



Strategic Plan for Software and Computing at the Laboratory

James Amundson

2022 June PUBLIC PAC Meeting

June 23, 2022

Computational Science and Artificial Intelligence at Fermilab

(Not Core Computing)

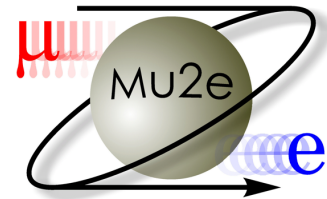
Objective: Support the scientific mission of the laboratory



Maximizing Fermilab's Scientific Output

By the end of the decade, Fermilab's experimental program will be dominated by DUNE and HL-LHC, with significant contributions from the Short Baseline Neutrino program and Mu2e.

- In addition:
 - small experiments
 - potential for new customers, especially cosmic



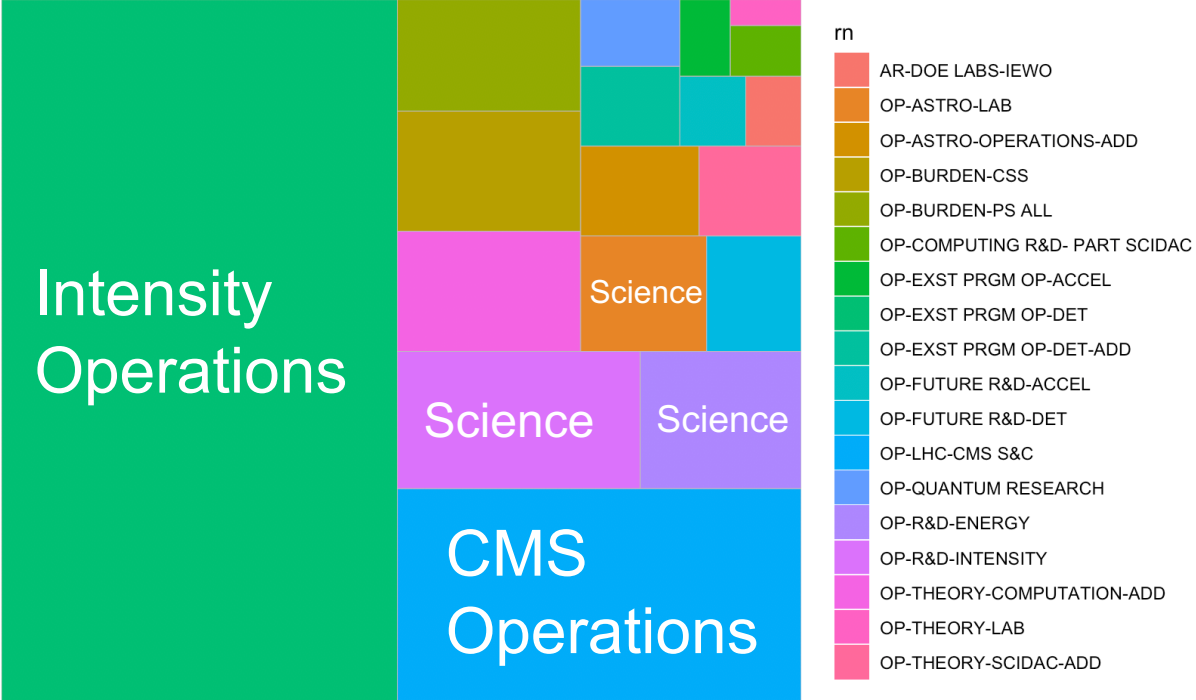
Computing Strategy

100,000 ft. view

- Maintain core Fermilab computing facility
 - Focus on things that cannot be done as well or better elsewhere
 - Mass storage is the core of Fermilab's computing facility
- Take maximum advantage of non-HEP resources
 - DOE Advanced Scientific Computing Research (ASCR)
 - Exascale/HPC Computing resources
 - Software
 - NSF Supercomputing resources
 - Other resources (e.g., Open Science Grid)
 - Commercial resources
- Embrace AI/ML developments
 - Enable scientific AI/ML applications
 - Utilize AI/ML across the lab
- Support computational science as appropriate
 - In particular, simulation
 - Assist the field in modernization of computing: GPUs, specialized services



Funding Support for Scientific Computing at Fermilab

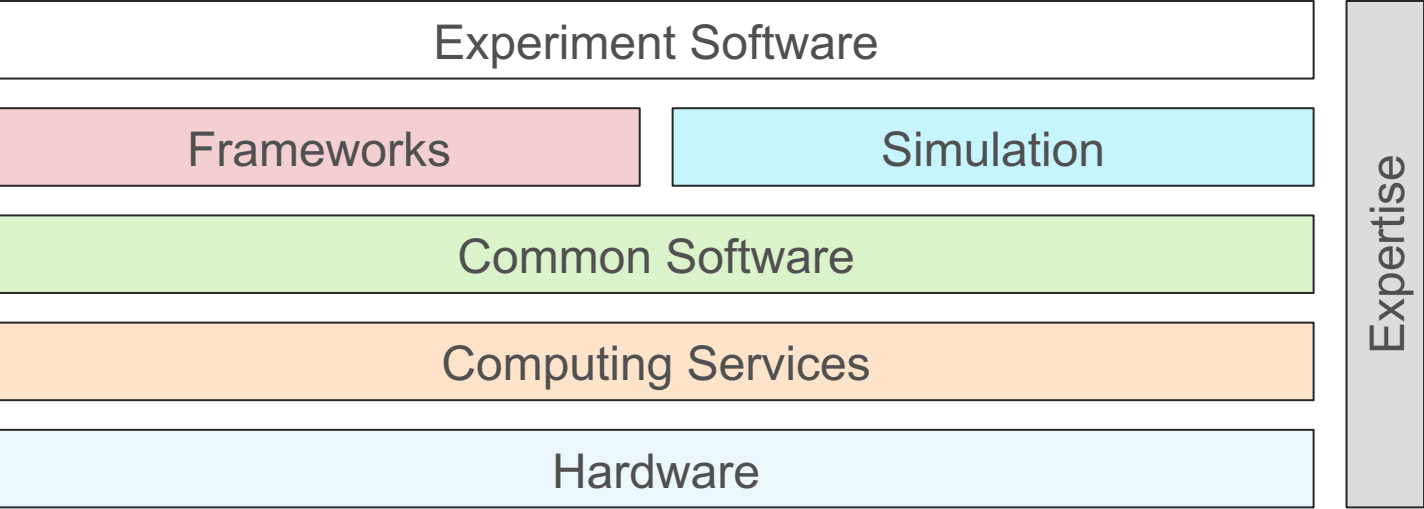


**n.b. not FY22*

Computing Across Frontiers

- Intensity Frontier
 - NOvA, MicroBooNE, etc.
 - DUNE, SBN
 - Muon g-2, Mu2e
- Energy Frontier
 - CMS
- Cosmic Frontier
 - DES, Rubin, etc.
- Theory
 - Lattice QCD
 - Generators
- Accelerator
 - Simulation
- Other
 - Primarily R&D

Layers



Foundational Layers

Hardware

- CPU, GPU and Storage Resources
 - Both on- and off-premises
- Accessing resources requires...

Computing Services

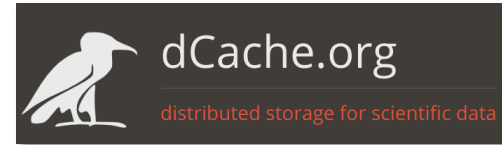
- FIFE for traditional resources
 - Not discussed today
- HEPCloud for Cloud and HPC

Hardware

	Compute	Storage
Fermilab	CPU, some GPU	Tape, Disk
HPC Center	CPU, GPU	
Open Science Grid	CPU, some GPU	
Collaborators	CPU, some GPU	Disk
Cloud	CPU, GPU	

Storage Research and Development

- Fermilab has selected CERN's CTA as a replacement for Enstore in tape layer
 - Informal agreement to collaborate with CERN
 - Formal agreement in the works
- Evaluating multiple technologies in the disk layer
 - dCache
 - Existing collaboration with DESY
 - Integrating with CTA
 - EOS
 - Currently required by CTA, used by EOS
 - ceph
 - Broad usage in multiple industries
 - Could replace very expensive NAS storage
- Emphasizing Rucio within software layer
 - Broad community support
 - Provides mechanism to enforce data lifetimes
 - Experience shows that manual lifetime management is not realistic



Storage R&D: Tape

- Following two general thrusts: tape/archival storage evolution and disk storage evolution
- Tape: **replacing legacy Enstore system with CERN Tape Archive (CTA)**
- Development/changes to CTA necessary
 - Enable CTA to read tapes with **CPIO wrapper** that most Enstore tapes are written in
 - Develop metadata migration for Enstore->CTA
 - **Small File Aggregation (SFA)** functionality replacement
 - Joint with DESY: dCache frontend for CTA
- Current status
 - Running CTA on partition of Fermilab tape library
 - Able to read/write Enstore formatted tapes with CTA
 - Fermilab team has made CTA code commits
 - Remaining significant item for all libraries: metadata migration
 - Remaining item for Public Enstore: SFA solution with CTA
 - Development of this will continue in parallel with migration beginning for CMS



Tape Evolution Timeline

- Approximate Enstore->CTA migration plan
 - Intend to start with CMS tape library (fewer tape families and no SFA)
 - 3Q'22
 - Implement and test CTA metadata scheme for Enstore tapes
 - Architect/procurement of servers for production CTA
 - 4Q'22
 - Configure/test CTA servers
 - Begin metadata migration from Enstore to CTA
 - 1Q'23
 - New data ingest to CTA only for CMS
 - Begin process for Public Enstore

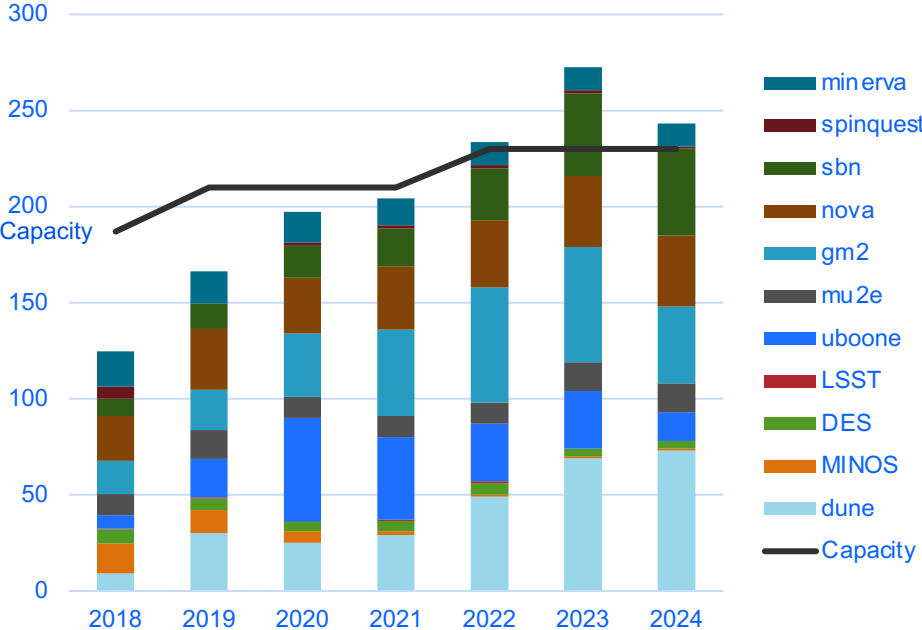
Very preliminary

Storage R&D: Disk

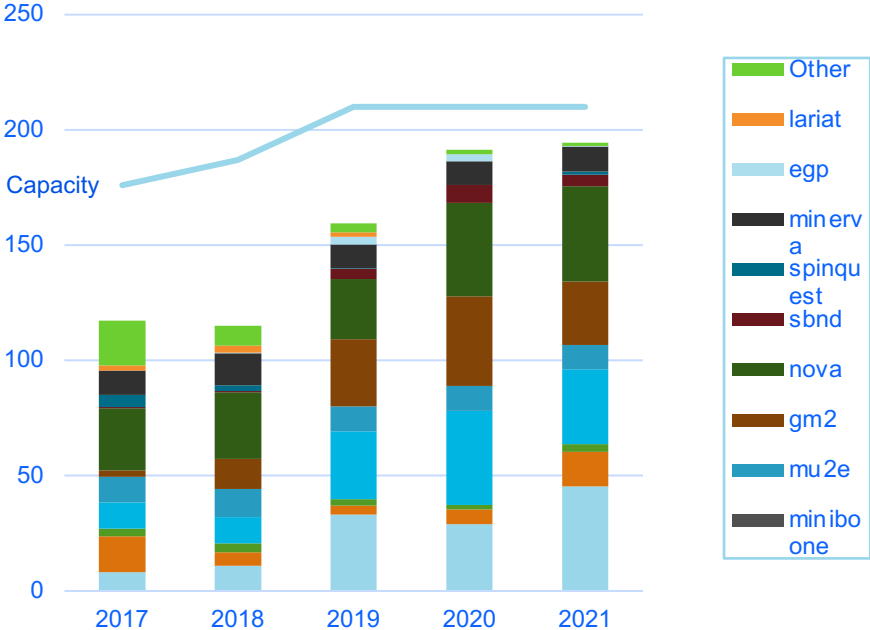
- Evaluation of **Ceph** as a disk-based storage solution with several applications
 - CephFS-based replacement for NAS storage currently used for interactive computing
 - Ceph Object Store for analysis computing and other applications
- Current status
 - 1PB Ceph test cluster configured for CephFS and Object store use
 - CMS object store evaluation (USCMS Ops program funded)
 - Interactive computing evaluation of CephFS (starting with DESI/LSST)
- Future
 - Test cluster with newer hardware
 - Evaluating SSD/HDD mixture and Erasure Coding schemes

Non-CMS Experiment CPU Requests and Usage

Non-CMS MCore Hour Requests



Actual Usage



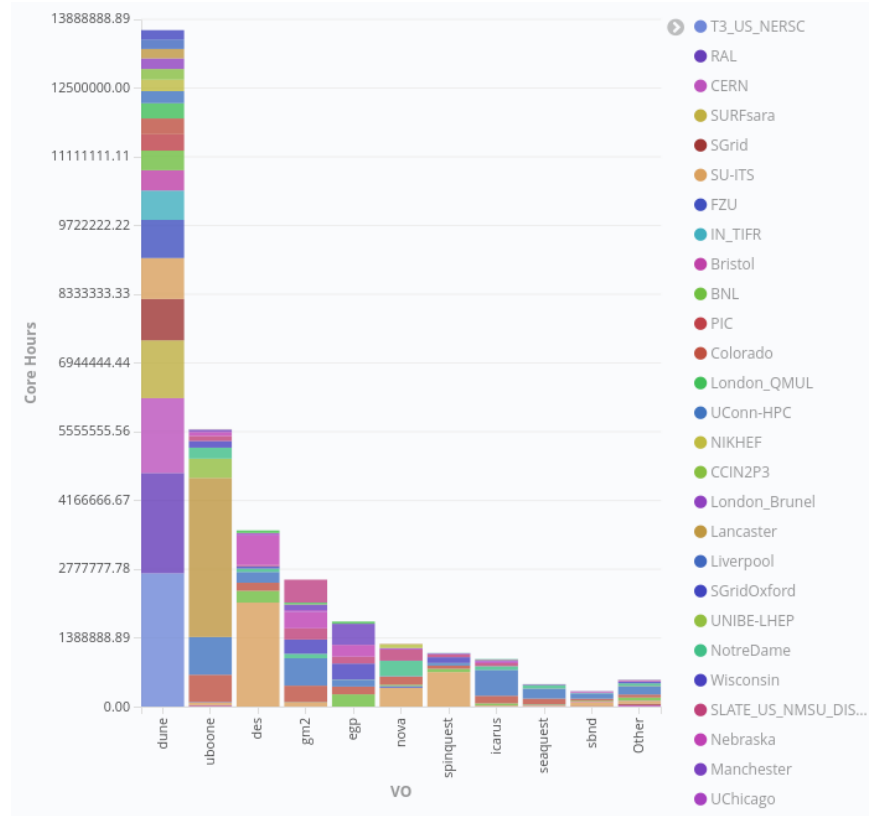
Fermilab GPU Resources

- Bootstrap problem
 - Experiments do not request GPUs if they think we do not provide them
 - We do not provision GPUs if experiments do not request them
- Annual request to DOE for funding for initial production GPU facility
 - Positive feedback
 - No funding to date
- Slowly increasing GPU purchases with portion of existing funds
 - 12 A100 GPUs purchased with FY21 funds

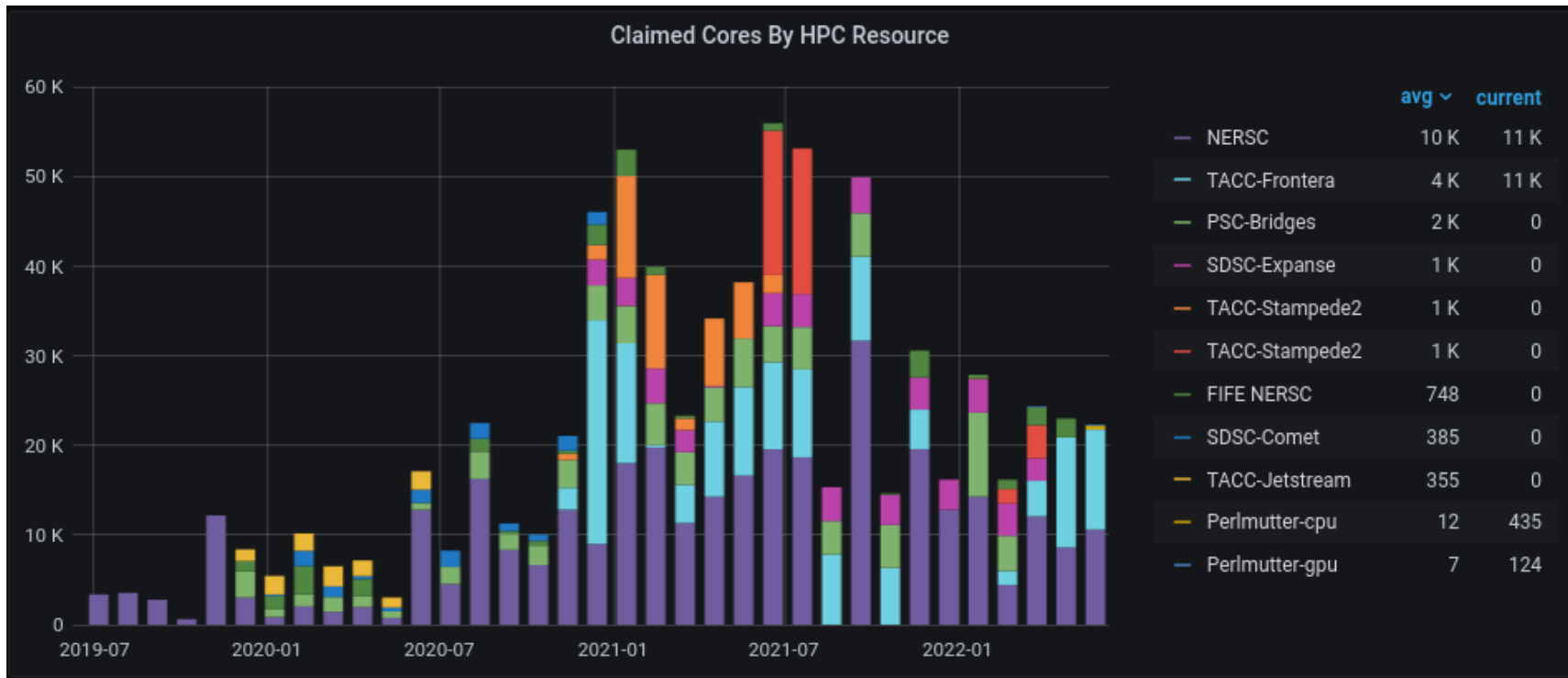
	FY21	FY22	FY23	FY24
Planned GPU acquisitions (thousands of NVIDIA A100-equivalent hours)	105	35	70	70
Planned retirements (thousands of NVIDIA A100-equivalent hours)	2	20	30	25
GPU capacity (thousands of NVIDIA A100-equivalent hours)	250	265	305	350

Non-Fermilab CPU Usage

- HPC sites (allocations)
- OSG (opportunistic)
- GCE, AWS (paid)
- If experiments have special agreements with collaborating sites, we can enable access to their individual allocations
- Containers should limit issues at remote sites
- Some VOs could push more offsite

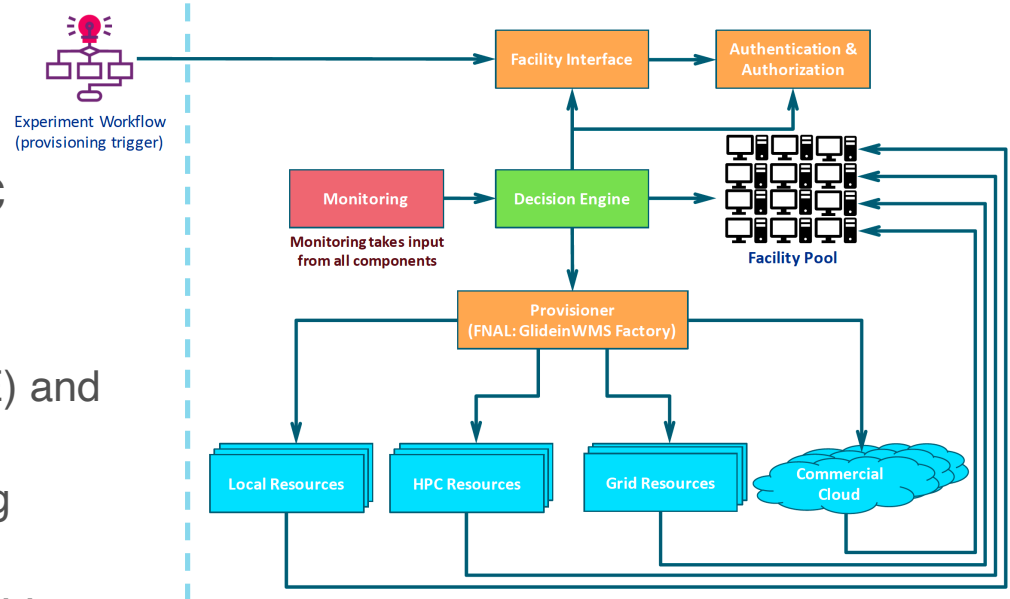


HPC Usage



HEPCloud

- HEPCloud is our solution for accessing a heterogeneous set of resources, including cloud and HPC
- HEPCloud is currently running in production
 - HPC centers including NERSC (DOE) and TACC (NSF)
 - Commercial cloud providers including Google and Amazon Web Services
- Progress on Leadership Class Facilities
 - Argonne (ALCF) and Oak Ridge (OLCF)
 - Two solutions for problem of contacting nodes isolated from general network
 - Still ramping up production



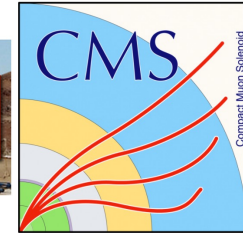
HEP-CCE

- Goal is to enable HEP on Exascale
- Funded by DOE CompHEP
- Multi-year project
- Multi-lab project
 - Fermilab
 - Argonne
 - Brookhaven
 - Lawrence Berkeley
- Multi-thrust project
 - Platform Portability
 - Device-independent approaches to GPUs
 - I/O
 - Workflows
 - Generators

High Energy Physics - Center for Computational Excellence

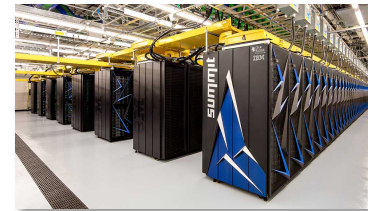
HEP-CCE

<https://www.anl.gov/hep-cce>



DOE HPC/Exascale Resources

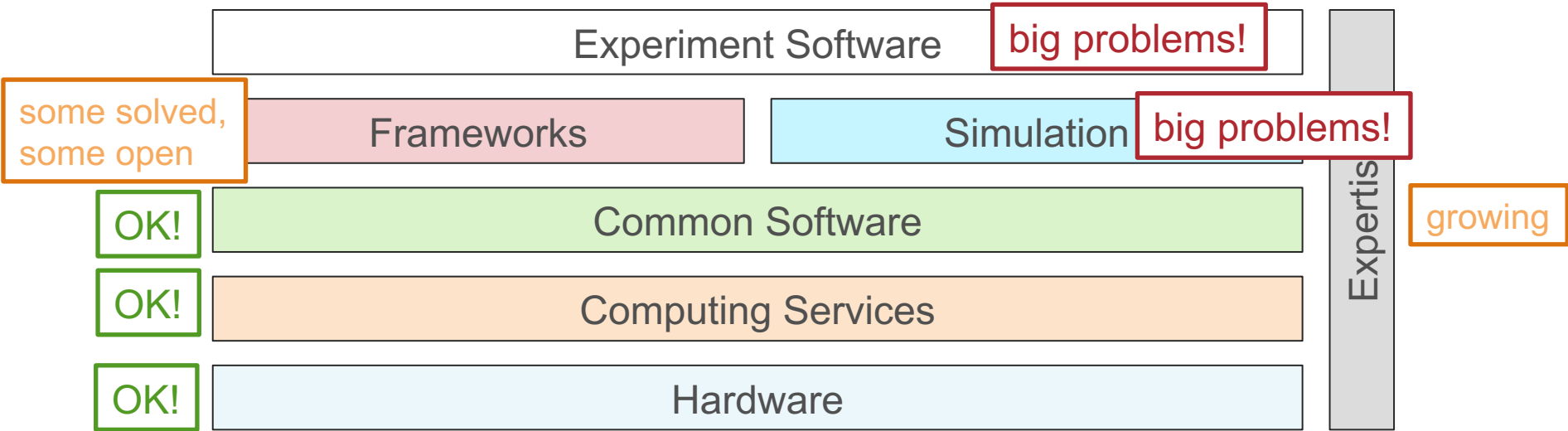
- NERSC
 - Cori
 - Haswell, 2,388 nodes, 2.81 PFlops
 - KNL, 9,688 nodes, 29.5 PFlops
 - Perlmutter (Phase 1)
 - AMD + NVIDIA, 1,536 nodes
 - 3.9 PFlops CPU
 - 59.9 PFlops GPU (94%)
- ALCF (Argonne)
 - Current: Theta (also ThetaGPU)
 - KNL, 4,392 nodes, 11.7 PFlops
 - Next: Aurora (exascale!)
 - Intel CPU + GPU, > 9,000 nodes, > 1,000 PFlops
- OLCF (Oak Ridge)
 - Summit
 - IBM Power9 + NVIDIA, 4,608 nodes, 200 PFlops
 - Frontier (exascale!)
 - AMD CPU + GPU, 1,100 PFlops



Barriers to Exascale for HEP

- Allocations
 - LCF Allocation mechanisms are not compatible with HEP computing
 - Political problem, not a technical problem
 - First ALCC grant for CMS recently awarded
 - tiny
- Job submission/authentication, etc.
 - HEPCloud!
- Workflow management
 - HEPCloud!
- Data access
 - Not yet limiting
 - Addressed by CCE
- GPU Utilization
 - Addressed by CCE

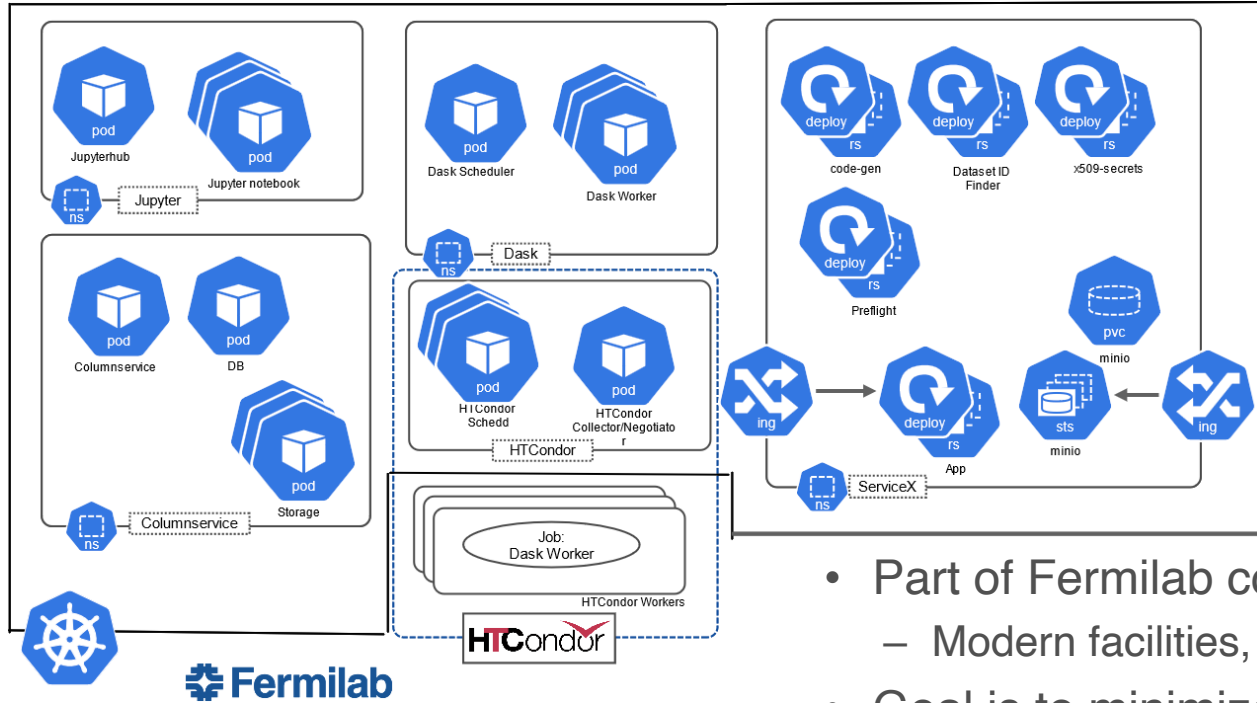
Problems with GPUs



What if HEP on Exascale Fails?

- Biggest problem is for CMS
 - CMS is making good progress
 - CMS Software and Computing is its own budget
- Intensity Frontier problems not imminent
 - Worst case scenario: more compute hardware would have to come out of Computing and Detector Operations budget
 - Assuming a flat budget, would require a corresponding reduction in staffing
 - In FY23, a computing professional costs \$422K
 - The best way to ensure success of Intensity Frontier on Exascale would be to find more funding for LArSoft

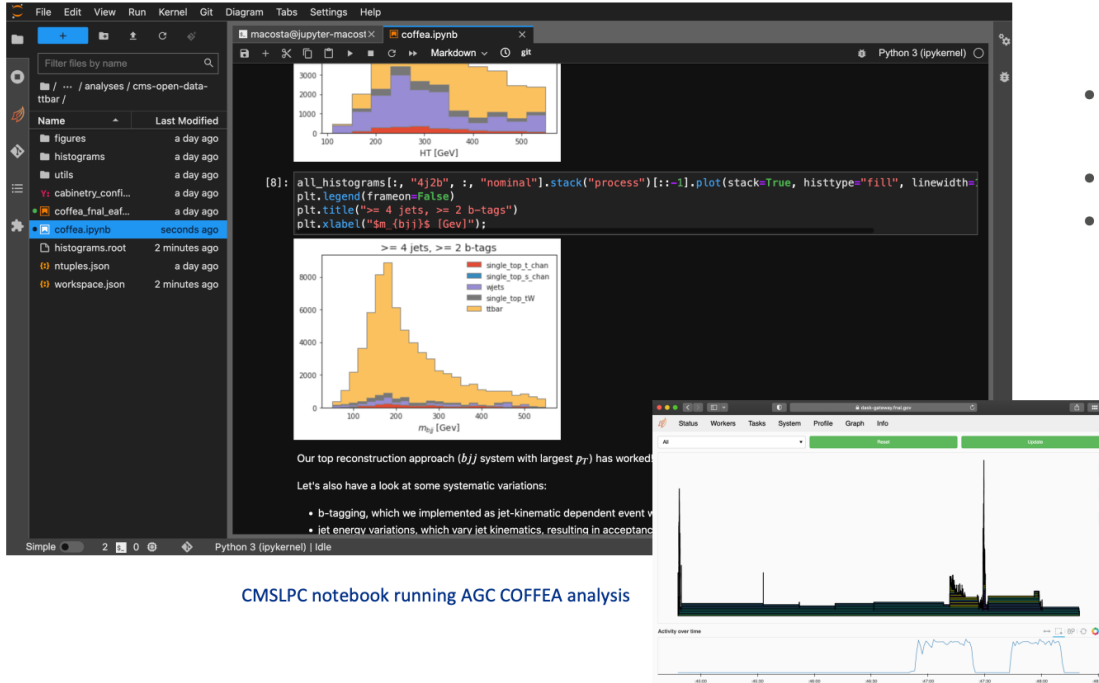
Elastic Analysis Facility



- Part of Fermilab computing strategy
 - Modern facilities, modern tools
- Goal is to minimize time to scientific insight
- Facility is available as a beta release

Elastic Analysis Facility Beta Release

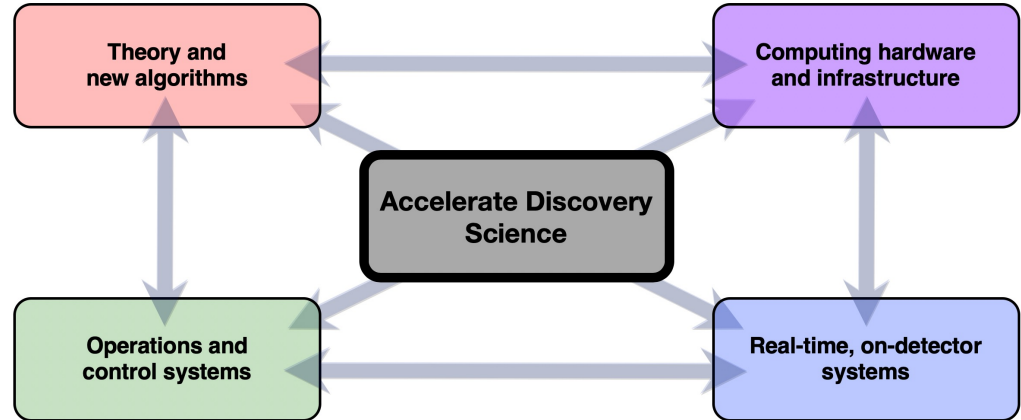
A JupyterHub-based deployment (on Beta) <https://analytics-hub.fnal.gov>



- 43 Beta users (thank you!)
- 22 Notebook flavors
- 1.2 Tb Ceph persistent storage allocated (of 45TB)

AI/ML at Fermilab

- AI/ML strategy at Fermilab extends across divisions
 - New dedicated AI/ML professional hired for accelerator operations
- Recent GPU purchases aimed at AI/ML use
- AI Associate program created for staffing purposes
 - Term positions
 - Do not require physics deliverables
 - Associates can spend time at the lab for career development
 - Recently had first permanent staff member come from associate program



AI/ML Projects at Fermilab

- Primarily driven by smallish funding opportunities
- Dominated by “Experiment Software” layer
- Long-awaited major funding opportunity recently materialized
 - DE-FOA-0002705: 2022 Artificial Intelligence Research for High Energy Physics
 - Not a game-changing level of funding
 - Fermilab led three proposals
 - Fermilab participated in six proposals led by other labs
 - Results pending

Support for Simulation at Fermilab

- Detector Simulation
 - Geant4 requires ongoing support
 - Computing
 - Physics Models
 - Experiment integration
 - Funding falls between the cracks in DOE OHEP
- Collider Generators
 - Robust program in both theory and computing
 - SciDAC4 Support
 - Mature community
- Neutrino Generators
 - Fundamentally more complicated than Collider Generators
 - Less mature community
 - Requires support in many areas
 - Nuclear and Particle Theory
 - Model integration
 - Experiment integration
 - Programming issues
 - Release management
 - Effort in Theory Division
 - Computing effort focused on GENIE
 - Steven Gardiner recently hired as associate scientist

Our Computing Strategy is the Conclusion

100,000 ft. view

- Maintain core Fermilab computing facility
 - Focus on things that cannot be done as well or better elsewhere
 - Mass storage is the core of Fermilab's computing facility
- Take maximum advantage of non-HEP resources
 - DOE Advanced Scientific Computing Research (ASCR)
 - Exascale/HPC Computing resources
 - Software
 - NSF Supercomputing resources
 - Other resources (e.g., Open Science Grid)
 - Commercial resources
- Embrace AI/ML developments
 - Enable scientific AI/ML applications
 - Utilize AI/ML across the lab
- Support computational science as appropriate
 - In particular, simulation
 - Assist the field in modernization of computing: GPUs, specialized services

