# A Word on DUNE Computing

Ken Herner, Fermilab

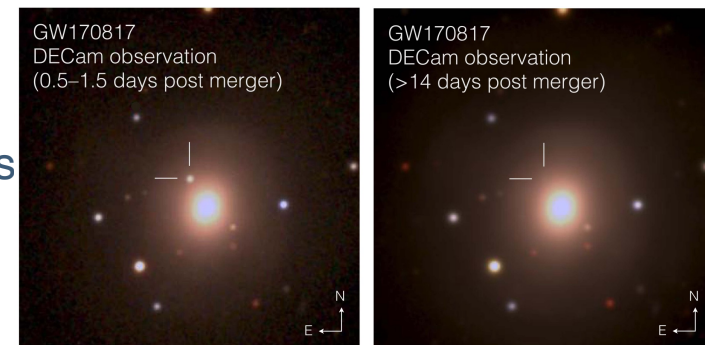Computational HEP Traineeship Summer School

24 May 2024

# Outline

- Brief Introduction to DUNE and its science goals

- Set the scale

- DUNE's uniqueness (and unique challenges)

- A Big and Unique problem: supernova burst processing

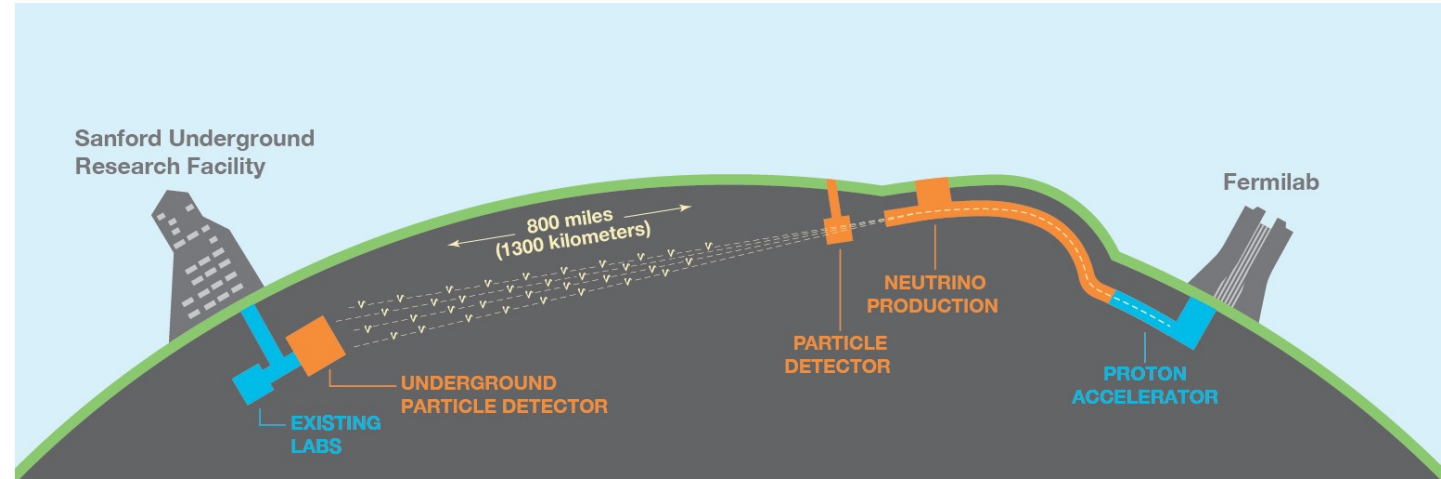2024-05-24     K. Herner | DUNE Computing

# Introduction to Ken

- Ph.D. on the D0 (Tevatron) experiment; most of postdoc as well (centered on Higgs boson searches)

- At Fermilab in Scientific Computing Division (now Computational Science and AI Directorate) since late 2012

  - Started with long-term preservation of D0 dataset

  - Distributed computing support for nearly all lab experiments

  - Fabrlc for Frontier Experiments Project lead 2018-2021 (modular SW toolkit for non-CMS experiments)

- DUNE

  - Production Group coordinator 2018-2022: ProtoDUNE data processing, MC generation, new compute site integration, infrastructure R&D (e.g. GPUaaS); Since October 2023: Host Lab Technical Lead and DUNE-US software and computing lead

- Some other projects (not exhaustive list)

  - Various projects on Dark Energy Survey, including optical followup of GW events

    - Currently co-convener of Transient and Moving Objects Science Working Group

  - Vera Rubin Observatory

    - Work in Data management on alert production systems and campaign management



GW170817
DECam observation
(0.5–1.5 days post merger)

GW170817
DECam observation
(>14 days post merger)

[Soares-Santos et al., ApJL 848:L16 (2017)](#)
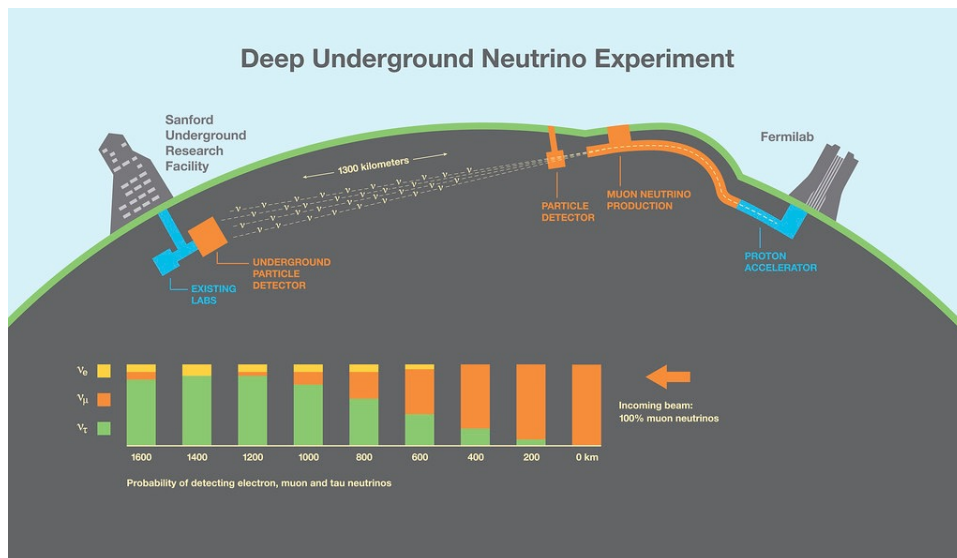
**Fermilab**  DUNE

# The Deep Underground Neutrino Experiment

- Future flagship experiment of Fermilab; primary focus is extending our understanding of the neutrino sector

- Consists of a Near Detector Complex (FNAL) and a Far Detector Complex (SURF in SD)
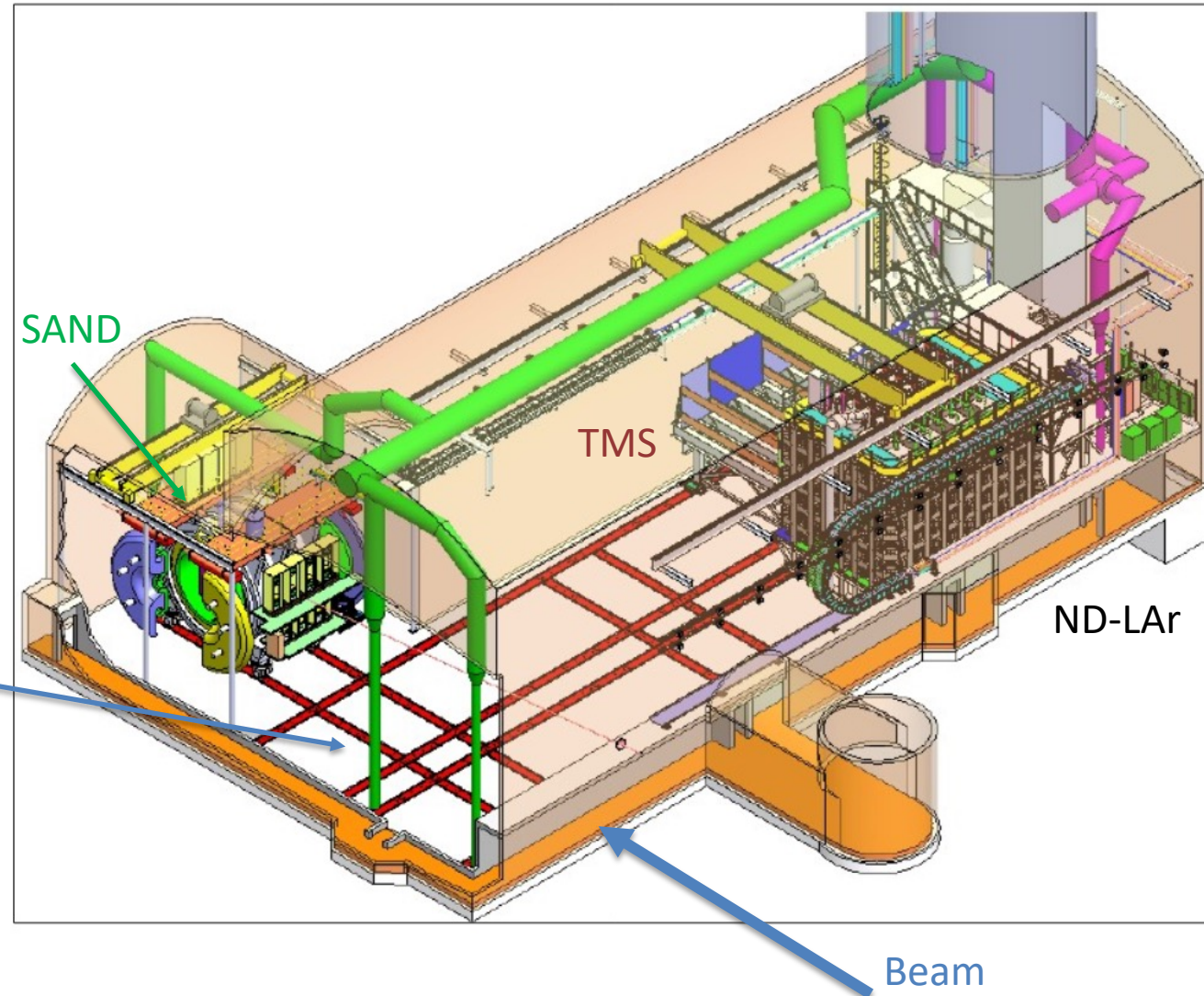
# Science Goals

- Three main physics goals:

  - Origin of matter (neutrino oscillation parameters, mass ordering, ...)

  - Unification of forces (BSM physics, e.g. proton decay)

  - Black Hole Formation (neutrinos from core-collapse SN)

- Learn more at https://www.dunescience.org



Deep Underground Neutrino Experiment

2024-05-24    K. Herner I DUNE Computing
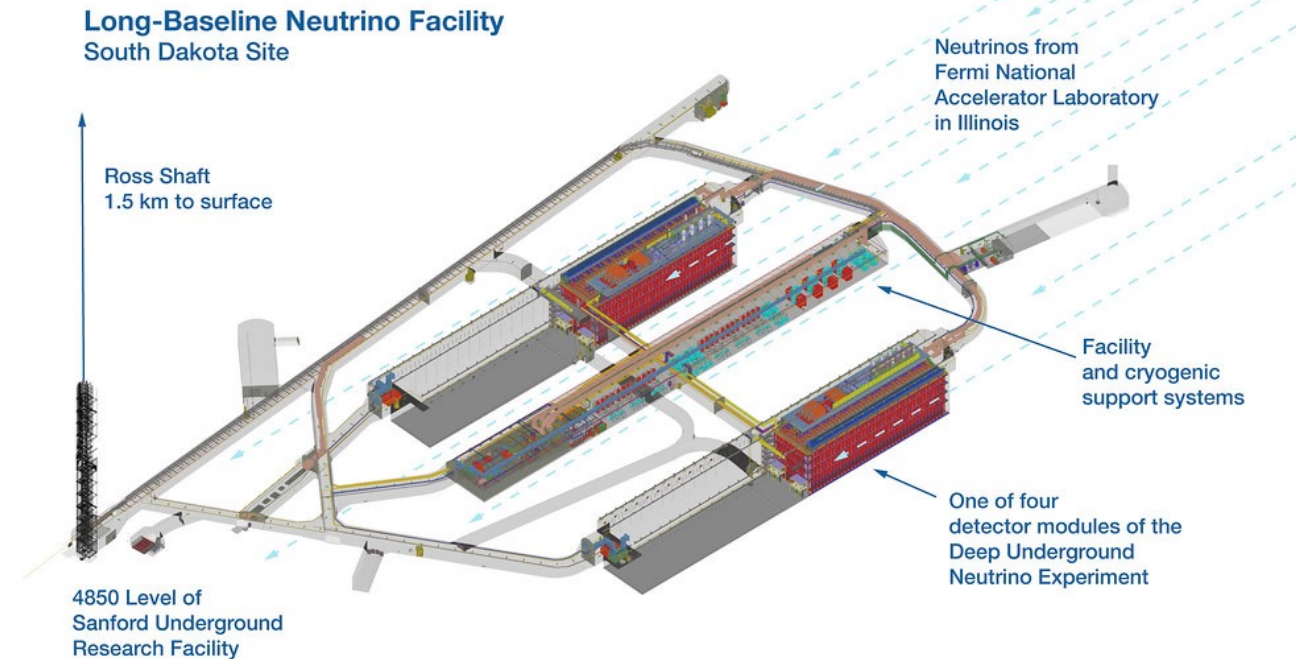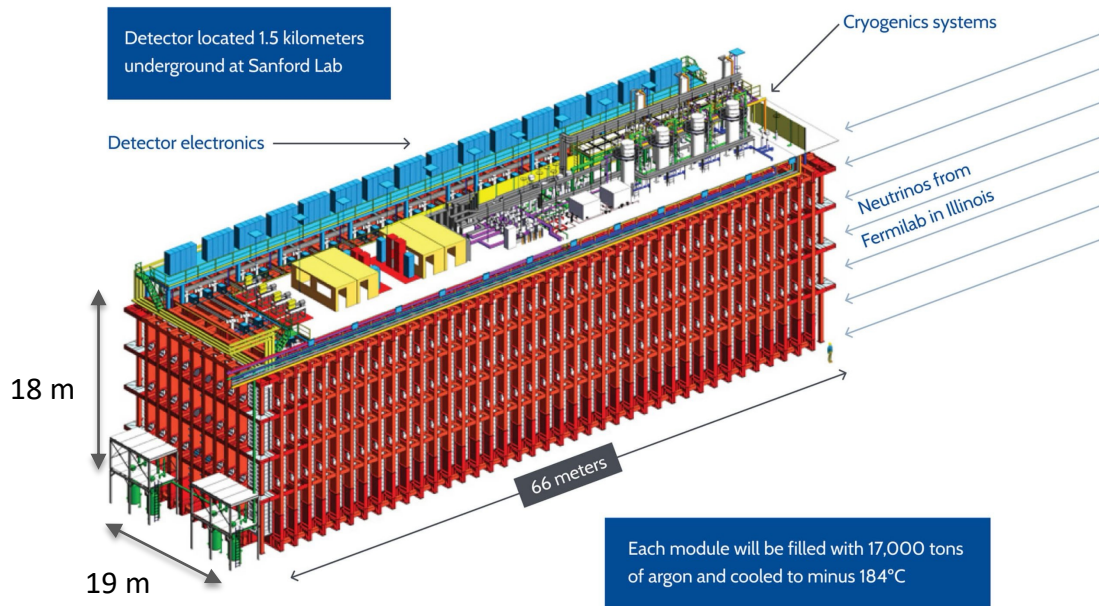
# Near Detector Complex

- Really, three detectors:

  - ND-LAr: Liquid argon TPC

  - TMS: Temporary Muon Spectrometer

  - SAND: System for on-Axis Neutrino Detection

- PRISM system allows ND-LAr and TMS to move off-axis to perform measurements at different angles to beam (SAND fixed).

- TMS to be replaced



SAND

TMS

ND-LAr

Beam

**Fermilab** DUNE

# Far Detectors

- Far Detector will eventually consist of four LAr-based* TPC modules

- Installed in two stages (FD1 and FD2 for DUNE Phase 1)

\* fourth module is still officially the "module of opportunity"

# Prototype Detectors

- ProtoDUNE horizontal and vertical drift detectors (Far detector technology; CERN-Prevessin)

- NDLAr **2x2** prototype (Near detector technology; FNAL)



[Link to image](#)

# Rough DUNE Timeline and Goals

- 2024
  - ProtoDUNE II beam run @ CERN; ND 2x2 prototype run @ FNAL
    - Aside: DAQ outputs HDF5 files (new to neutrino experiments)

- ca. 2025-2028
  - Phase 1 far detector installation

- ca. 2029-2031
  - ND installation, FD commissioning, first science runs with all Phase 1 detectors (ND and FD1/2)

- Late 2030s
  - Phase 2: full far detector, upgraded near detector, upgraded neutrino beam (all together: 4x Phase 1 event rate)

| Experiment Stage | Physics Milestone | Exposure (kt-MW-years) | Years (Staged) |
|---|---|---|---|
| Phase I | $5\sigma$ MO ($\delta_{CP} = -\pi/2$) | 16 | 1-2 |
| | $5\sigma$ MO (100% of $\delta_{CP}$ values) | 66 | 3-5 |
| | $3\sigma$ CPV ($\delta_{CP} = -\pi/2$) | 100 | 4-6 |
| Phase II | $5\sigma$ CPV ($\delta_{CP} = -\pi/2$) | 334 | 7-8 |
| | $\delta_{CP}$ resolution of 10 degrees ($\delta_{CP} = 0$) | 400 | 8-9 |
| | $5\sigma$ CPV (50% of $\delta_{CP}$ values) | 646 | 11 |
| | $3\sigma$ CPV (75% of $\delta_{CP}$ values) | 936 | 14 |
| | $\sin^2 2\theta_{13}$ resolution of 0.004 | 1079 | 16 |

Snowmass 2021

5 sigma needs a LOT of pseudoexperiments....

🟣 Fermilab   DUNE

# Solving the Computing Problem

- Computing CDR written in 2022

- https://arxiv.org/abs/2210.15665

- Covers offline model, use cases, resource estimates

- Must eventually integrate information from seven distinct detectors

- Of course, things evolve even on the timescale of two years

- **Every single item you've discussed this week is relevant to DUNE in some way, shape, or form (many are absolutely critical to its success)**
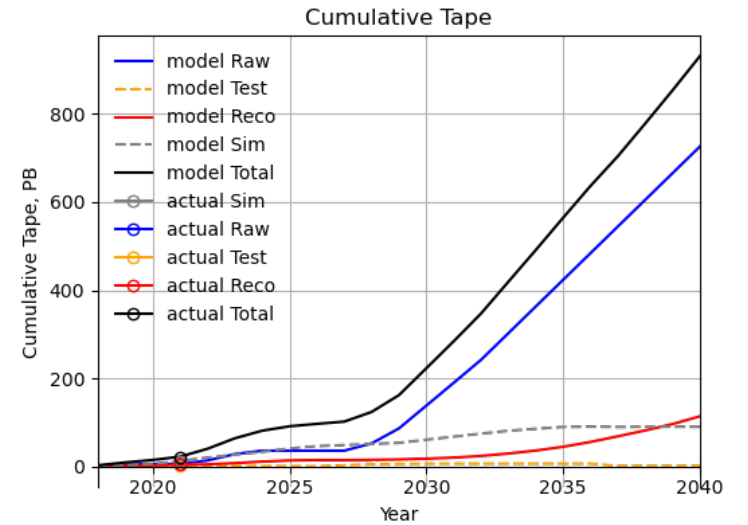
DUNE Offline Computing
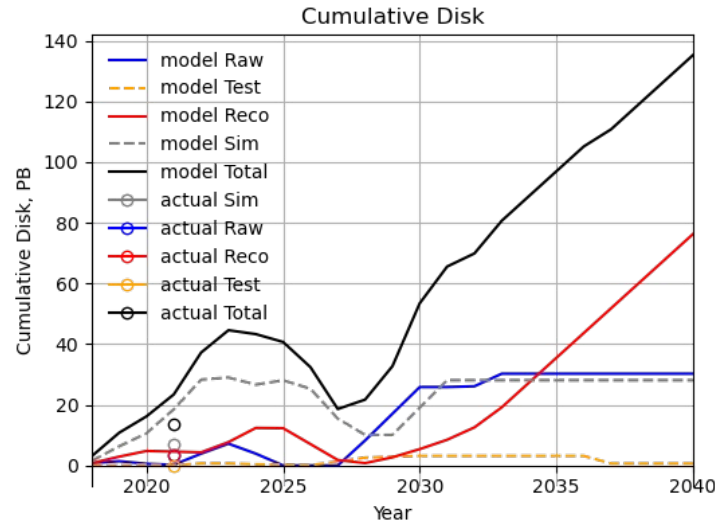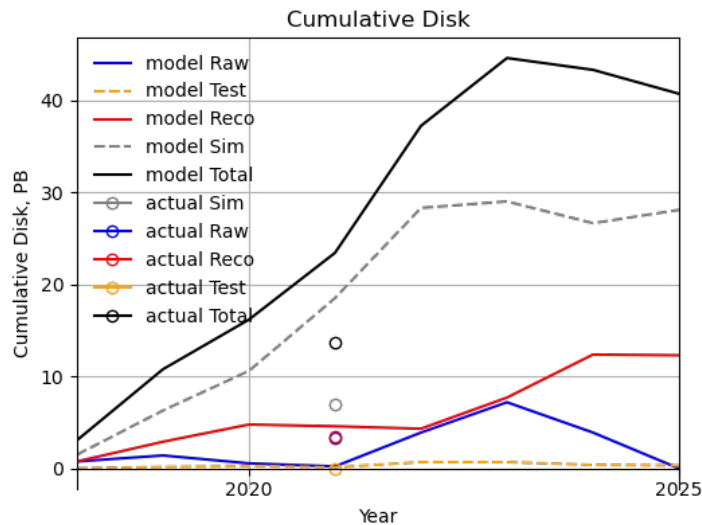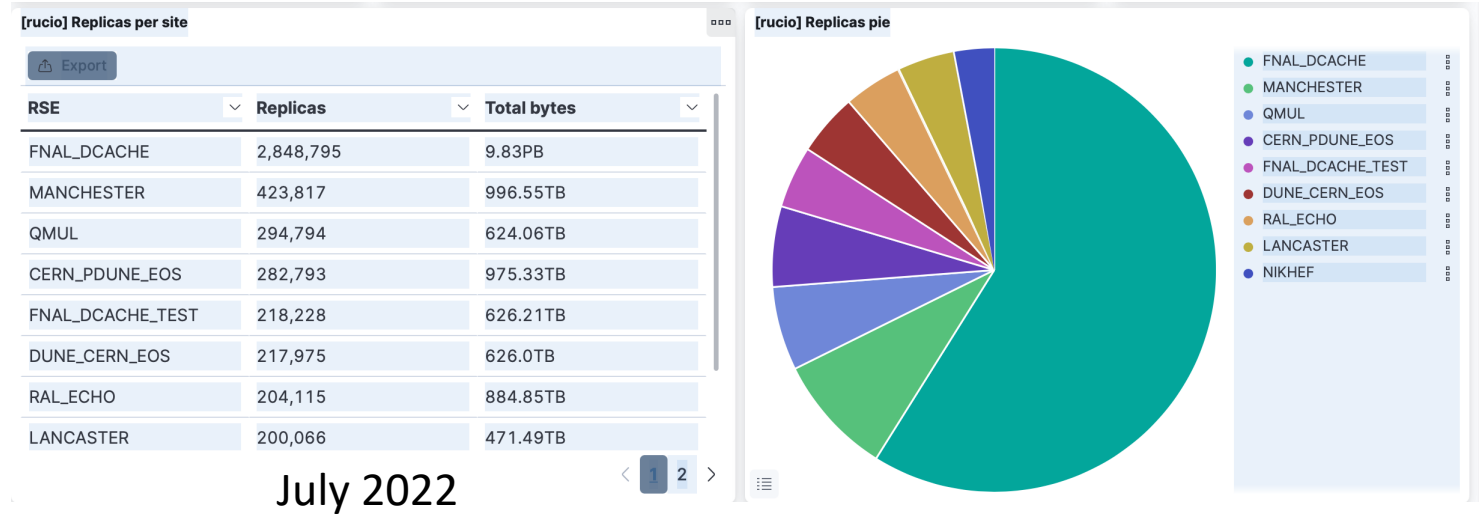
Conceptual Design Report



October 31, 2022

The DUNE Collaboration

arXiv:2210.15665v1 [physics.data-an] 28 Oct 2022

# Storage Status and Requirements

- Combination of disk cache and tape for archiving

- Roughly 10% of LHC experiments in 2030s (also true of CPU needs)

- 2 geographically separate copies of raw detector data

| [rucio] Replicas per site | | |
|---|---|---|
| ⬆ Export | | |
| RSE | Replicas | Total bytes |
| FNAL_DCACHE | 2,848,795 | 9.83PB |
| MANCHESTER | 423,817 | 996.55TB |
| QMUL | 294,794 | 624.06TB |
| CERN_PDUNE_EOS | 282,793 | 975.33TB |
| FNAL_DCACHE_TEST | 218,228 | 626.21TB |
| DUNE_CERN_EOS | 217,975 | 626.0TB |
| RAL_ECHO | 204,115 | 884.85TB |
| LANCASTER | 200,066 | 471.49TB |

July 2022

[rucio] Replicas pie
- FNAL_DCACHE
- MANCHESTER
- QMUL
- CERN_PDUNE_EOS
- FNAL_DCACHE_TEST
- DUNE_CERN_EOS
- RAL_ECHO
- LANCASTER
- NIKHEF

# Distributed Computing Tools in DUNE

- Mostly High Throughput Computing-style (many single-core) jobs so far, similar to Large Hadron Collider experiments

  - However, more of a service-based model than a tiered model. Sites provide one of more services (standard compute, High Performance Computing, storage/compute, archiving, user analysis, etc.)

- Use variety of sites; mix of dedicated/pledged resources and opportunistic access through OSG

- Mostly stream input data over network with XRootD

- Rucio for dataset transfer and replication

- CVMFS for core software distribution; jobs run in containers (image also distributed via CVMFS)

- New software developed for bookkeeping and workflow management, known as justIN

- Good network design and performance are critical to success, especially in a "flatter" architecture
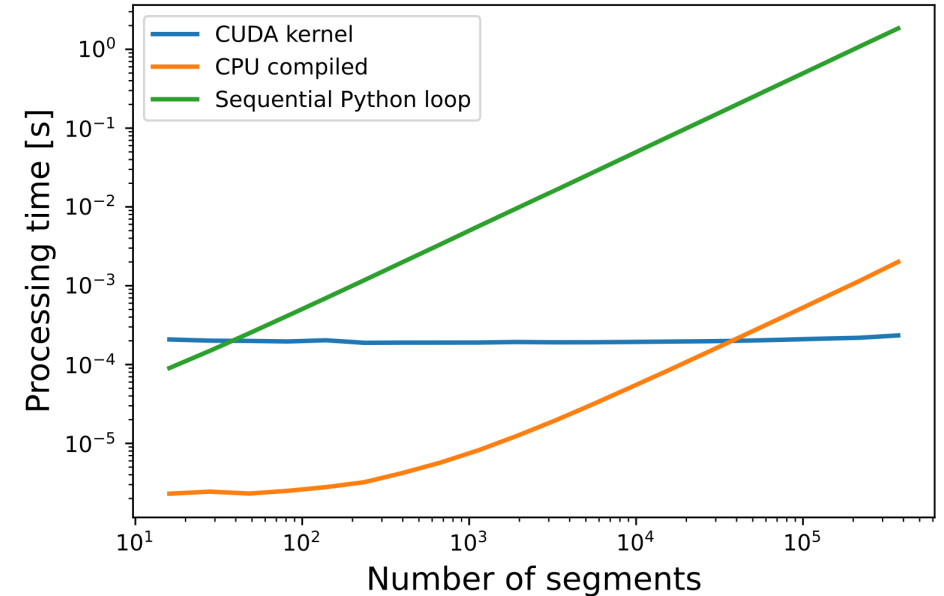
# Some DUNE-ique items

- There are some differences wrt LHC experiments

  - "Events" don't really mean the same thing ("Trigger records")

  - Trigger records are *much* bigger than LHC events: O(100 MB) vs. O(1 MB)

  - Current and future DAQs creating HDF5 files (before, essentially root files).

- Earlier: "Mostly High Throughput Computing-style (many single-core) jobs so far..."

- HPC becoming increasingly important, especially with GPU availability

  - Also critical for performing final oscillation parameter fits and systematic variations (mentioned a bit earlier); phase space is large and requires MPI processing; exception to above rule of thumb (but a vital one!)

- Speaking of GPUs...

**Fermilab** DUNE

# GPU needs example: ND-LAr

- ND needs to be able to make measurements close to the neutrino production point and contain a large fraction of the (hadronic side) of the
- events to fulfill basic requirements like constraining the flux*cross
- section model in a way that can be transferred to the FD.
- Design (~5x5mm pixels, segmentation) driven by energy resolution requirements (see ND conceptual design report) and need to disentangle the ~50 events per beam spill.
  - Leads to roughly 12 million channels.
- Large number of independent channels lends itself very well to parallelization, which is why GPUs were adopted pretty early on (see larnd-sim paper; shown on this slide)
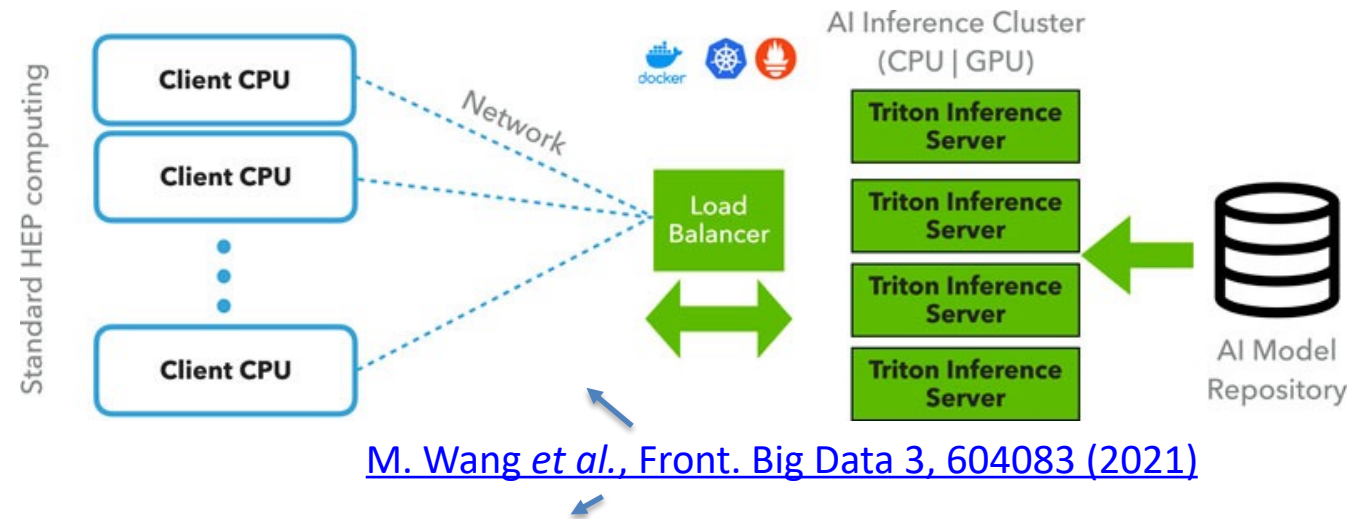- **Also significant use of ML in reconstruction**



| Calculation | Loop over | Quantity | GPU speed-up factor | |
|---|---|---|---|---|
| Recombination factor | Segments | $10^5$ | shown in plot | ×3 |
| Induced current | Pixels | $10^3$ | most expensive task | ×7314 |
| Charge electronics response | Pixels | $10^3$ | | ×985 |
| Light time profile | Time ticks | $10^5$ | | ×228 |
| Scintillation profile | Time ticks | $10^3$ | | ×568 |
| Light electronics response | Time ticks | $10^3$ | | ×1883 |

DUNE Collaboration, JINST 18 P04034 (2023)

🔷 Fermilab  DUNE
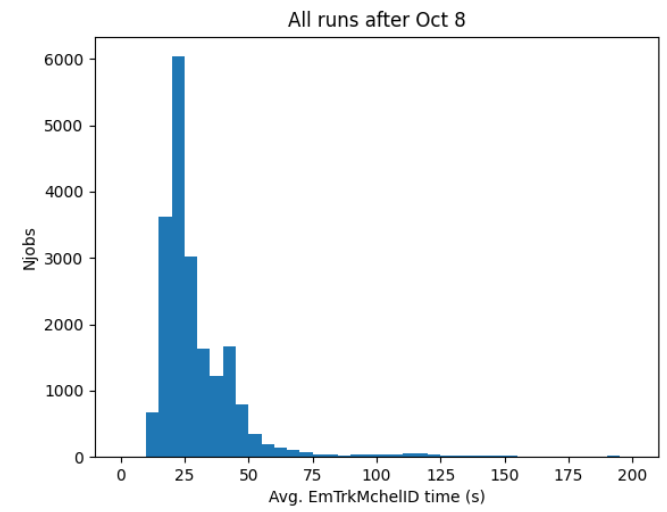
# GPUaaS Studies



- Tests with MC show large decrease in processing time for part of track reconstruction if run on a GPU

  - Use Triton inference server and gRPC in job for communication, allows many-to-1 model of CPU jobs to single GPU (the SONIC Model)

  - Rest of stack uses standard CPU, can thus run on any site with ext. network without needing local GPU

- ProtoDUNE beam data reprocessing campaign in 2021 utilized approach at scale (few thousand concurrent jobs), used cloud-based GPU server. See Cai et al. 2023 for details

  - Clear overall speed increase wrt CPU-only version, but overall amount of data movement per job increased *significantly* (10x) wrt CPU-only versions

  - Must take care not to saturate network capacity or time gains will be lost

  - Have also run with local servers on NERSC Perlmutter nodes

M. Wang *et al.*, Front. Big Data 3, 604083 (2021)

|  | ML module | non-ML modules | Total |
|---|---|---|---|
|  | Wall time (s) | | |
| CPU only | 220 | 110 | 330 |
| CPU + GPUaaS | 13 | 110 | 123 |



Further reading on SONIC model:
Duarte et al. (2019), Krupa et al. (2021), Pedro et al. (2019)
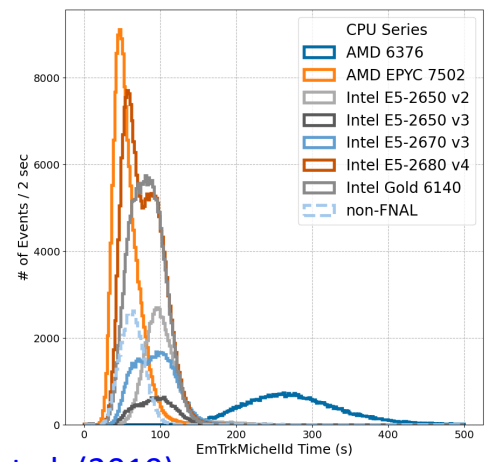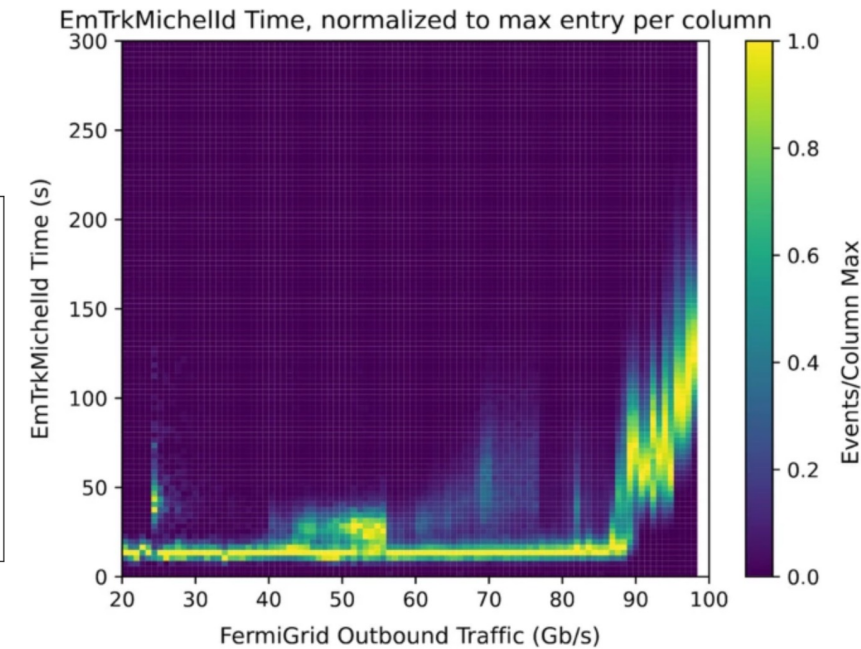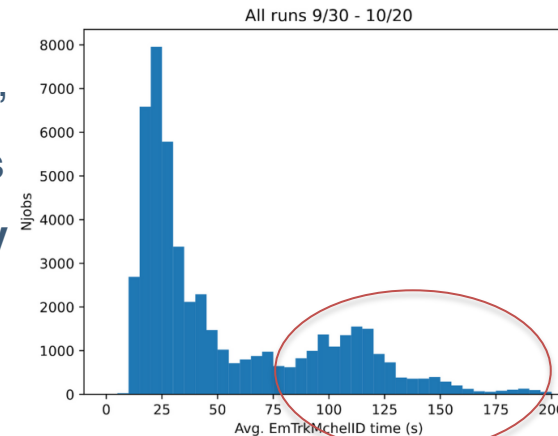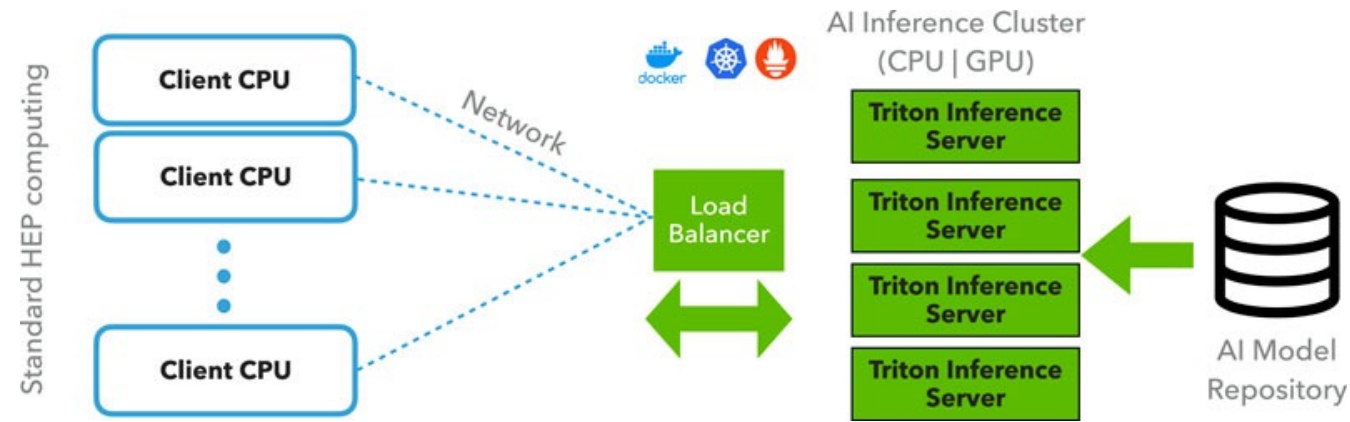
🟦 Fermilab DUNE

# GPUaaS Studies



- Tests with MC show large decrease in processing time for part of track reconstruction if run on a GPU

  – Use Triton inference server and gRPC in job for communication, allows many-to-1 model of CPU jobs to single GPU (the SONIC Model)

  – Rest of stack uses standard CPU, can thus run on any site with ext. network without needing local GPU

- ProtoDUNE beam data reprocessing campaign in 2021 utilized approach at scale (few thousand concurrent jobs), used cloud-based GPU server. See Cai et al. 2023 for details

  – Clear overall speed increase wrt CPU-only version, but overall amount of data movement per job increased *significantly* (10x) wrt CPU-only versions

  – **Must take care not to saturate network capacity or time gains will be lost**

  – Have also run with local servers on NERSC Perlmutter nodes


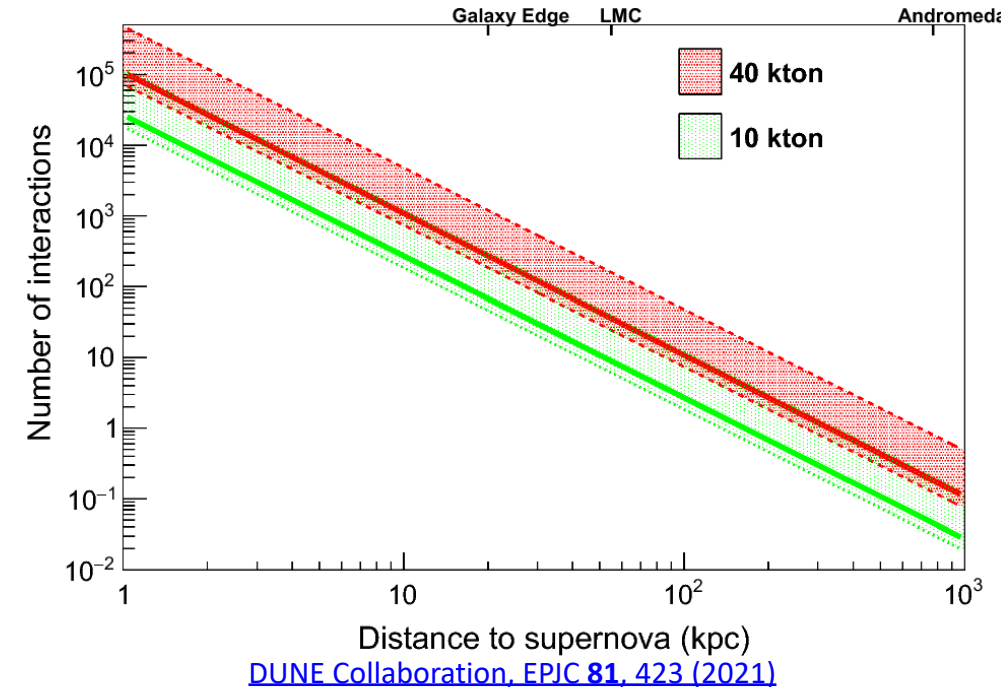


Further reading on SONIC model:
Duarte et al. (2019), Krupa et al. (2021), Pedro et al. (2019)

# Core Software Framework

- ProtoDUNE and FD software based on the Art framework (descendant of CMSSW). Remember, however, that

    – DUNE and collider detector events differ in size and structure; readout not always synced to a beam timing

    – DUNE needs easy access to heterogenous computing resources

    – Some external tools and packages lend themselves more easily to Python than C++

- Therefore, DUNE has commissioned the design of a new core software framework to support both ND and FD workflows. Key features include:

    – Support for algorithms/modules in multiple languages (C++ and Python for sure)

    – Flexible, user-defined, data groupings (must sync up different detector data in differing ways)

    – More "eager writing" (don't wait until the end to write entire event to disk)

    – Configurable, composable, framework-agnostic algorithms will "plug in" to the core framework which will schedule data/workflow/hardware access (e.g. GPU) and handle data provenance tracking and other metadata, and data persistency and I/O.

- Planning to deliver during FY27 (FY starts October 1, 2026)

🎄 **Fermilab** DU(VE

# The Big Challenge: a supernova

- Sensitive to neutrinos from a (relatively) nearby supernova

- Would read detector out in a continuous mode for ~ 100s

- Expect approximately 400 TB (uncompressed) for four-module readout; **still 184 TB at current compression levels**

- Prelim studies of event reconstruction could lead to a source location accurate to < 5 degrees– **important for follow-up by other instruments (optical/near IR telescopes, etc.)**

  - **But time is everything in this game. 4+ hours to transfer data on a dedicated 100 Gb network, plus processing time…**

- Expect one SN trigger per month (mostly false alarms)

DUNE Collaboration, EPJC **81**, 423 (2021)

| Supernova: | | |
|---|---|---|
| SINGLE CHANNEL READOUT | 300 MB | Uncompressed 100 s |
| FOUR MODULE READOUT | 460 TB | Uncompressed 100 s (assumption) |
| TRIGGER RATE | 1 per month | |

**Fermilab**  DUNE

# Processing supernova data

- Limited space and infrastructure (i.e. cooling) at the far site means bulk processing on a local farm is impractical to impossible

- Don't have to transfer full dataset before processing starts

  - Workflow mgmt systems must deal with dynamic datasets and shifting resource availability

- 10,000 – 40,000 present-day CPUs needed for reconstruction (30k likely) to finish within a few hours of event (goal: preliminary direction before event rises in optical bands)

  - HPC centers can probably fit the bill, but data transfer in and out is the major problem to solve

  - **Must consider redundancy as well**

- Must be able to handle large input stream as well as output at a similar rate

  - Run standard data reco or make a slimmed-down, faster version? Speed vs. accuracy tradeoff?

- Implications for additional network paths? What are those requirements? Costs?

🟁 **Fermilab**  DUⁿE

# Summary

- DUNE will significantly advance our understanding of one of the least understood elementary particle types

- Potential for significant paradigm-shifting discoveries (e.g. proton decay)

- A supernova in our galaxy will lead to a wealth of information, but processing detector data in a timely fashion will be one of DUNE's biggest challenges

- DUNE will not succeed without efficient, robust computing! Affects all aspects of the experiment

- General strategy: use common tools wherever possible, share information with other experiments, but don't be afraid to break new ground where needed.

- Plenty of room to get involved and work on interesting projects!

Thank You!

🔬 **Fermilab** DUNE

# Backup

     K. Herner | DUNE Computing

🔷 Fermilab  DUNE

# Some further reading

- [dunescience.org](dunescience.org)

- [https://arxiv.org/abs/2210.15665](https://arxiv.org/abs/2210.15665)

- [ML reconstruction for LAr TPCs](ML reconstruction for LAr TPCs)

- [DUNE Physics Snowmass Whitepaper](DUNE Physics Snowmass Whitepaper)

- [ProtoDUNE performance paper](ProtoDUNE performance paper) (Figure 47 is a good example of how to distinguish between the different types of events)

# DUNE Computing is global

- Significant efforts to add international presence began in 2018

- Compute sites in 12 countries + CERN on four continents

- Record time for adding DUNE support to a grid site (start to successful batch jobs) is **2h**.

- For October 2021-October 2022, > 50% of production wall hours outside USA; < 25% at Fermilab

- Numerous sites have also pledged storage support (19 PB in 2022)



Legend:
- US
- UK
- CERN
- NL
- ES
- CZ
- FR
- RU
- BR
- CA
- IN

US_FermiGrid, US_NERSC_Cori, US_BNL, US_Colorado, US_Wisconsin, US_WSU, US_UCSD, US_NotreDame, US_PuertoRico, US_NERSC, US_MWT2, US_Michigan, US_SU-ITS, UK_RAL-Tier1, UK_RAL-PPD, UK_Bristol, UK_Lancaster, UK_QMUL, UK_Manchester, UK_Oxford, UK_Imperial, UK_Brunel, UK_Sheffield, UK_Liverpool, UK_Edinburgh, NL_SURFsara, NL_NIKHEF, ES_PIC, ES_CIEMAT, CZ_FZU, FR_CCIN2P3, RU_JINR, BR_CBPF, CA_Victoria, IN_TIFR

Intelligent choices about matching jobs and file replicas reduces overall latency bottlenecks

🔶 Fermilab  DUNE

# Networking Setup and Timeline

- Planning 100 Gb/s primary link between SURF and FNAL

  - Enables 1 wk DAQ backlog to clear in 1 day in case of connection outage. DAQ has 99% uptime

  - Primary path currently 10 Gb

  - Secondary path now 1 Gb; will go to 100 Gb when commissioning starts

  - Tertiary path capable of meeting DAQ requirements during normal physics operations

- Networking between compute sites varies; leveraging existing setups at LHC compute sites where possible (e.g. LHCONE), other large-scale infrastructure (e.g. ESnet, GÉANT, NRENs, etc.)

  - Expect larger sites w/storage to have 40+ Gb/s ext. connectivity; medium sites 20+ Gb/s

- Monitoring tools include perfSONAR



Not shown:
10 Gb Tertiary path through Denver and KC

DUNE FD WAN Bandwidth Timeline Projections:

| Date | Stage of the experiment | Primary Path | Secondary Path | Tertiary Path |
|------|------------------------|--------------|----------------|----------------|
| Now | Cavern excavation | 10GE | < 1GE via SURF | none |
| 2025 | Detector construction | 10GE | < 1GE via SURF | none |
| 2027 | Computing/DAQ deployment | 100GE | 10GE | < 1GE via SURF |
| 2028 | Cryo deployment completed | 100GE | 100Gb/s | 10GE |
| 2029 | Start of science | 100GE | 100Gb/s | 10GE |

*REED: South Dakota Higher Education Network*

vLAN service provided by REED/GPN (shared)
Dedicated circuit Ross Dry Bldg. to Chicago
Dedicated circuit Yates Complex to Denver

# Physics sensitivities

- Sensitivity to CP violation for 50% of deltacp values, as a function of time in calendar years. The width of the bands shows the impact of potential beam power ramp up; the solid upper curve is the sensitivity if data collection begins with 1.2 MW beam power and the lower dashed curve shows a conservative beam ramp scenario where the full power is achieved after 4 years. The green band shows the Phase I sensitivity and the red band showss the Phase II sensitivity. The cyan band shows the Phase II sensitivity if the beam upgrade does not occur.

- From DUNE Snowmass physics whitepaper

🔬 Fermilab DUNE