

# AI/ML computations on SDSC's Expanse Cluster

Presented by Paul Rodriguez and Mahidhar Tatineni

Fifth National Research Platform (5NRP) Workshop

March 19, 2024

*University of California San Diego & San Diego  
Supercomputer Center*



# EXPANSE

COMPUTING WITHOUT BOUNDARIES  
5 PETAFLOP/S HPC and DATA RESOURCE

## HPC RESOURCE

13 Scalable Compute Units  
728 Standard Compute Nodes  
52 GPU Nodes: 208 GPUs  
4 Large Memory Nodes

## LONG-TAIL SCIENCE

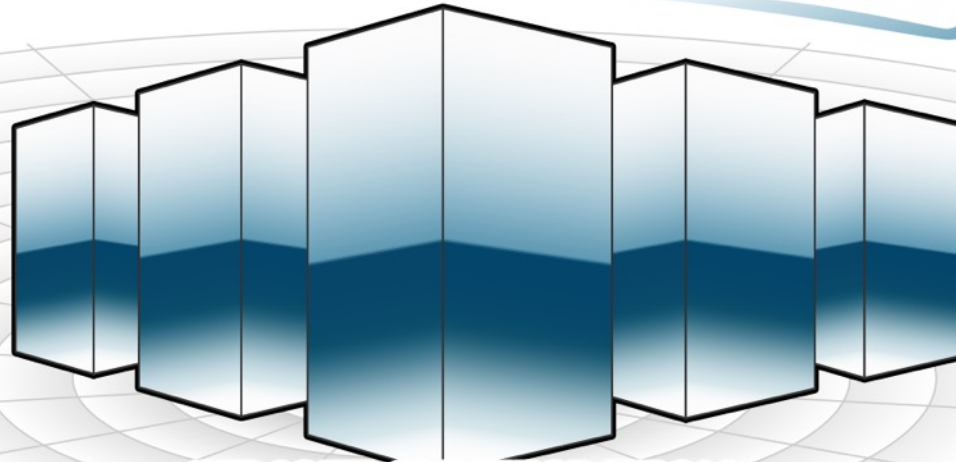
Multi-Messenger Astronomy  
Genomics  
Earth Science  
Social Science

## DATA CENTRIC ARCHITECTURE

12PB Perf. Storage: 140GB/s, 200k IOPS  
Fast I/O Node-Local NVMe Storage  
7PB Ceph Object Storage  
High-Performance R&E Networking

## INNOVATIVE OPERATIONS

Composable Systems  
High-Throughput Computing  
Science Gateways  
Interactive Computing  
Containerized Computing  
Cloud Bursting



REMOTE CI INTEGRATION

CLOUD

Open Science Grid

Heterogeneous Resources

NSF Award # 1928224

PIs: Mike Norman (PI), Ilkay Altintas, Amit Majumdar, Mahidhar Tatineni, Shawn Strande

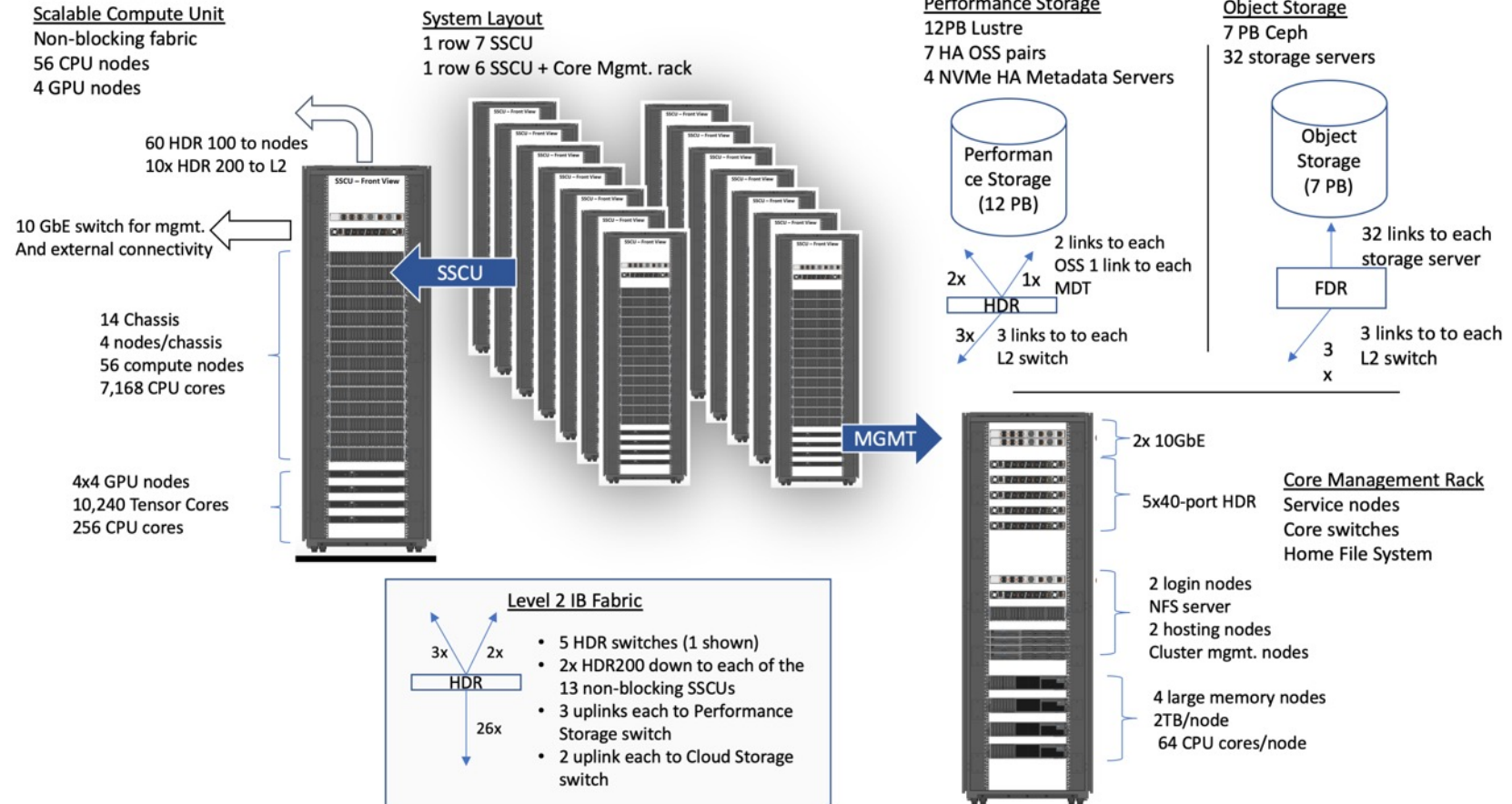
# Expanse Overview

- **Original Configuration (13 racks):**
  - Category 1: Capacity System, NSF Award # 1928224
  - 728, 2-socket AMD Rome-based compute nodes (2.25 GHz EPYC; 64-core/socket). 93,184 compute cores in total.
  - 52 4-way GPU nodes with V100 GPUs w/NVLINK
  - Allocated via ACCESS
- **Partnership to Advanced Throughput Computing (PATH) racks (2):**
  - 112 2-socket AMD Milan-based compute nodes; 512 GB of memory per node
  - 8 4-way GPU nodes w/ A100 GPUs
- **Industry rack (funded by UCSD/SDSC):**
  - 56 2-socket AMD Rome-based compute nodes
  - 4 4-way GPU nodes based on V100 w/NVLINK
- **System integrated by Dell**

# Expanse is a heterogeneous architecture designed for high performance, reliability, flexibility, and productivity

## System Summary

- 14 SDSC Scalable Compute Units (SSCU)
- 784 x 2s Standard Compute Nodes
- 100,352 Compute Cores
- 200 TB DDR4 Memory
- 56x 4-way GPU Nodes w/NVLINK
- 224 V100s
- 4x 2TB Large Memory Nodes
- HDR 100 non-blocking Fabric
- 12 PB Lustre High Performance Storage
- 7 PB Ceph Object Storage
- 1.2 PB on-node NVMe
- Dell EMC PowerEdge
- Direct Liquid Cooled



# The SSCU is Designed for the Long Tail Job Mix, Maximum Performance, Efficient Systems Support, and Efficient Power and Cooling

## Standard Compute Nodes

- 2x AMD EPYC 7742 @2.25 GHz
- 128 Zen2 CPU cores
- PCIe Gen4
- 256 GB DDR4
- 1 TB NVME

## GPU Nodes

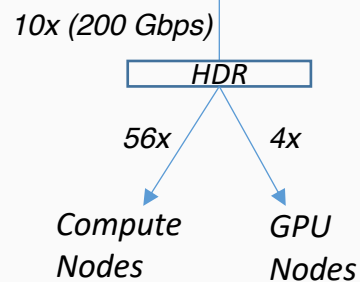
- 4x NVIDIA V100/follow-on
- 10,240 Tensor Cores
- 32 GB GDDR
- 1.6 TB NVMe
- Intel CPUs

## SSCU Components

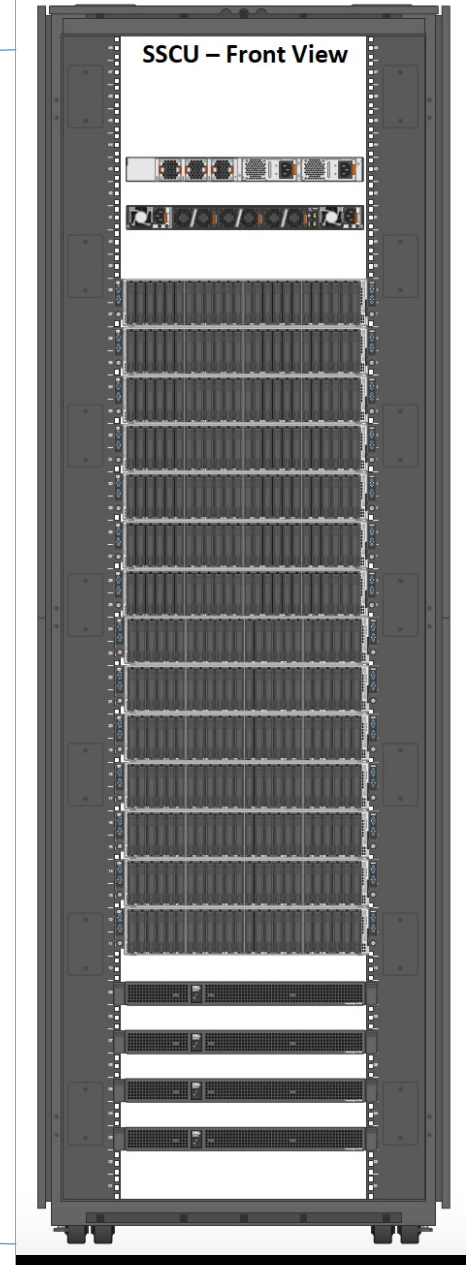
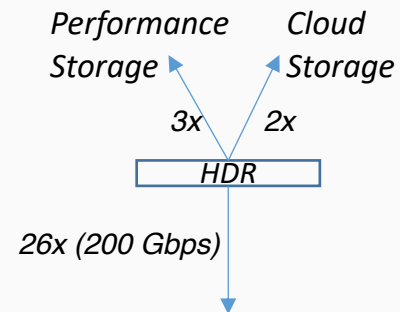
- 56x CPU nodes
- 7,168 Compute Cores
- 4x GPU nodes
- 1x HDR Switch
- 1x 10GbE Switch
- HDR 100 non-blocking fabric
- Wide rack for serviceability
- Direct Liquid Cooling to CPU nodes

## Non-blocking Interconnect

### 1 HDR Switch/SSCU



### 5 Level 2 switches





# Running Jobs, Software on Expanse

- Expanse uses the **Simple Linux Utility for Resource Management (SLURM)** batch environment
- **Primary method to run jobs:** Submit batch scripts from the login nodes
- **Expanse User Portal:** Useful for interactive MATLAB, Rstudio GUI enabled jobs; Jupyter Notebooks (enabled using our Satellite service and Galyleo)
- **Software Stack**
  - Modules environment for applications/libraries installed by SDSC staff
  - Singularity Containers - SDSC staff maintain containers for some applications (e.g. TensorFlow, PyTorch)
  - User installed - using miniconda3 for example

# TensorFlow and PyTorch on Expanse

- Two main approaches:
  - **Singularity container images** with all the python packages included, along with GPU drivers, CUDA libraries. Examples:
    - /cm/shared/apps/containers/singularity/tensorflow/tensorflow-latest.sif
    - /cm/shared/apps/containers/singularity/pytorch/pytorch-latest.sif
  - **Conda/Miniconda based installs** (typically users have their own custom versions).
    - If you install a GPU version - make sure the version of cuda, cuda toolkit in the miniconda install is compatible with our driver.

# Machine Learning/Deep Learning on Expanse via Singularity

- **Machine learning/deep learning applications on Expanse primarily made available via Singularity images.**
  - These packages are constantly being upgraded and the dependency list is difficult to update in the standard Expanse environment.
- **Install options**
  - **Singularity image provides dependencies and user can compile actual application from source.**
  - **Entire dependency stack and the application is in the image.**
- **Run options**
  - **Most cases are run on single GPU nodes (4 GPUs at most)**
  - **Can access this via Jupyter notebooks**
  - **Multi-node options are possible via singularity.**