

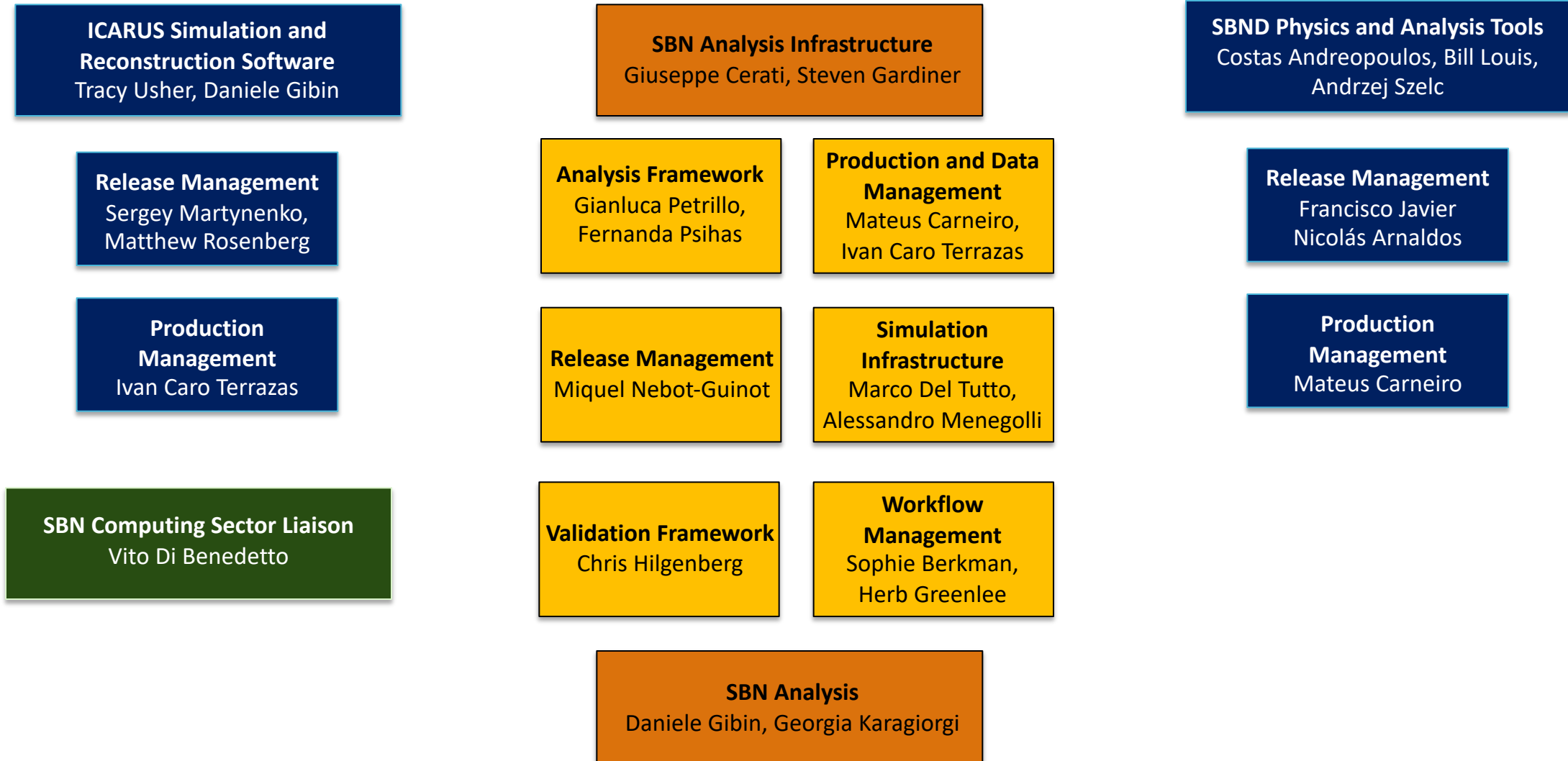


SBN Presentation for FCRSG 2023

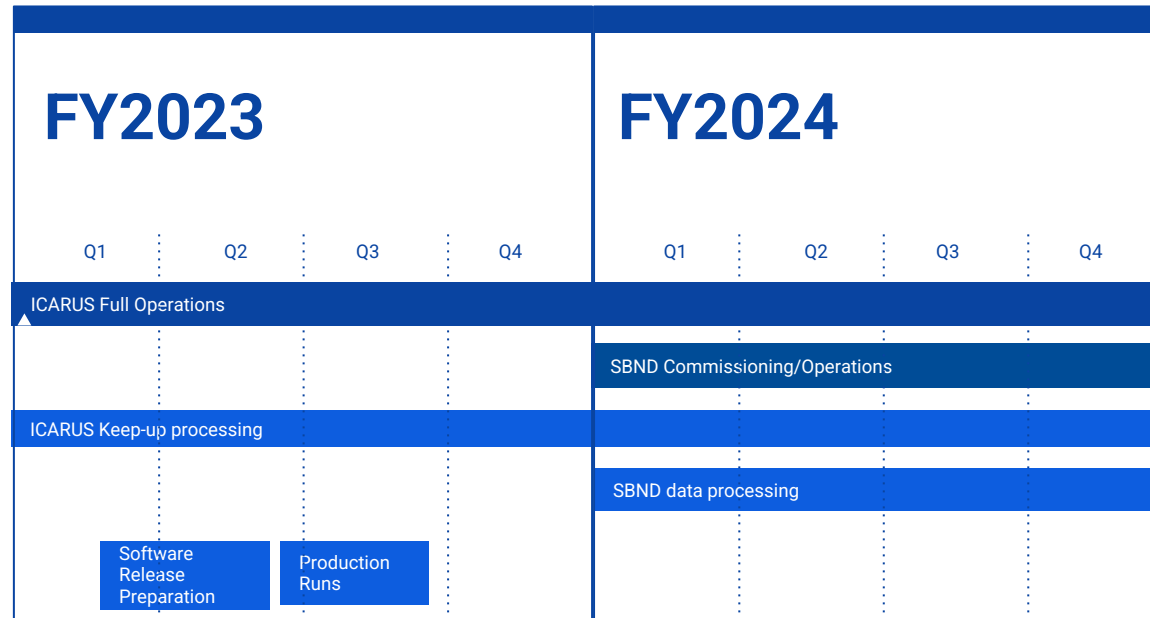
Giuseppe Cerati, Steven Gardiner, Wes Ketchum

15 February 2023

Organization Chart for Offline Computing



Overview Timeline

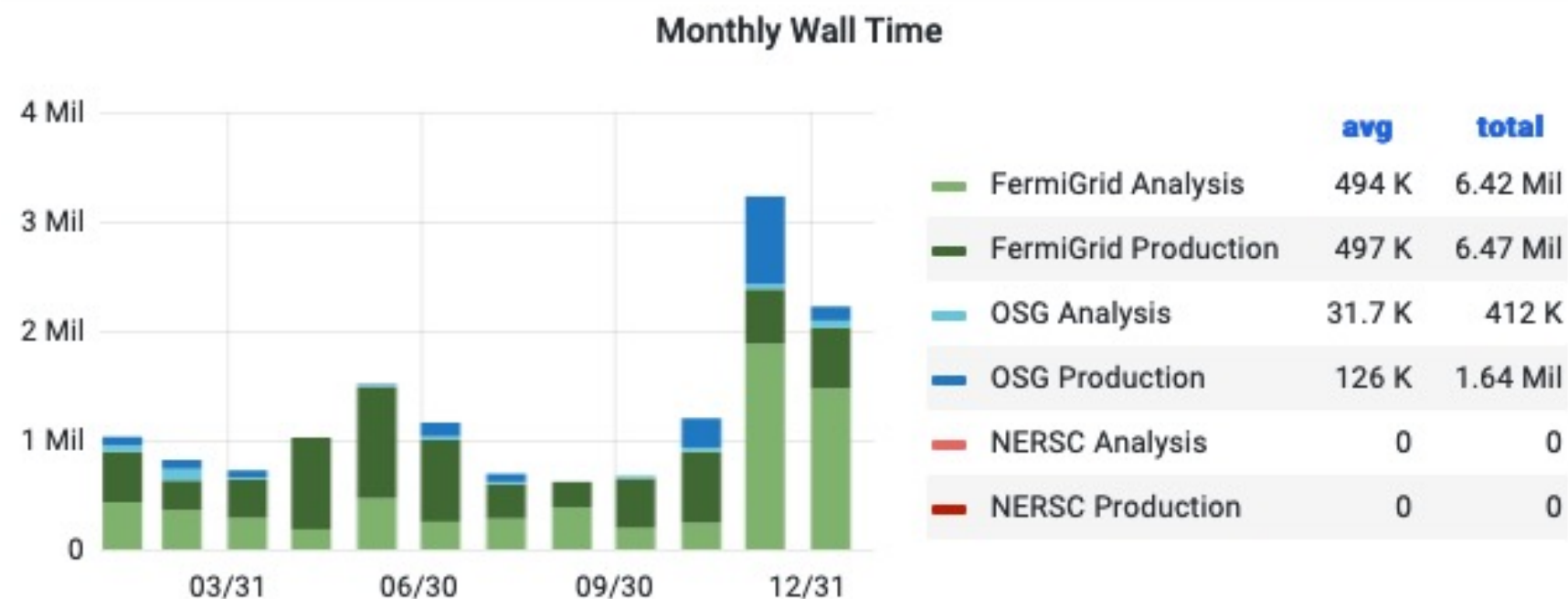


- ICARUS started physics data taking at the end of calendar year (CY) 2022
 - First physics results as early as FY2023
- SBND to enter commissioning in FY2024 and move (hopefully rapidly) to full operations

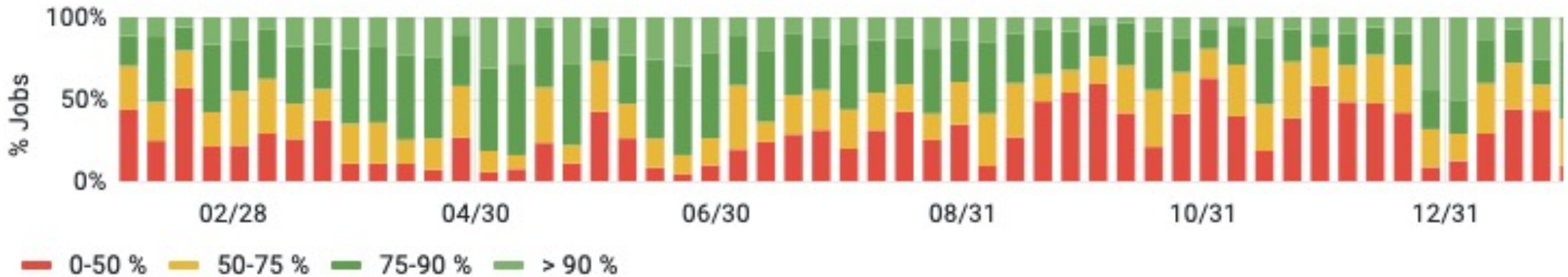
Computing model updates

- Run keep-up processing and save output of first stage reconstruction on disk
 - size of stage0 output is 6.5x smaller than compressed RAW
 - Allows for re-processing of data without reading back from tape
- Run MC campaigns from scratch at least once per year
 - Avoid tape bottleneck
 - Allows for updates to interaction and simulation model
 - Drop most artroot files and keep analysis file format on disk
 - Store first pass reco on tape in case of reprocessing
 - Keep 0.1-1M events on disk after full reco for development
- Main use of tape is for perennial archival data
 - and for storing other samples for limited time before deletion

ICARUS CPU Usage over the past year

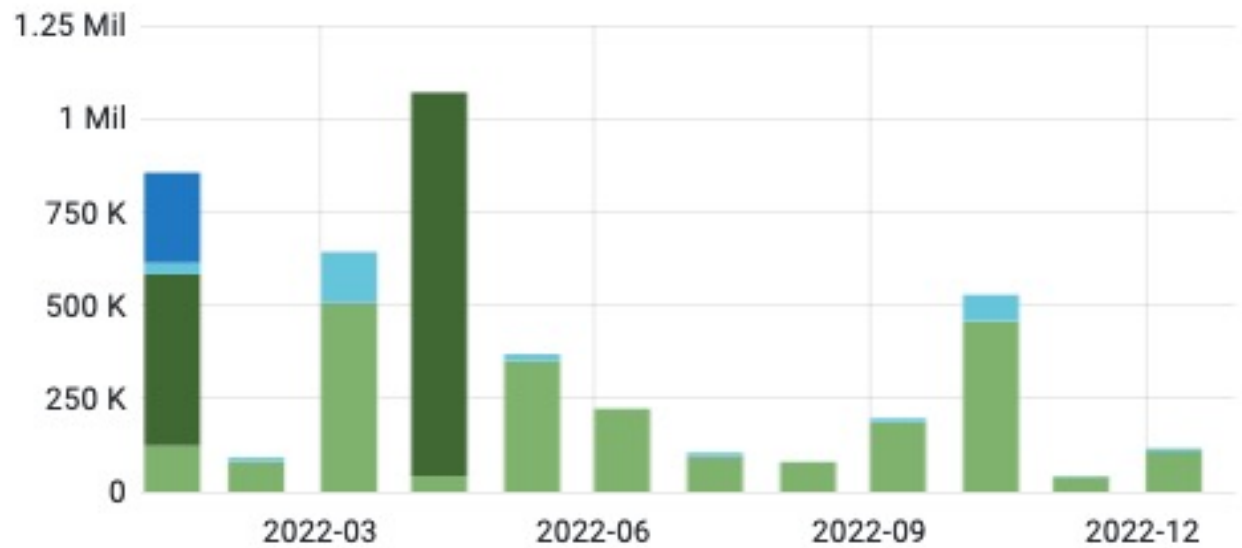


CPU Efficiency (CPU time / Wall time) (Combined Production and Analysis)



SBND CPU Usage over the past year

Monthly Wall Time



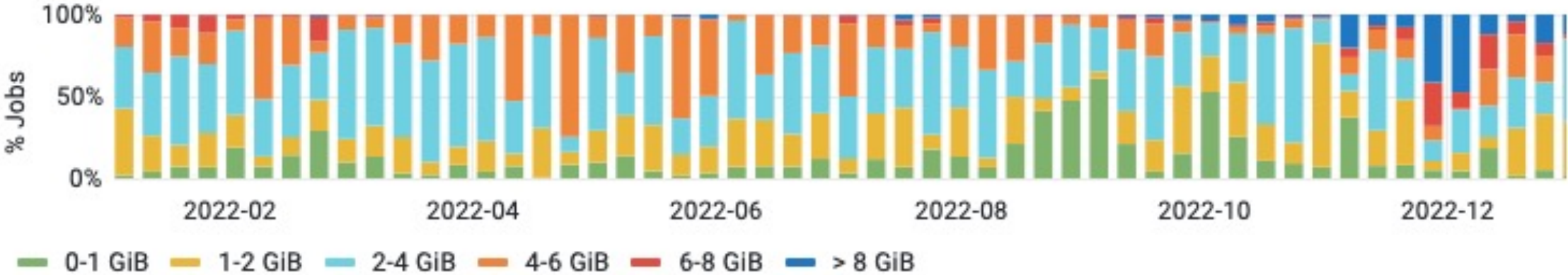
	avg	total
FermiGrid Analysis	176 K	2.28 Mil
FermiGrid Production	118 K	1.53 Mil
OSG Analysis	23.1 K	300 K
OSG Production	19.8 K	257 K
NERSC Analysis	0	0
NERSC Production	0	0

CPU Efficiency (CPU time / Wall time) (Combined Production and Analysis)

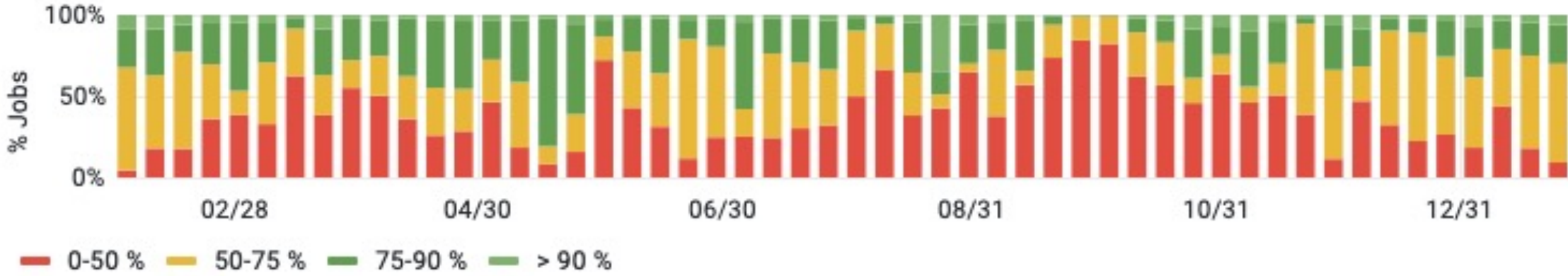


ICARUS memory usage over the past year

Memory Usage (Combined Production and Analysis)

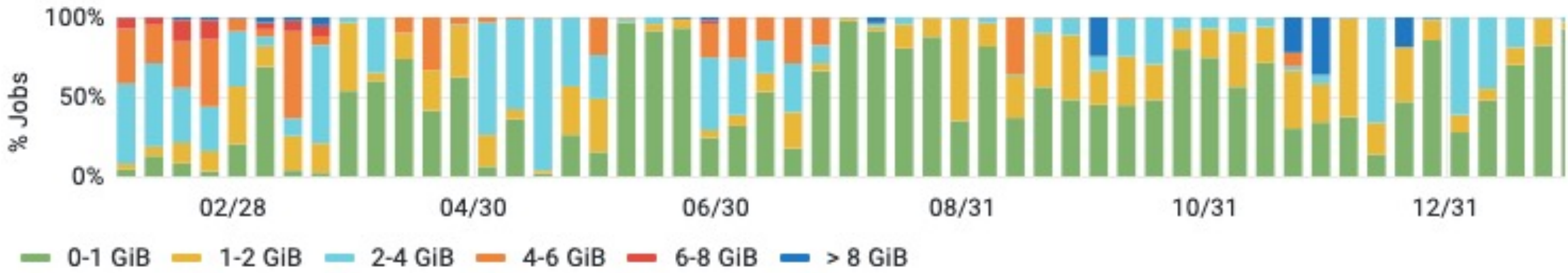


Memory Efficiency (Usage/Request) (Combined Production and Analysis)

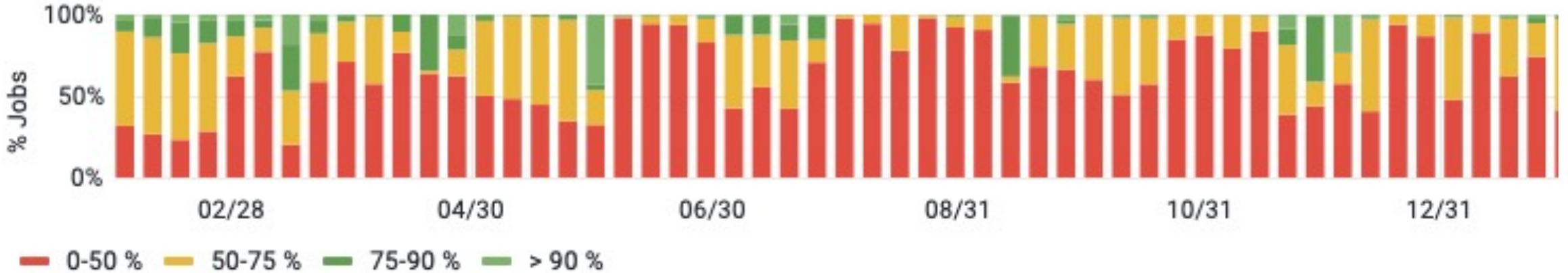


SBND memory usage over the past year

Memory Usage (Combined Production and Analysis)



Memory Efficiency (Usage/Request) (Combined Production and Analysis)



What do you want to achieve in computing over the next 5 years?

Goals	Where does the experiment need to contribute	Where does CSAID needs to contribute
Deploy PMT-CRT matching as software filter for early cosmic rejection	Experiment goal	N/A
Enable data overlay simulation	Experiment goal	May need support for optimization of workflow (e.g. if cosmic data needs to be read from tape)
Improve resource usage in production with multithreading and/or reduced memory footprint	Joint effort	Joint effort
Integration of external resources (CNAF) through RUCIO	Integration of non-FNAL resources in prod. workflows	Need to enable data transfers to FNAL, and integrated data management
Better integration of AI inference in production workflows	Joint effort	Centrally supported mechanism for running inference in production
Harmonize infrastructure between ICARUS and SBND (simulation, AI, etc.)	Experiment goal	N/A
Usage of HPC resources for production jobs	Need to transition from R&D to integrated resource	Need to transition from R&D to integrated resource
Enable usage of Elastic Analysis Facility for analysis and AI training	Define use cases, advertise in collaborations	User support (docs, tutorials, support of experiment resources)

ICARUS Campaign Schedules

	2023	2024	2025	2026	2027
Processing campaigns (start month-end month if known). Include when you expect to be prestaging	Data: All Year keepup process MC Campaign 1: Mar-May MC Campaign 2: Aug-Oct No significant prestaging	Data: All Year keepup process MC Campaign 1: Jan-Mar MC Campaign 2: Aug-Oct Prestaging of MC1 for MC2	Data: All Year keepup process Overlay Camp. 1: Jan-Mar Overlay Camp. 2: Aug-Oct Prestage OffBeam data Prestaging of MC1 for MC2	Data: All Year keepup process Overlay Camp. 1: Jan-Mar Overlay Camp. 2: Aug-Oct Prestage OffBeam data Prestaging of MC1 for MC2	Data: All Year keepup process Overlay Camp. 1: Jan-Mar Overlay Camp. 2: Aug-Oct Prestage OffBeam data Prestaging of MC1 for MC2
Storage + CPU estimates (call out any special resource needs if known, e.g. HPC or GPU). Include amount(s) to be prestaged and file families, in addition to space needed for new outputs. Note: currently not including deletions from disk nor tape (see later slide)	Data: CPU 4.5M hrs Write 3.7 PB to tape Write 1.9 PB to disk MC1: CPU 5.7M hrs No significant prestaging No write to tape Write 50 TB to disk MC2: CPU 5.7M hrs No significant prestaging Write 50 TB to tape Write 50TB to disk	Data: CPU 5.7M hrs Write 3.7 PB to tape Write 0.9 PB to disk MC1: CPU 10.4M hrs No significant prestaging Write 4.4 PB to tape Write 80 TB to disk MC2: CPU 5.5M hrs Prestage 4.4 PB (SIM) Write 80 TB to tape Write 80 TB to disk	Data: CPU 4.3M hrs Write 3.7 PB to tape Write 0.9 PB to disk MC1: CPU 9.8M hrs Prestage 4.9 PB (rawoffbeam) Write 6.7 PB to tape Write 100 TB to disk MC2: CPU 2.8M hrs Prestage 6.7 PB (SIM) Write 100 TB to tape Write 100 TB to disk	Data: CPU 4.9M hrs Write 3.7 PB to tape Write 0.9 PB to disk MC1: CPU 12.8M hrs Prestage 6.5 PB (rawoffbeam) Write 8.8 PB to tape Write 130 TB to disk MC2: CPU 3.7M hrs Prestage 8.8 PB (SIM) Write 130 TB to tape Write 130 TB to disk	Data: CPU 5.5M hrs Write 3.7 PB to tape Write 0.9 PB to disk MC1: CPU 15.7M hrs Prestage 8.1 PB (rawoffbeam) Write 11 PB to tape Write 150 TB to disk MC2: CPU 4.5M hrs Prestage 11 PB (SIM) Write 150 TB to tape Write 150 TB to disk
Conference or result targets (month if known)	First ICARUS results	Neutrino 2024			

SBND Campaign Schedules

	2023	2024	2025	2026	2027
Processing campaigns (start month-end month if known). Include when you expect to be prestaging	No Data MC Campaign 1: Mar-May (or later) No significant prestaging	Data: All Year keepup process MC Campaign 1: Jan-Mar MC Campaign 2: Aug-Oct No significant prestaging	Data: All Year keepup process Overlay Camp. 1: Jan-Mar Overlay Camp. 2: Aug-Oct No significant prestaging	Data: All Year keepup process Overlay Camp. 1: Jan-Mar Overlay Camp. 2: Aug-Oct Prestage OffBeam data Prestaging of MC1 for MC2	Data: All Year keepup process Overlay Camp. 1: Jan-Mar Overlay Camp. 2: Aug-Oct Prestage OffBeam data Prestaging of MC1 for MC2
Storage + CPU estimates (call out any special resource needs if known, e.g. HPC or GPU). Include amount(s) to be prestaged and file families, in addition to space needed for new outputs. Note: currently not including deletions from disk nor tape (see later slide)	MC1: CPU 6.2M hrs No significant prestaging No write to tape Write 50TB to disk	Data: CPU 10M hrs Write 4.2 PB to tape Write 2.1 PB to disk MC1: CPU 6.2M hrs No significant prestaging Write 50 TB to tape Write 50 TB to disk MC2: CPU 6.2M hrs No significant prestaging Write 50 TB to tape Write 50 TB to disk	Data: CPU 12M hrs Write 3.2 PB to tape Write 1.6 PB to disk MC1: CPU 7.7M hrs No significant prestaging Write 50 TB to tape Write 60 TB to disk MC2: CPU 7.7M hrs No significant prestaging Write 60 TB to tape Write 60 TB to disk	Data: CPU 15M hrs Write 3.2 PB to tape Write 0.8 PB to disk MC1: CPU 10.5M hrs Prestage 5.9 PB (rawoffbeam) Write 1.3 PB to tape Write 100 TB to disk MC2: CPU 5.8M hrs Prestage 1.3 PB (SIM) Write 100 TB to tape Write 100 TB to disk	Data: CPU 18M hrs Write 3.2 PB to tape Write 0.8 PB to disk MC1: CPU 15.0M hrs Prestage 8.5 PB (rawoffbeam) Write 1.8 PB to tape Write 130 TB to disk MC2: CPU 8.3M hrs Prestage 1.8 PB (SIM) Write 130 TB to tape Write 130 TB to disk
Conference or result targets (month if known)	MC studies for summer/fall conferences	Neutrino 2024			

SBN CPU @ Fermilab Prediction Going Forward and Accuracy of the Predictions [units of Million (1 CPU, 2GB) wall hours per CY]

	2019	2020	2021	2022	2023	2024	2025	2026	2027
Requested (could have multiple values for different MWC combinations)	13.2	16.5	19.5	27	29	53±15	48±15	58±15	74±15
Actual Used	4.46	7.7	11.4	15	N/A	N/A	N/A	N/A	N/A
Efficiency	33.8%	46.7%	58.5%	55.6%	N/A	N/A	N/A	N/A	N/A

Assumed that efficiency = requested / actual



CPU – non-FNAL HTC Resources Going Forward and Accuracy of Your Predictions [units of Million (1 CPU, 2GB) wall hours per CY]

We plan to use OSG in opportunistic mode.

We've been using it for up to 20% for production jobs.

ICARUS has allocation at CNAF (Italy), re-negotiated each year.

Currently 3M CPU hours (1 PB disk, 2 PB tape), but needs deployment of RUCIO transfers from CNAF to FNAL to be fully integrated in workflows.

Short term plan is to use this resource for reprocessing of data from runs before summer 2022.

CPU – HPC Resources Going Forward and Accuracy of Your Predictions [units of Million (1 CPU, 2GB) wall hours per CY]

HPC resources used by individual groups, not by SBN experiments themselves (no centralized accounting at present). Current R&D efforts may lead to centralized requests in future planning. Notable examples:

ICARUS ML workflow: plan to run end-to-end ML reconstruction at HPC centers (GPU) using allocation from a collaborator's grant. Dedicated input files need to be transferred from FNAL to offsite storage during or shortly after production. Analysis file format eventually needs to be accessible to the full collaboration.

SBNfit: Run Feldman-Cousins fits for sensitivity studies of SBN oscillation analyses at NERSC, using development version of the code specific for HPC. Leverages SciDAC4 allocation.

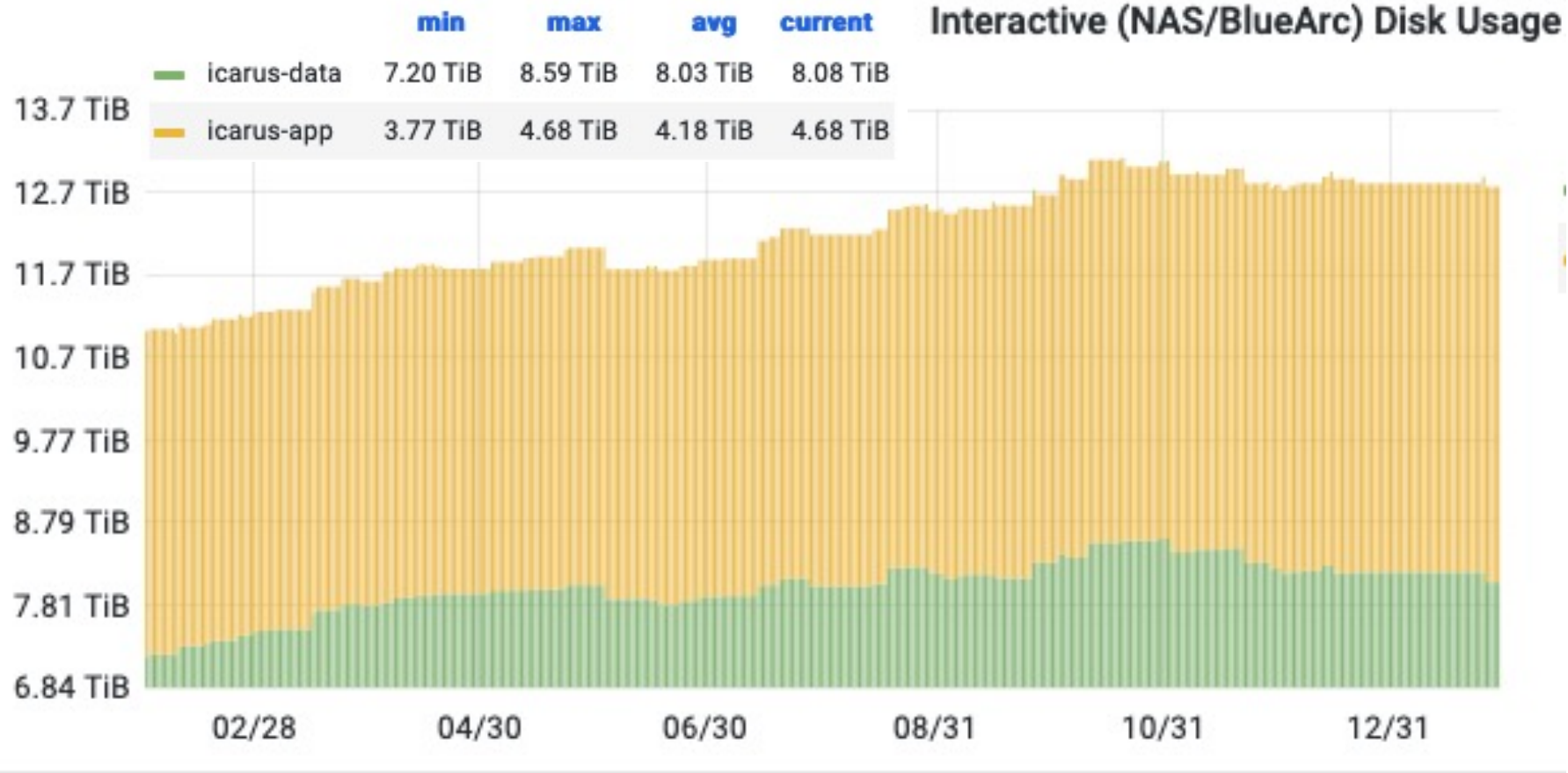
HPC LArSoft Workflow: development of HPC processing workflow based on Spack built central LArSoft code is under test on Theta@ALCF. Part of SciDAC projects as well.

CPU – GPU Resources Going Forward and Accuracy of Your Predictions [units of Million (1 CPU, 2GB) wall hours per CY]

No current roadmap to evolving central SBN software to run on GPUs. Support of GPU offloading for common tasks in LArSoft (e.g. Geant4?) can be a game changer.

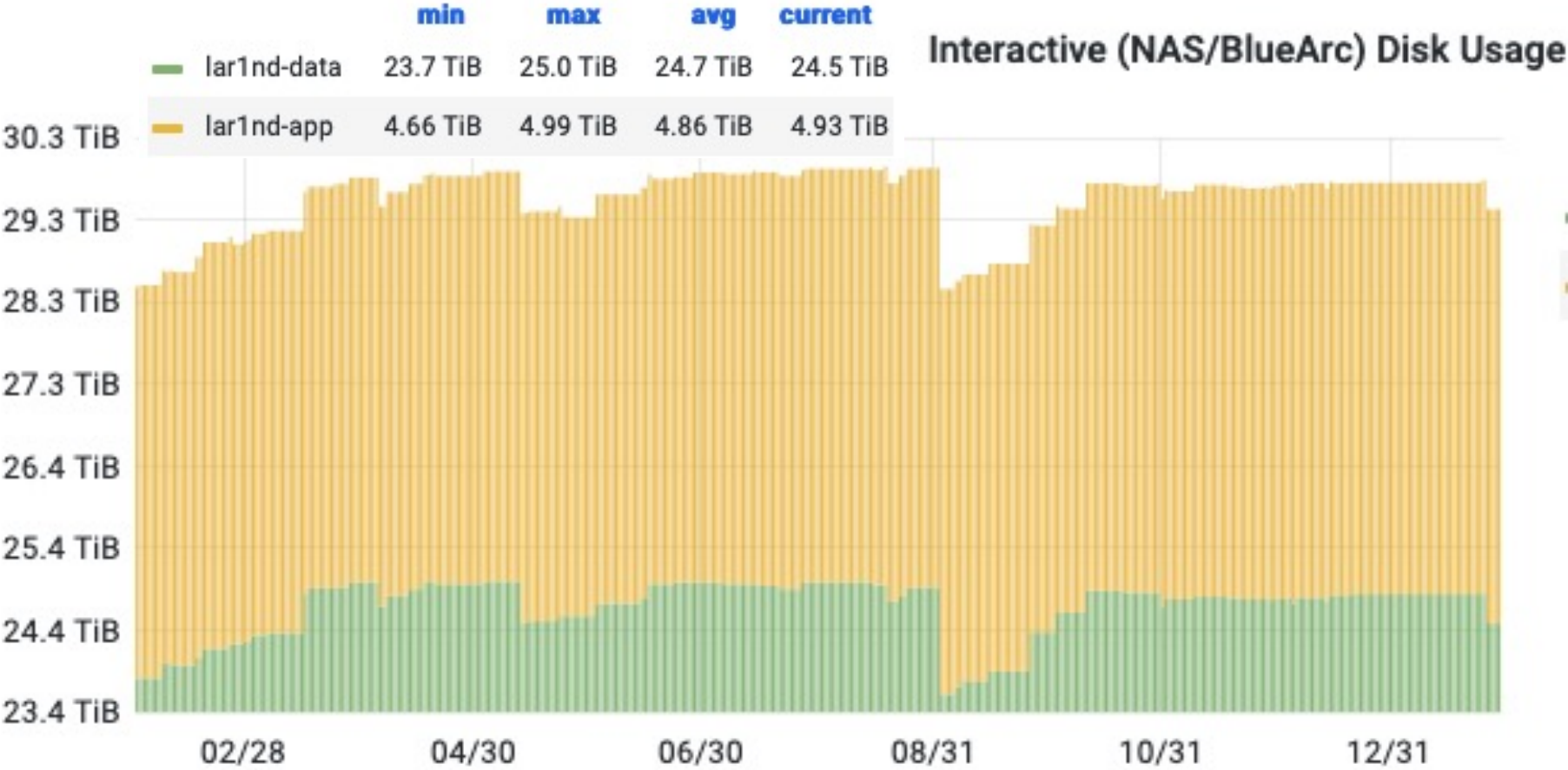
GPUs are used for ML training and inference, both onsite (see also later slide on Elastic Analysis Facility) and offsite.

ICARUS NAS+Ceph Usage and Projections



	App [TB]	Data [TB]
2022	5	9
2023	5	25
2024	6	30
2025	7	35
2026	8	40
2027	8	40

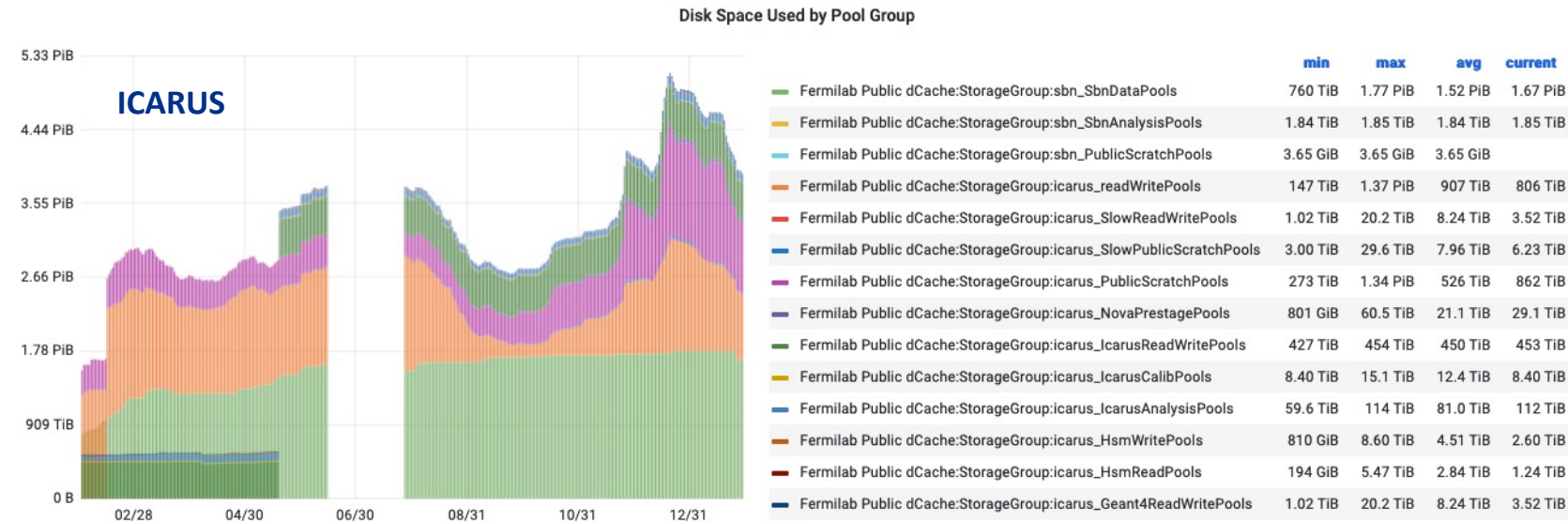
SBND NAS+Ceph Usage and Projections



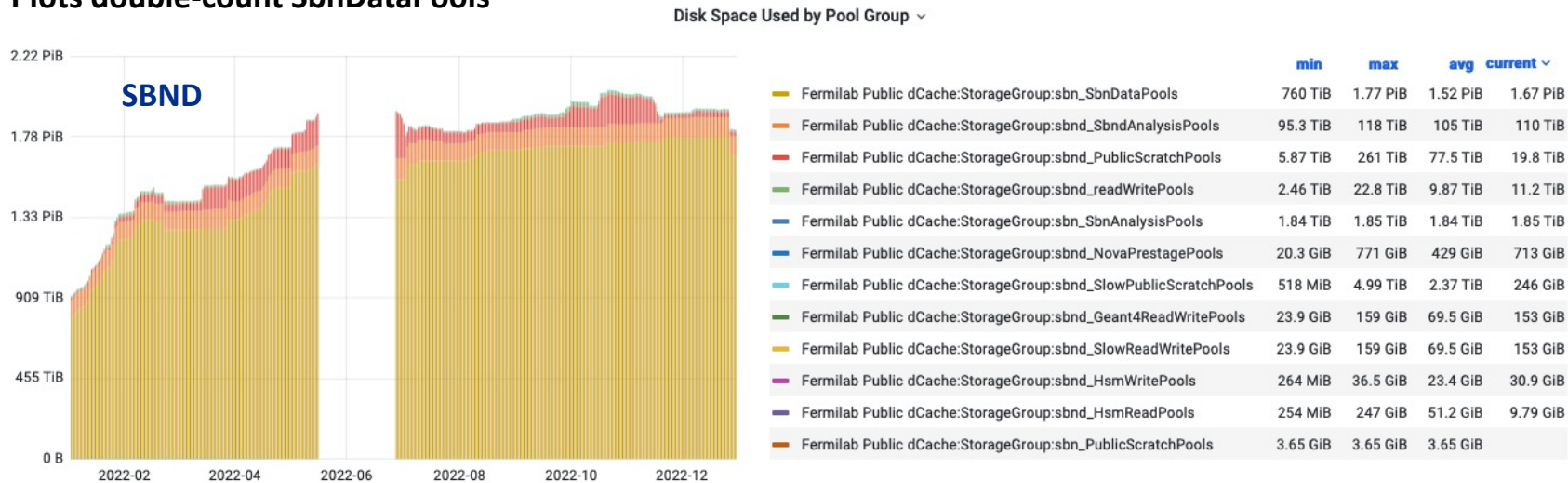
	App [TB]	Data [TB]
2022	5	25
2023	7	30
2024	8	35
2025	8	40
2026	9	45
2027	9	45

ICARUS + SBND dCache Usage and Predictions

ICARUS has 1 PB allocated at CNAF (currently 60% used)



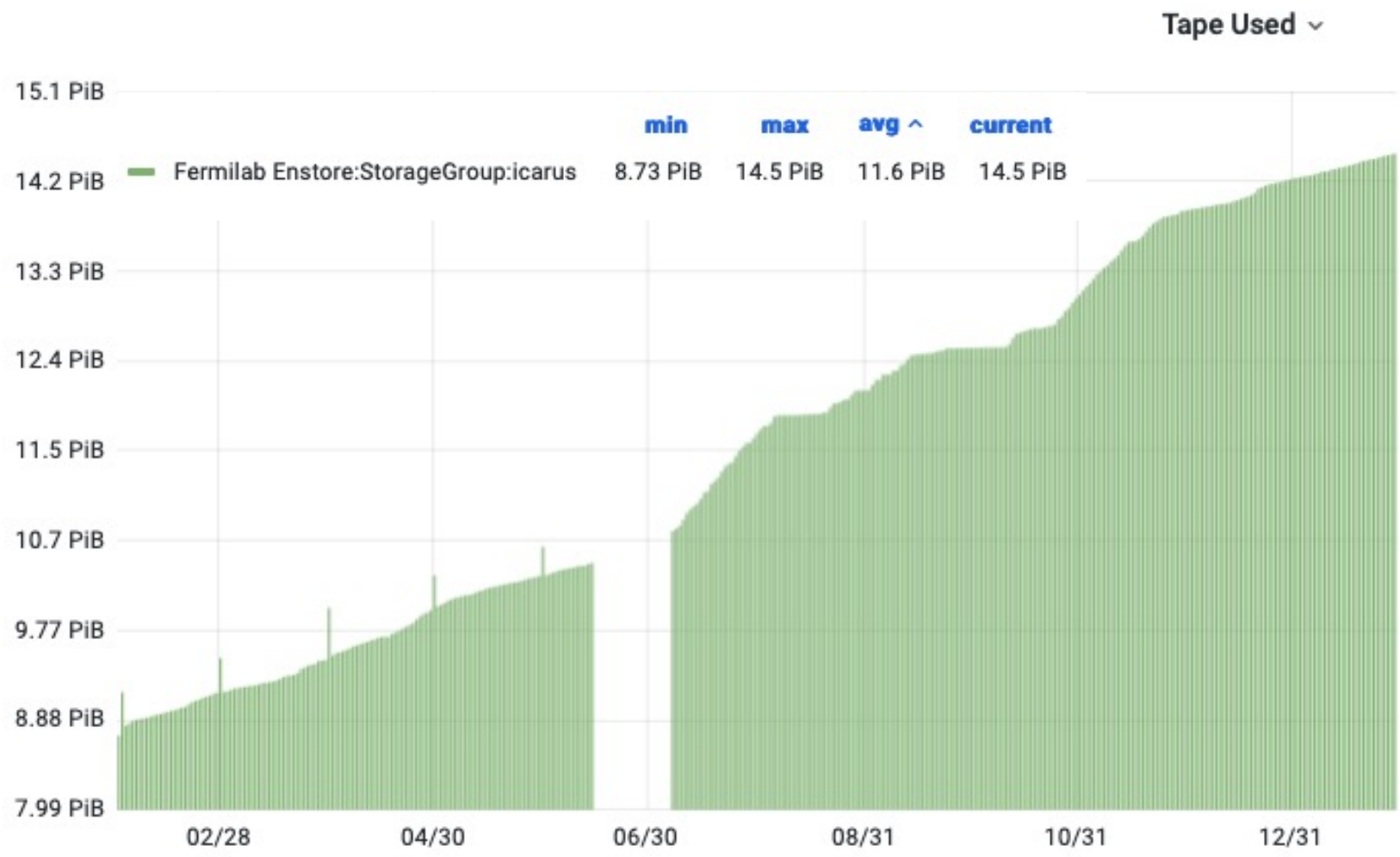
Plots double-count SbnDataPools



	Analysis (Persistent)	Other Dedicated (Write)
Current	1.9 PB (actual)	0.8 PB (actual)
2023	3.3 PB	1 PB
2024	5.7±1 PB	2 PB
2025	5.9±1 PB	2 PB
2026	6.2±1 PB	2 PB
2027	7.3±1 PB	2 PB

Analysis = IcarusAnalysis + SbnAnalysis + SbnData pool
Other = readWrite pool

ICARUS Tape usage and predictions



	Total Added By End of Year
At end 2021	8.7 PB
2022	+3.9 PB
2023	+3.8 PB
2024	+8±2 PB
2025	+11±5 PB
2026	+4±5 PB
2027	+1±5 PB

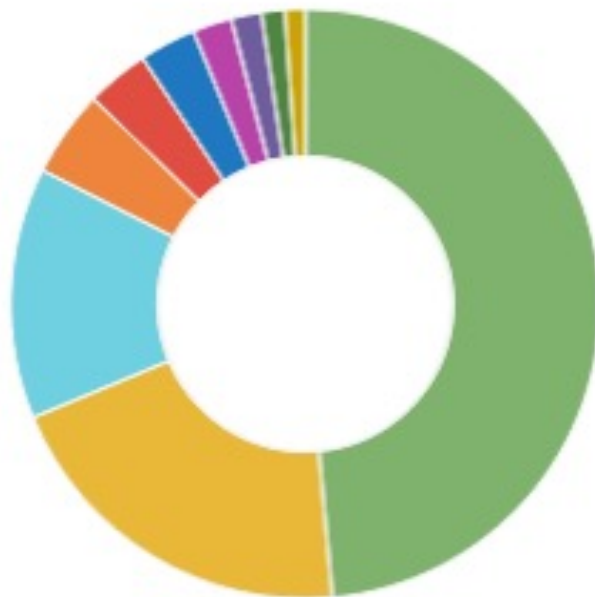
ICARUS Tape usage and predictions

File families organized by data stream.

Half of volume is in generic family

We need more granular definitions, e.g.:

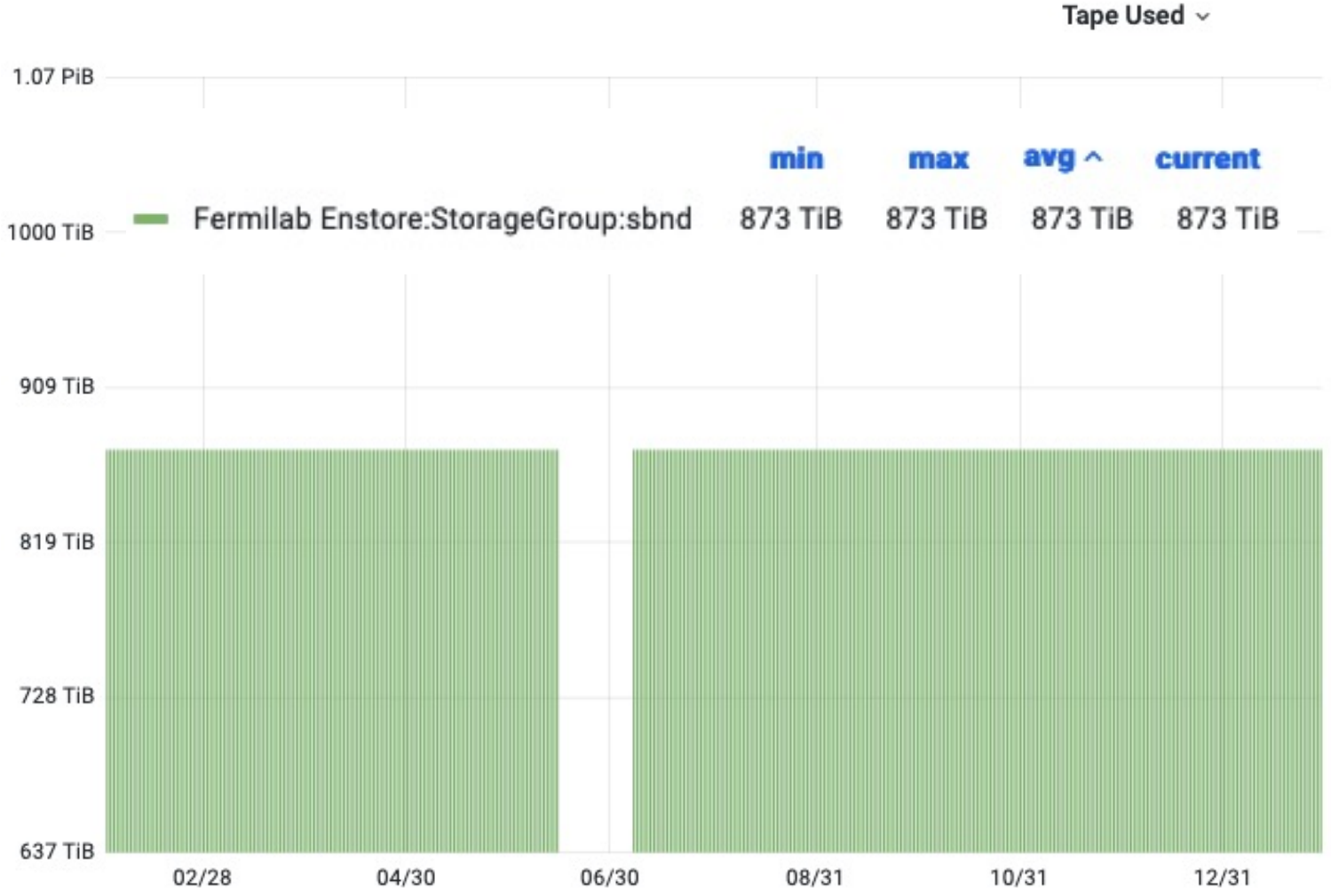
- split by run?
- artroot MC?
- analysis files?



Data Volume by File Family

	current	percentage
icarus	7 PB	49%
data_artroot_raw_offbeambnbminbias	3 PB	20%
raw	2 PB	14%
data_artroot_reconstructed_bnb	729 TB	5%
data_artroot_raw_bnb	527 TB	3%
data_artroot_raw_unknown	477 TB	3%
data_artroot_raw_offbeambnb	329 TB	2%
data_artroot_raw_num1	254 TB	2%
data_artroot_raw_offbeamnum1	195 TB	1%
icarusdata_artroot_reconstructed_num1	171 TB	1%

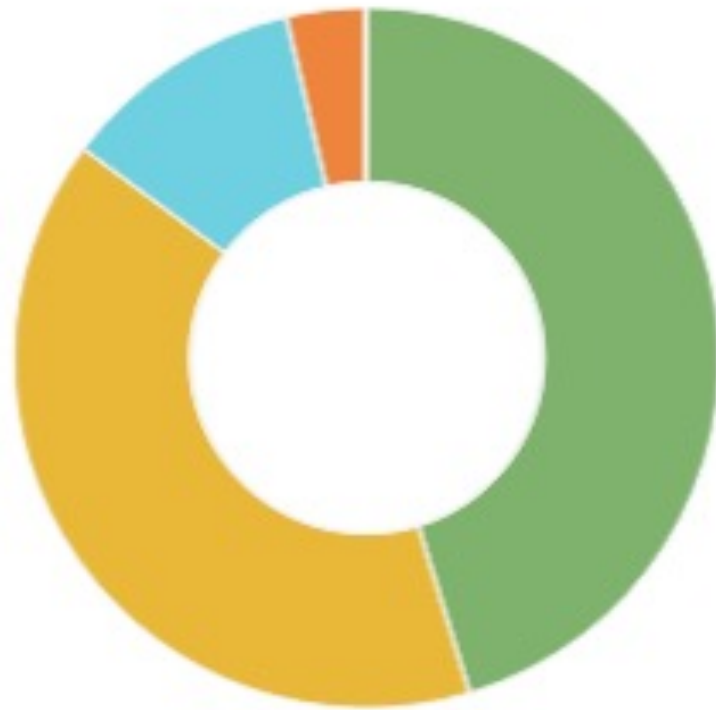
SBND Tape usage and predictions



	Total Added By End of Year
At end 2021	+0.9 PB
2022	+0 PB
2023	+0±1 PB
2024	+4.3±2 PB
2025	+3.4±2 PB
2026	+4.6±2 PB
2027	+5.2±2 PB

SBND Tape usage and predictions

Data Volume by File Family



	current	percentage
— sbnd	400 TB	45%
— mc_reco	353 TB	40%
— mc_detsim	99 TB	11%
— mc_g4	32 TB	4%
— mc_anatree	318 GB	0%
— mc_gen	16 GB	0%
— mc_root_small	4 GB	0%

Data Lifetimes

Summary of plan for data lifetimes:

- Raw (detector) data: permanent on tape
 - At the beginning of a new run, keep a few months of RAW data on disk to allow fast reprocessing in case calibrations are not ready
- First pass reco data: keep permanently on disk for reprocessing
- Other derived data: 1-2 year lifetime
- Simulation: 1-2 year lifetime

Analysis Facility Use

- A few SBN users have started using the facility for ML training
 - In the longer term we envision a combination of analysis work + ML training
- Mounting /pnfs would be great
 - xrootd already works, but not sure if there are bandwidth concerns for massive usage
 - Our analysis files take $O(100)$ TB per campaign
- Likely combination of NVMe (solid state) + HDD storage would be optimal
 - Just a guess at this point
- Mixed CPUs + GPUs needed given the different applications

Backup