# Scientific Computing Division: Strategy and Updates

James Amundson

SCD All-Hands

November 10, 2021

# It has been a while

*I was hoping for a time when Covid policy would become crystal clear…*

- Three all-hands meetings
  - Today
    1. Division Strategy
    2. Miscellany
    3. Questions
  - December 10
    - Division accomplishments
  - Early 2022
    - Current experiments
      - Running
      - Analyzing

🔷 **Fermilab**

# 1. Division Strategy

Fermilab

# Based on presentation to the Physics Advisory Committee



**Strategic Plan for Software and Computing at the Laboratory**

James Amundson
2021 June PAC Meeting
June 8, 2021

# Charge

*We ask the committee to review the strategic plan for software and computing at the laboratory and the status of the recommendations made at the July 2020 PAC meeting:*

- "SCD should start thinking about the data lifecycle and its relation to long term data management and update the Committee at the next meeting. This report should include all relevant aspects, e.g., data storage, transfer, and compression".

- "The PAC continues to recommend [...] the development of a long-term computing transition plan" (to HPC)

Our responses to these two recommendations are integral parts of our overall strategic plan for computing at Fermilab, which I will describe today

🐝 **Fermilab**

# Charge

*We ask the committee to review the strategic plan for software and computing at the laboratory and the status of the recommendations made at the July 2020 PAC meeting:*
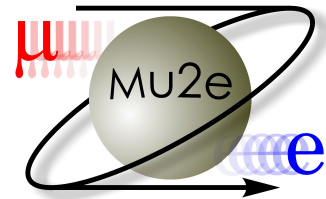
- "SCD should start thinking about the data lifecycle and its relation to long term data management and update the Committee at the next meeting. This report should include all relevant aspects, e.g., data storage, transfer, and compression".

- "The PAC continues to recommend [...] the development of a long-term computing transition plan" (to HPC)

Our responses to these two recommendations are integral parts of our overall strategic plan for computing at Fermilab, which I will describe today

🟦 **Fermilab**

# Computing Goal: Maximize 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
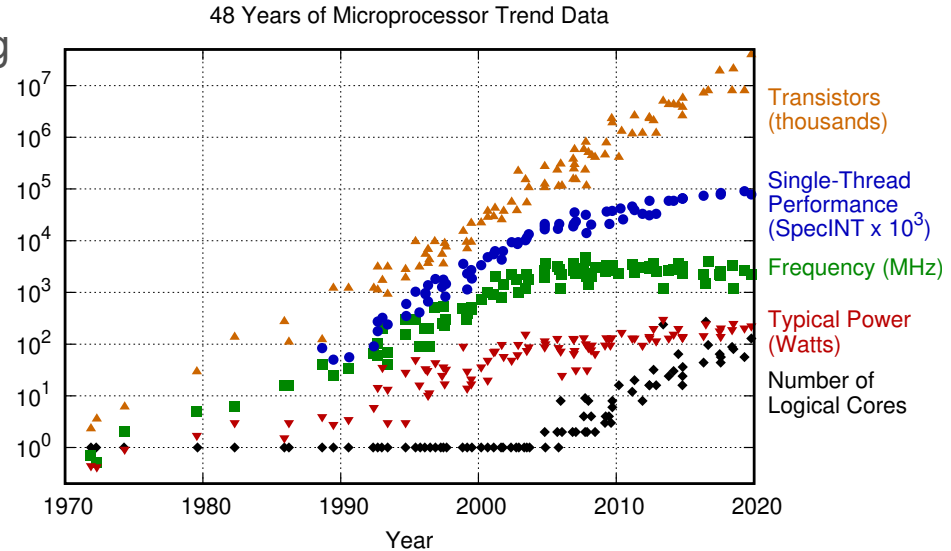
# Computing Strategy

## 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
  - Other academic/HEP resources
  - Commercial resources
- Embrace AI/ML developments
  - Enable scientific AI/ML applications
  - Utilize AI/ML for computing operations
    - AI/ML for accelerator operations out of the scope of this talk
- Adapt to evolving computing technologies while doing all the above

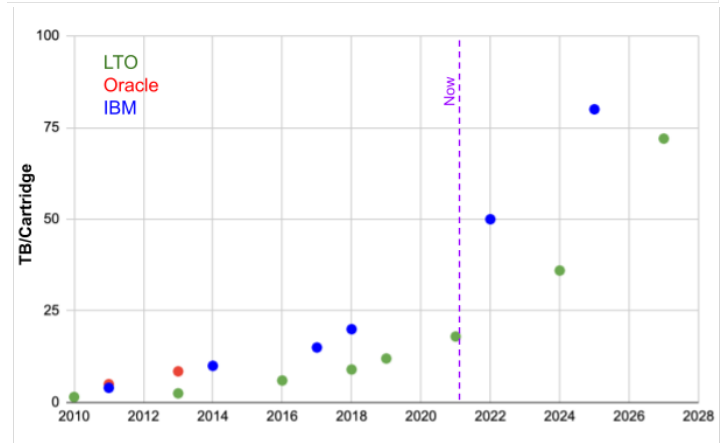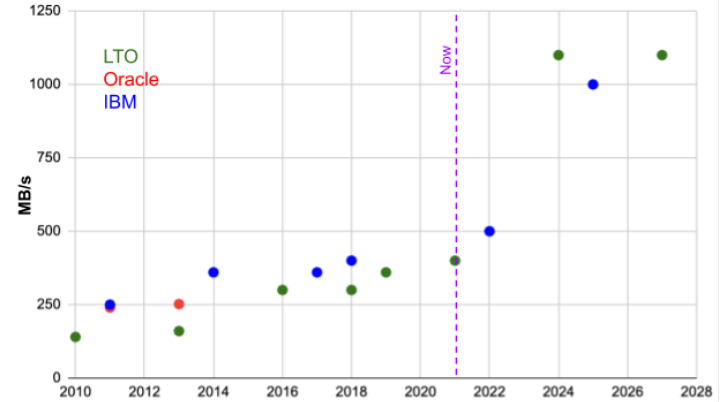🟦 **Fermilab**

# Evolving Computing Technologies

- CPU evolution radically different than before ~2010
  - Multicore now standard
    - Mostly accepted by community
  - GPUs now dominate numerical computing
    - Community acceptance slowly growing
    - 2014: Synergia2 100k turn simulation
      - 144 cores x 14 days
    - 2021: Synergia3 100k turn simulation
      - 1 GPU x 20 hours
  - Further developments on the horizon
    - Specialized AI hardware

### 48 Years of Microprocessor Trend Data

Transistors (thousands)

Single-Thread Performance (SpecINT x $10^3$)

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2019 by K. Rupp

🔷 Fermilab

# Evolving Storage Technologies: Tape

- Roadmaps show a large jump in drive rate past current generation
  - PAST predictions not reached, LTO8 predicted to be 500 MB/s, actual 360 MB/s
  - LTO9 @ 400 MB/s
  - Drive rates plateau or only increase slowly --> large number of drives required to meet HL-LHC rates (hundreds = multiple libraries)
- Cartridge capacity gives some indication of $/TB, cartridge cost holds around O($100)
  - LTO8 litigation debacle of 2018-9
  - LTO9 roadmap prediction (24 TB) not reached, actual is 18 TB.
    - LTO traditional "doubling/generation" prognostication continues – LTO10 @ 36 TB
  - Industry continues to fall behind TB/cartridge predictions
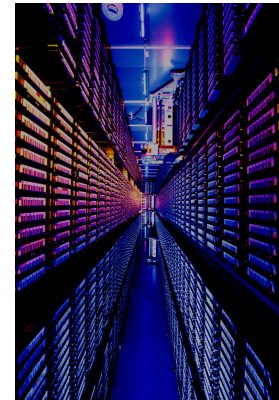- <span style="color:darkred">Throughput will be the new limiting factor in tape-based systems</span>



Tape Cartridge Capacity in TB vs year of introduction

🧲 Fermilab

# Fermilab Computing Strategy

## Detailed View

- Compute
  - Work towards getting a substantial fraction of US HEP computing from the Exascale/HPC machines
  - High-throughput computing at Fermilab will remain important
  - Utilize cloud resources when cost effective
    - Increase peak capacity
    - Access to non-standard hardware
- Storage
  - Mass storage will remain the foundation of Fermilab's computing capability
    - Tape libraries
    - Disk systems
    - Full-stack storage software
      - Including support for data lifetime management
- Analysis
  - Build an elastic analysis facility taking advantage of industry tools and Fermilab storage
- AI/ML
  - Provide GPU resources
  - Develop AI-enhanced operations
- Software
  - Pursue community-wide solutions
  - Engage ASCR partners
  - Collaborate with CERN
  - Leverage industry-standard tools where available
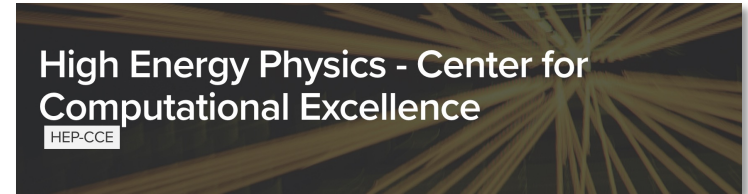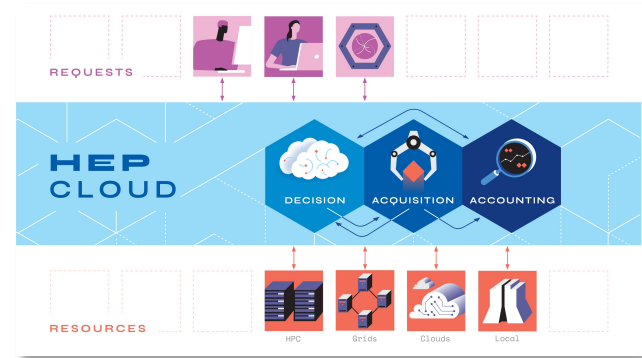
🔅 Fermilab

# DOE HPC/Exascale Resources

- NERSC
  - Current: Cori
    - Haswell, 2,388 nodes, 2.81 PFlops
    - KNL, 9,688 nodes, 29.5 PFlops
  - Next: 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)
  - Current: Summit
    - IBM Power9 + NVIDIA, 4,608 nodes, 200 PFlops
  - Next: Frontier (exascale!)
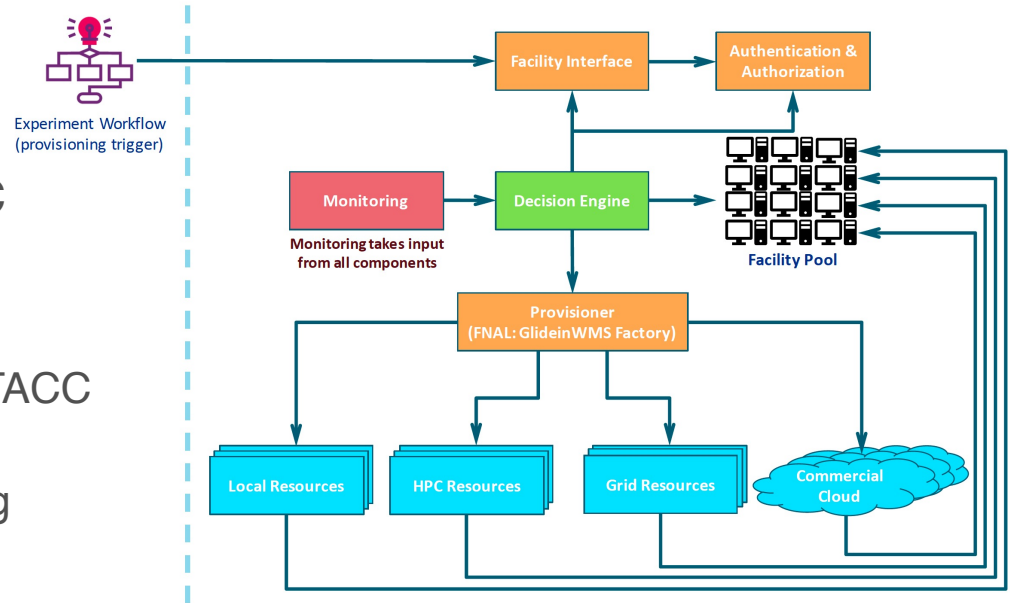    - AMD CPU + GPU, lots o' nodes, > 1,500 PFlops



**Fermilab**

# Barriers to Utilizing Exascale/HPC for HEP

- Job submission
  - Need ability to direct jobs to machines, authenticate, etc.
- Workflow management
  - HPC batch systems are designed for large, monolithic jobs
    - Grid-based systems are very different
- Data access
  - Data needs to flow in and out of HPC centers
  - LCF nodes do not have access to the open internet
  - HPC I/O systems were not designed with HEP in mind
    - HEP I/O was not designed with HPC in mind
- Compute architecture
  - Haswell: great!
  - KNL: OK, but needs some adaptation
  - NVIDIA GPUs: utilizable by a tiny fraction of HEP code, CUDA
  - Intel GPUs: new territory, not seen in wild, no CUDA
  - AMD GPUs: new territory, not seen in wild, HIP is like CUDA
- Allocations
  - In order to use the machines, experiments must have allocations



HEP CLOUD

REQUESTS

DECISION    ACQUISITION    ACCOUNTING

RESOURCES

HPC    Grids    Clouds    Local



**High Energy Physics - Center for Computational Excellence**

HEP-CCE

🔷 **Fermilab**

# 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 and TACC (NSF)
  - Commercial cloud providers including Google

- Work in progress on LCF facilities
  - Facility cooperation necessary to work around lack of network connectivity

🔷 Fermilab

# 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
- Related work ongoing in SciDAC4 projects
  - SciDAC5 projects will be complementary, but not overlapping
  - SciDAC includes both HEP and ASCR contributions

High Energy Physics - Center for Computational Excellence

HEP-CCE

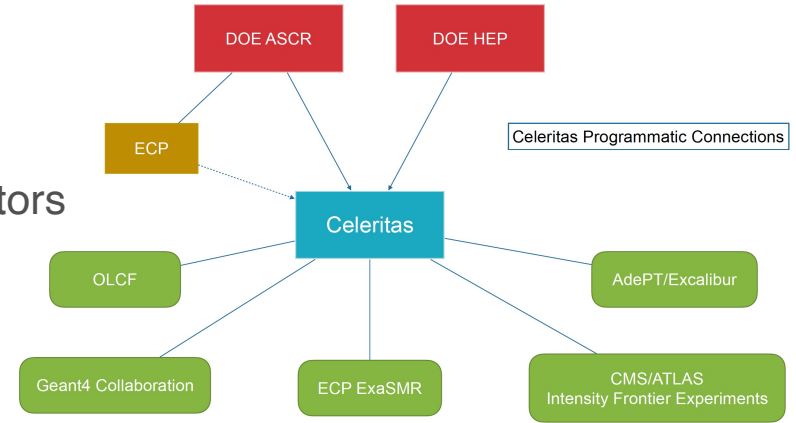https://www.anl.gov/hep-cce

🔷 Fermilab

# More HEP Software on GPUs

- Celeritas: detector simulation on GPUs
- Requirements
  - Utilize leadership class hardware (GPUs)
  - Read events from HEP community event generators (Pythia/HEPMC3, etc.)
  - Use community Geant4 geometry models (VecGeom/GDML)
  - Include (ultimately) complete physics models for HEP detector simulation in Geant4
    - Preliminary focus is on high-energy EM physics
  - Target most compute-intensive component of HEP detector simulation workflow: time-dependent, detector energy deposition (hit generation)
    - Complements and is part of standard Geant-driven LHC simulation workflow
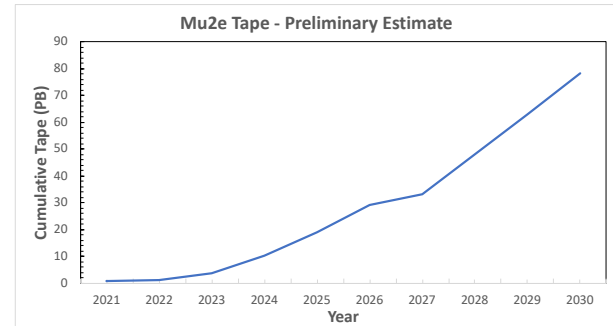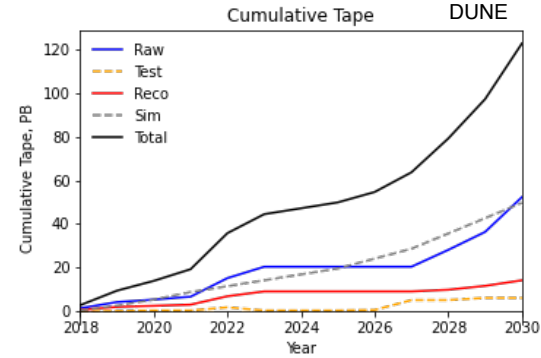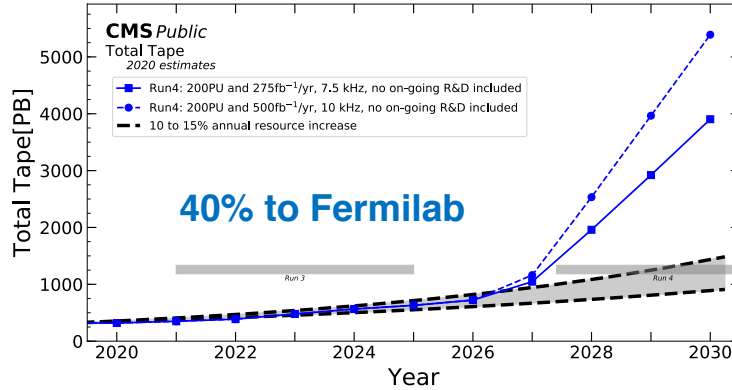
# DOE HPC Allocations

- NERSC
  - ERCAP
    - Request sent to OHEP manager
- ALCF, OLCF
  - Directors Allocation
    - Smallest
  - ASCR Leadership Computing Challenge (ALCC)
    - Less than half of allocated time
    - Competitive proposals
  - Innovative and Novel Computational Impact on Theory and Experiment (INCITE)
    - Principal source of allocated LCF time
    - Competitive, science-based proposals
    - Geared toward traditional HPC
      - (I have personal experience...)
  - CMS INCITE proposal went in from Fermilab this year
    - Expect this to be the start of a long process

🟛 **Fermilab**

# Fermilab Mass Storage Requirements

- Fermilab's data volume on tape today is 269 TB (225 TB active)
  - Two main categories
    - CMS
    - Public

- CMS during HL-LHC
  - ~1.6 EB by 2030
- DUNE
  - ~120 PB by 2030
- SBN
  - ~120 PB by 2030
- Mu2e
  - ~80 PB by 2030
- Small experiments
  - Lacking detailed plans, but small co
  - Adds to support load and complexit
    - complexity grows *more slowly*
    - complexity grows *faster than li*
- Legacy experiments
  - The support load of these experime



**40% to Fermilab**



Mu2e Tape - Preliminary Estimate

2/8/21    Jayatilaka I Paradigm Evolution

2/8/21    Jayatilaka I Paradigm Evolution

# Building our Next-Generation Mass Storage System

- Current system: tape + disk
  - Tape: Enstore
    - Fermilab product
  - Disk: dCache
    - Collaboration with DESY and Nordic e-Infrastructure collaboration, NeIC
  - Built for Tevatron Run II
  - Designed for Petabytes
- Next generation: tape + multiple disk types + data management systems
  - Need an integrated system
  - Tape is terrible, but no realistic alternative
  - Designing for Exabytes

🔷 **Fermilab**

# Storage Research and Development

- Starting work on 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
  - EOS
    - Collaborate with CERN
  - ceph
    - Broad usage in multiple industries
- Emphasizing Rucio within software layer
  - Broad community support
  - Provides mechanism to enforce data lifetimes
    - Experience shows that manual lifetime management is not realistic











🔷 **Fermilab**

# Ramping Up Storage Research and Development

- Storage research and development is the major focus of this year's International Computing Advisory Committee (ICAC)
  - February 2021 meeting focused on storage R&D plans
    - https://indico.fnal.gov/event/47365/
  - Winter 2022 meeting to focus on implementation



- Storage developers are Scientific Computing Division's highest hiring priority
  - One new hire has arrived
  - More in the pipeline

- Continuing with storage-related collaborative activities
  - WLCG DOMA, Rucio, dCache, IRIS-HEP (analysis)

🟁 **Fermilab**

# Elastic Analysis Facility



- Part of Fermilab computing strategy
- Take advantage of industry big data tools
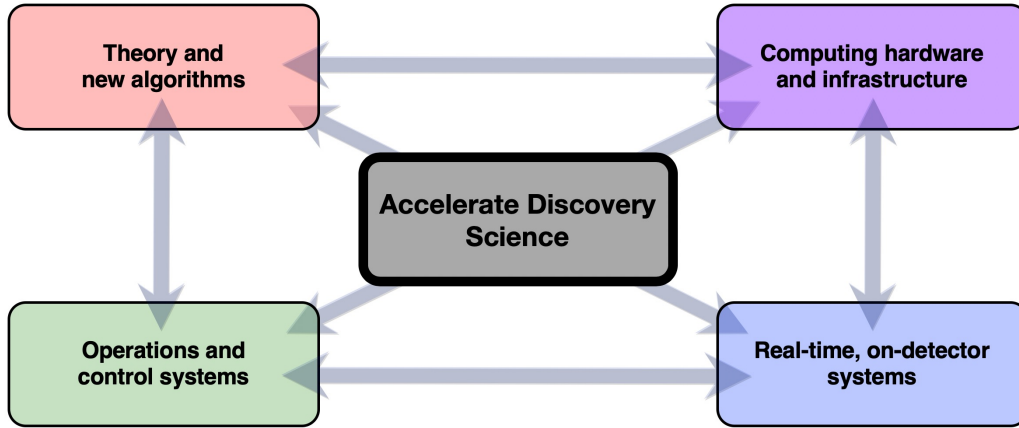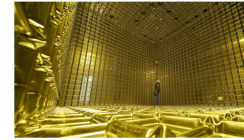- Fast data access is a key ingredient
- Goal is to minimize time to scientific insight

🐾 Fermilab

# AI/ML at Fermilab

```
┌──────────────────┐                           ┌──────────────────┐
│  Theory and      │                           │ Computing hardware│
│  new algorithms  │                           │ and infrastructure│
└──────────────────┘     ┌──────────────────┐  └──────────────────┘
                         │ Accelerate Discovery│
                         │    Science         │
┌──────────────────┐     └──────────────────┘  ┌──────────────────┐
│  Operations and  │                           │ Real-time, on-detector│
│  control systems │                           │    systems        │
└──────────────────┘                           └──────────────────┘
```

- AI/ML strategy at Fermilab extends across divisions
- Computing items
  - Purchasing 12 A100 GPUs using FY21 funds
    - In progress
  - Effort to use AI/ML to optimize computing operations just getting started
    - Seeking funding
  - "Expect" to get support to hire dedicated AI/ML professional in computing



AI / DEEP LEARNING | HPC                         Apr 30, 2021

Scaling Inference in High Energy Particle Physics at Fermilab Using NVIDIA Triton Inference Server

By Shankar Chandrasekaran, Lindsey Gray, Farah Hariri, Kevin Pedro, Vartika Singh, Nhan Tran, Mike Wang and Tingjun Yang

Tags: featured, kubernetes, NGC, physics, Triton

## ML Commons
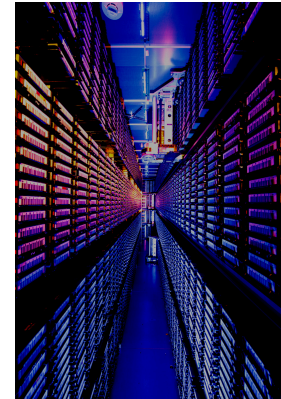
## Recent newsworthy results
- Inference results highlighted by NVIDIA
- MLCommons
  - AI/ML public benchmarking site
  - Fermilab TinyML submission
    - low-power, edge devices
  - Press release in June

🟦 Fermilab

# Fermilab Computing Strategy

## Detailed View

- Compute
  - Work towards getting a substantial fraction of US HEP computing from the Exascale/HPC machines
  - High-throughput computing at Fermilab will remain important
  - Utilize cloud resources when cost effective
    - Increase peak capacity
    - Access to non-standard hardware
- Storage
  - Mass storage will remain the foundation of Fermilab's computing capability
    - Tape libraries
    - Disk systems
    - Full-stack storage software
      - Including support for data lifetime management
- Analysis
  - Build an elastic analysis facility taking advantage of industry tools and Fermilab storage
- AI/ML
  - Provide GPU resources
  - Develop AI-enhanced operations
- Software
  - Pursue community-wide solutions
  - Engage ASCR partners
  - Collaborate with CERN
  - Leverage industry-standard tools where available

**Fermilab**

# 2. Miscellany

‡ Fermilab

# Hiring

- We have several openings and there are a few more in the works
  - Not computing professionals (mostly external and project funding)
    - AI Research Associates
    - AI Associate/AI developer with FPGA familiarity
    - Real-time FPGA Developer/Staff Engineer
    - Guest engineer
    - Cosmic postdocs
  - Computing professionals (mostly base and CompHEP funding)
    - Multiple positions
      - Ramping up storage research and development
      - Replacing retirements

🔷 **Fermilab**

# IPv6

- The lab must adhere to an OMB mandate requiring a PLAN by end of FY21 to accomplish for systems and applications/services:
    - By FY23 all new systems IPv6 enabled ("system" can be a group of devices)
    - By end FY23 20% of assets IPv6-only
    - By end FY24 50% of assets IPv6-only
    - By end FY25 80% of assets IPv6-only
    - Identify, justify, and schedule to replace/retire any that cannot be converted
- In response to a data call in 5/2021 the lab reported that only 21% of assets could be made to be compliant in this timeframe
- A new data call issued 10/20/21, due 12/1/2021 requires updates with more details
    - Various SCD people are helping with this now – please cooperate!
- Lots of questions remain:

  Interaction with IPv4-only resources?    "Science" exemptions?
  Dual-stack devices allowed?              Funding for replacements?

🎄 Fermilab

# Travel (really "Travel")

- Virtual "travel"
  - Domestic events without associated fees to not require travel paperwork
  - Domestic events with fees do require travel paperwork
  - Foreign events with or without fees require pre-approval and subsequent paperwork
    - Lead times are long! Submit your request as soon as possible.
- Physical travel
  - In person attendance at foreign or domestic conferences/workshops/seminars is not currently allowed
  - Mission critical travel for other work is possible

- Reimbursements
  - Auditing of travel-related expenses has become much more stringent than it was in the past
    - Funds committed before approval are not allowable
    - Rule of thumb: do not spend your own money and expect to be reimbursed
      - Applies outside of travel!

🔶 **Fermilab**

# 3. Questions

🎇 **Fermilab**

## (Anonymous)

There was word recently that DOE has suspended pretty much all in-person travel to conferences, workshops, seminars, etc. Have you heard whether that also applies to experiments' collaboration meetings?

🔷 **Fermilab**

# From Art Kreymer

I have had one escorted visit to my Fermilab office since March 2020,to turn off the desktop computer.

Would it be appropriate for people working at home to schedule a monthly office visit, unescorted but vaccinated, masked, and distanced, for the purpose of moving materials between home and office?

🔷 **Fermilab**

**(Anonymous)**

How are we accounting for time lost sitting at the gate when selected for a randomized test (or if stuck behind someone who was)?

🐝 **Fermilab**

# More Questions?

# The floor is open

Fermilab