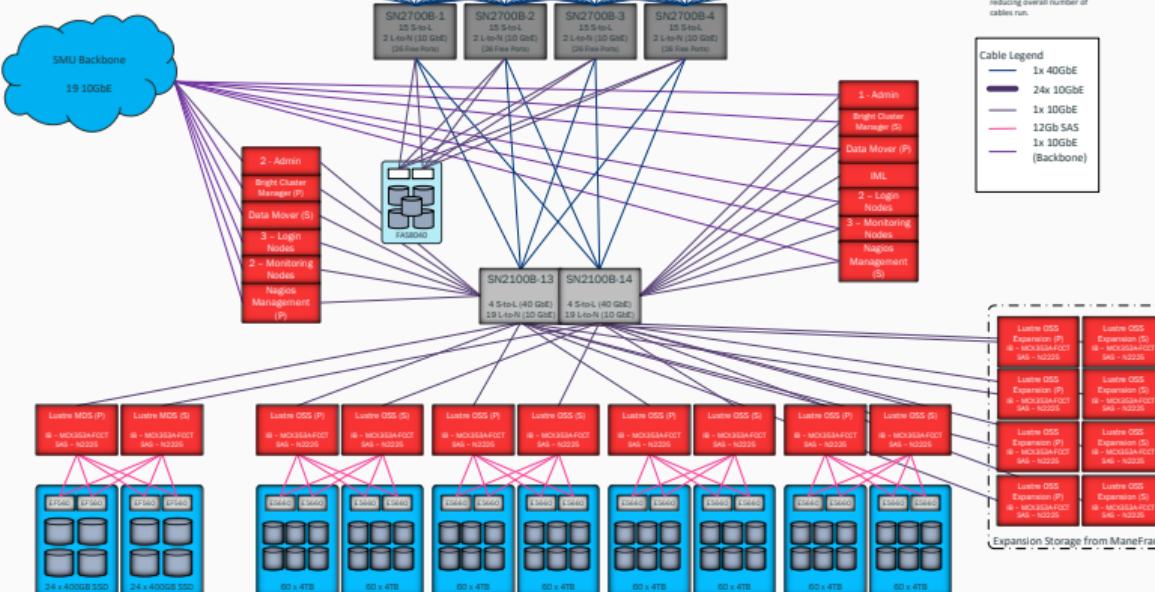
# K-Computer and Fugaku Processors







# M2 Network and Storage



## Group Discussion

---



## Items to Note for Path

- File systems
- Networks
- Interconnect
- Memory



## Trace the Path

1. Run Python script from **\$HOME**
2. Python script:
  - 2.1 Reads compressed data from files on **\$WORK**
  - 2.2 Decompresses the data the CPU
  - 2.3 Performs computation on the data that generates intermediate data stored to **\$SCRATCH**
  - 2.4 Launches computation on the GPU:
    - 2.4.1 Reads data “directly” from **\$SCRATCH**
    - 2.4.2 Sends results back to Python script
  - 2.5 Writes data to **\$HOME**

## Readings and Assignments

---



## Readings

- Eijkhout chapter 22



## Assignment

- Detail the CPU, GPU, memory, and hard drive for your own computer:
  - CPU
    - Model
    - Number of cores
    - Instruction set architecture (ISA)
    - Notable vector extensions available
    - Cache sizes
  - GPU
    - Model
    - Processing cores
    - Memory
    - Interconnect
  - Memory
    - Speed
    - Size
  - Hard drive
    - Size
    - Interconnect