

Jefferson Lab Scientific Computing Infrastructure Update

CLAS12 Collaboration Meeting
March 2026

Brad Sawatzky

Friday, March 13, 2026

Jefferson Lab



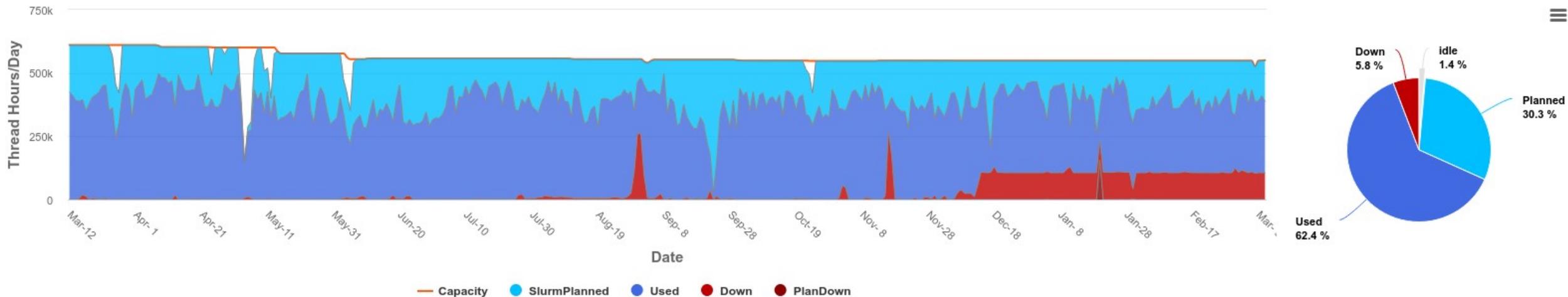
Jefferson Lab's High Throughput Computing – The Farm

- Farm has many components
 - ~30000 compute threads
 - ~11 PB Lustre
 - ~5 PB NFS/XRootD (ZFS)
 - ~100+ PB of (online) Tape
 - Consumes >400kW
 - GPU nodes available (sciml*)
 - Interactive nodes (ifarm240x)
- Growth is \$\$\$ and based on projections from Halls
 - Expenditures *usually* switch between storage + CPU every other year
 - JLab Scientific Computing projections assessed annually
 - » What actually happens is driven by provided numbers



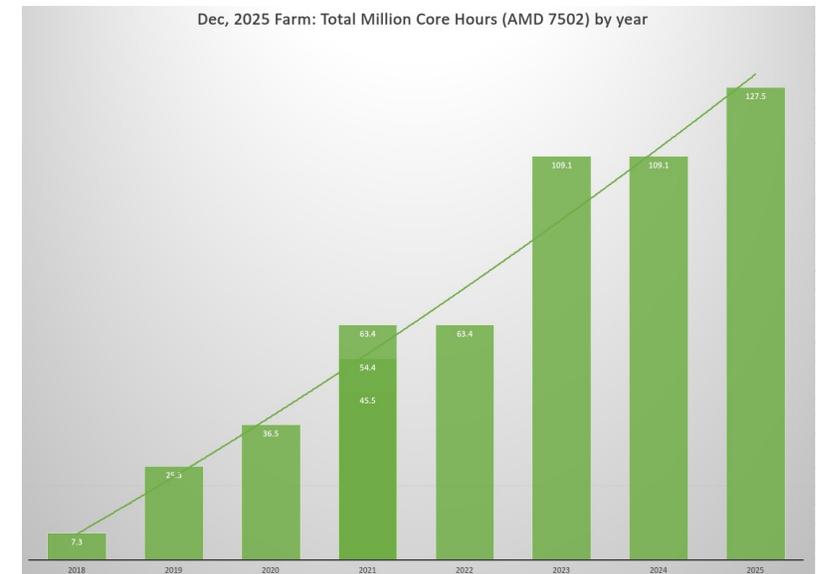
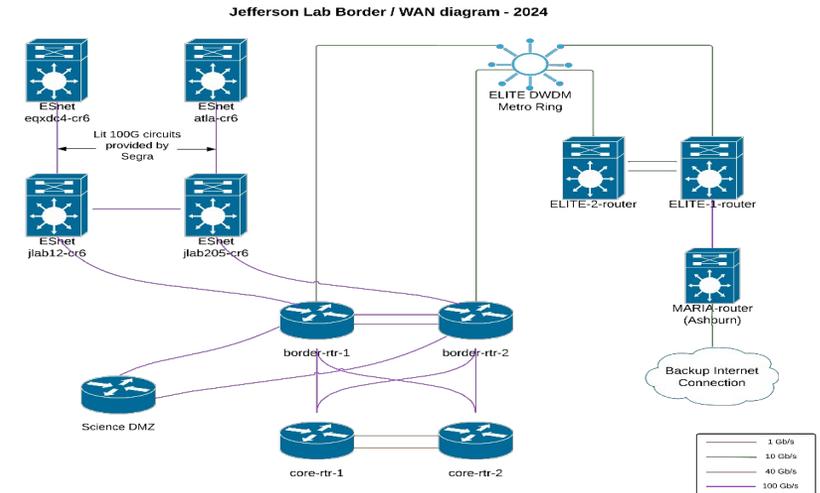
Jefferson Lab's High Throughput Computing – The Farm

- The Farm is routinely busy
 - General Utilization > 90% for years
 - Last 6+ months: ~100%
- Light Blue '*Planned*' are not 'idle' resources, they indicate thread-slots held by slurm for other optimizations:
 - Waiting for RAM or CPUs to free in order to allocate a priority job
 - Minimize hyperthread overhead (waste) for single thread jobs (most of the light-blue)
 - This 'reserved thread' is also billed against the single-threaded job (but that job also runs, roughly, 2x as fast as if it was sharing the core with a 2nd hyperthread)
- Red 'Down' is a problem
 - 20% of the Farm has been offline since December due to Purchasing/Budget choices by upper management. (PO to address underlying issue was 'stuck' for a year. Came to a head in Dec...)



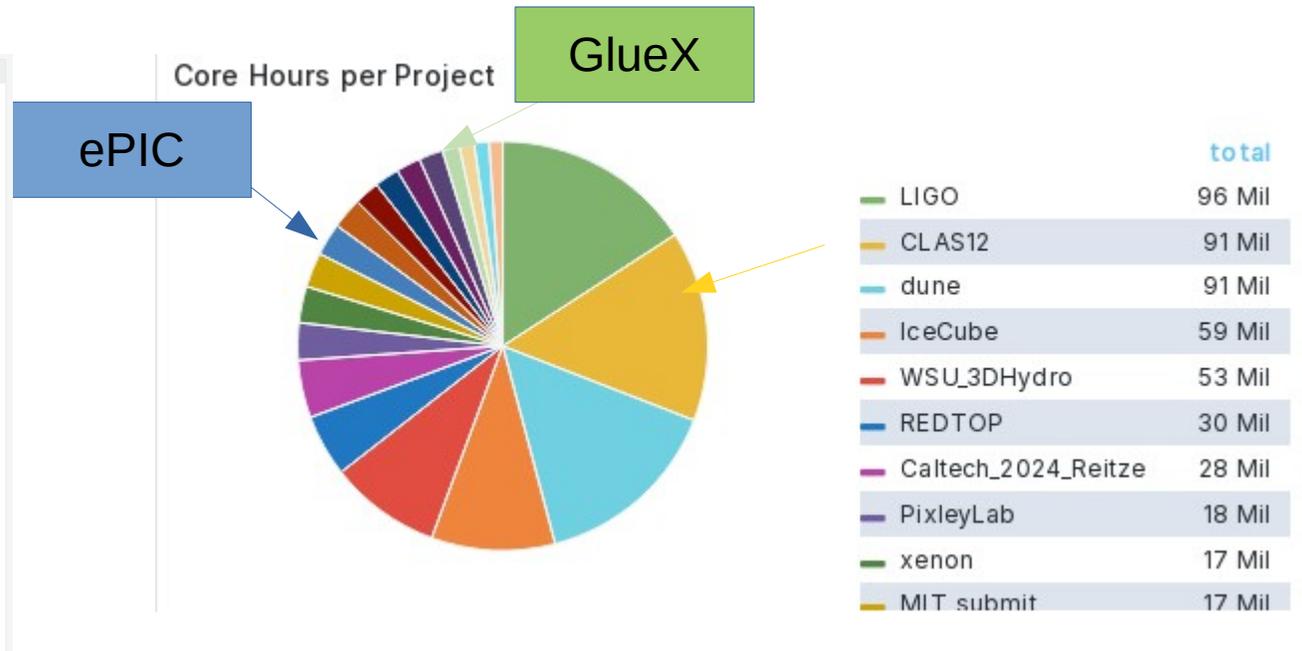
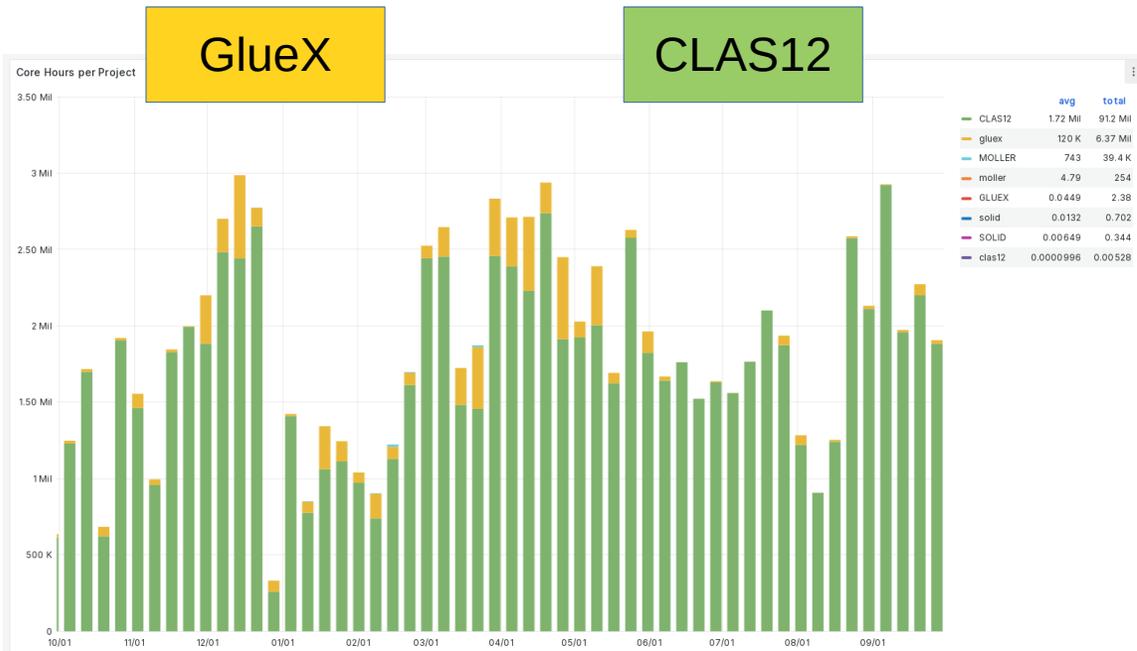
Recent and Near-term Infrastructure Updates (HW)

- “Farm25” CPU node installation (~~late summer 2025~~) *Now, finally...*
 - Held up behind procurement/budget delays...
- 20 new A800 GPUs installed (done)
 - Funded by Hall D and Data Sci LDRDs, but available to all
- Intel Gaudi2 nodes (~~late fall 2025~~) *End of March, finally...*
 - Also held up by procurement/budget delays...
- “/work” disk increase (done)
 - ~ 1 PB of new storage at /ceph24/... (cephFS backed)
- +4 tape drives brought online (done)
 - Tape system throughput is ~12.5 GB/sec
 - Significant optimizations to SW and HW backing Jasmine in May/July
- **Current plan is to buy CPU in FY26 “Farm26”**
 - Slight rev of Farm25 build; more compute but less RAM/thread (while maintaining 2GB/thread minimum)
 - Getting availability and quotes from vendors now



Open Science Grid Processing (2025 numbers)

- The Open Science Grid has been a significant resource for Monte Carlo Simulation Compute Cycles. GlueX and CLAS12 are significant consumers of CPU cycles
 - CLAS12 is in the top 3 OSG compute consumers (91 Million Core Hours [MCH])
 - ePIC is in the top 10 (16 MCH)
 - GlueX (8.7 MCH)
- OSG future funding is uncertain, was NSF (cut 9m ago; new proposal pending)



Infrastructure Updates (SW)

- code.jlab.org (GitLab instance)
 - CI/CD pipelines
 - Container and Model registries
 - JLab GitHub Org will remain while cost-effective
- [CVMFS 'Stratum 0' for JLab](#)
 - [/cvmfs/jlab.opensciencegrid.org/](https://cvmfs/jlab.opensciencegrid.org/) ← **NEW**
 - will replace [/cvmfs/oasis.opensciencegrid.org/jlab/](https://cvmfs/oasis.opensciencegrid.org/jlab/) at some point
- [Kubernetes for workflows that don't fit Batch model](#)
 - OpenShift 'enterprise' K8 platform rolled out in 2024; it was aggravating...
 - » required significantly more 'backend' work than advertised
 - » CI/CD now stable; other (internal) projects being deployed; adding some hardware to the cluster Mar 2026
 - **more conventional K8 deployment being rolled out in TestBed**
 - » This is probably where Users should start

Rucio

- Distributed (large-file) data management framework
- EIC simulation campaigns in full production
- GlueX and MOLLER in progress
- Swif+Rucio file URI integration in progress
 - » `rucio://.....`

[JLab Research DB](#)

- "One stop shop" to locate data, publications, workflow information, logbook references, etc...
- *Ties into DOE and JLab Data Management policy changes to be announced soon*

• New [/volatile deletion queue](#) page

- Can see what is 'next in the queue' for policy based auto-clean up



Infrastructure Updates (SW): 2025

- LLM frontend: <https://delphi.jlab.org>
 - liteLLM + vLLM + LibreChat infrastructure
 - » Based on LBL [CBORG](#) deployment
 - Provide a OpenAPI based service to a variety of LLM services
 - » chatBot, RAG frontend
 - » inference and training via API
 - Hosted at JLab
 - » allows for many different (open access) models to be used
 - » GPT-OSS* models (up to 120b params)
 - » Llama, Gemma, etc.
 - Framework allows the same infrastructure to provide metered access to commercial offerings as well
 - » ChatGPT, Gemini, etc..
 - » *(Still!!!) awaiting JLab + Site-office approval to provide this...*



- **Goals:**
 - increase ease of access to GPU resources
 - » can point your laptop/desktop at the service as well as slurm jobs (when appropriate)
 - 'hide' LLM hardware deployment complexity from the user
 - » better leverage non-CUDA HW assets like GaudiX architectures
 - support adoption of LLM supported workflows
- **Nothing comes free**
 - We are balancing demand on GPU resources (Farm + Jupyterhub)
 - Will replace/augment with new HW when available
 - » ie. Gaudi nodes (soon) and/or new nVidia nodes in FY26 (DataSci/EPSCI)

JLab NP Computing Situation

Core Problem

- We do not have the compute necessary to meet your projected requirements
 - Driven by procurement and budget delays
 - Farm25 nodes sat on the shelf for 8 months waiting on Facilities, power, and procurement/budget...
 - Farm26 procurement has been 'on the bubble'
 - Seems we have the budget restored (for now), but AI boom has driven costs up a lot – I'm expecting long lead and less compute/\$.
 - New Disk is going to have to wait until 2027. It will also have to compete with a tape library upgrade to LTO10
 - OSG future is complicated – a loss of those resources would be a huge problem
 - Seems unlikely, but I expect the situation to get harder before it gets better...
 - ***Elephant in the room is Data Center cooling and power (next slides)***
- **Upshot is we must use the capacity we have efficiently**
 - Greater oversight and training on the part of Hall Compute Coordinators to support effective use of the Farm could help
 - Low efficiency jobs: <https://scicomp.jlab.org/scicomp/slurmJob/cpuEff>
 - Identify inefficient workflows / code
 - Worth looking into profiling and optimization for core code
 - User error / Insufficient testing before large campaign
 - Users submitting O(100k) jobs that burn allocation and then notice a problem

Disk Storage Areas and Their Uses

- **We know disk space is an issue**
 - High performance, reliable disk and associated infrastructure is still expensive and has been long lead...
 - We attempted to get ahead of the curve on disk in FY24 but JLab canceled a planned Lustre disk purchase.
 - FY26 will be CPU again.
 - Next disk purchase will be FY27
 - NEW CephFS disk storage (standard POSIX) added to /work in Fall FY25 (~1 PB) – done

We need to use what we have effectively

- /volatile is good for large files, streaming, large block I/O, production farm runs.
 - Lustre is *not* good for small files, high IOPS, and frequent metadata operations (worst case: open, write 1kB, close, repeat)
- /work will not scale for large farm campaigns
 - use for small files; modest IOP loads, etc
- /home is for personal 'desktop' files; not source code or analysis data
- /group is for collaboration/team critical files
- Node-local /scratch is good for jobs with high IOPS to working files.
 - **Note: SWIF-declared MSS files are automatically copied to node-local working directory**

Path	Best Use	FS Type	Deletion	backup
/cache	Bulk I/O, Migration to tape	Lustre	Once on tape	/mss
/volatile	Bulk I/O Temporary storage	Lustre	auto	NO
/work	Source code, DB files, exe's, etc. User Managed	NFS+ZFS	manual	NO
/home	Dot files, personal documents, etc	NFS ssd	manual	YES
/farm_out	Farm job stdout/stderr	NFS ssd	auto	NO
/group or /scigroup	Source code Papers, thesis, analysis scripts	NFS ssd	Manual	YES
/scratch	Farm job I/O to node local disk	ssd	auto	NO
/u/scratch	CUE scratch. <i>Deprecated</i> (Unavailable on el9)			
/cvmfs	Software stack. Configuration.			

“Tips and Tricks” for New Collaborators

Getting Started on

→ <https://scicomp.jlab.org>

Scientific Computing username [Getting Started](#) [Support](#)

Jlab Scientific Computing

Welcome to the Jefferson Lab Scientific computing home page. [New users start here.](#)

Feb-04-26 Enforcement of /scratch Disk Requests on the JLab Compute Farm On the February 17th maintenance day, requests will be strictly enforced. Jobs will be limited to the amount of /scratch space they request. This change is intended to improve overall Farm reliability by preventing /scratch filesystems from filling up and causing ENOSPC ("no space left on device") errors for other jobs. As a result, jobs that do not request disk space—or request too little—may begin to fail if they previously relied on unreserved /scratch space being available. What you may need to do

- Swift users: Swift automatically calculates a disk request based on the size of input files. If your jobs use additional temporary space beyond the inputs and begin failing, you may need to request extra space using the `-disk-scratch` option.
- Slurm users: By default, Slurm jobs receive no /scratch allocation. If your job runs entirely in memory and shared filesystems, this may be fine. However, many applications implicitly use temporary scratch space. If you encounter ENOSPC errors, you will need to request disk explicitly, for example:

```
#SBATCH --gres=disk:1G
```

If your job requires /scratch space, declaring it will now be mandatory. Please review your workflows ahead of time to ensure appropriate disk requests are being made.

Slurm Job (Outstanding jobs)

Running	Pending	Held	Other
7,474	28,348	0	0

Slurm Job (past 24 Hrs finished jobs)

Success	Failed	Cancelled	Timedout	OverMemory	NodeFail
52,572	8,967	2,122	756	2	0

Cluster Node Status

Datamover Status

File System Status

Job Info Last 24 Hrs

Documentation

- [SciComp Docs](#)
- [User's Guide \(old\)](#)
- [Data Policy](#)
- [Unrecoverable File](#)

Getting Started

Wed, 10/27/2021 - 09:36 — amitoj

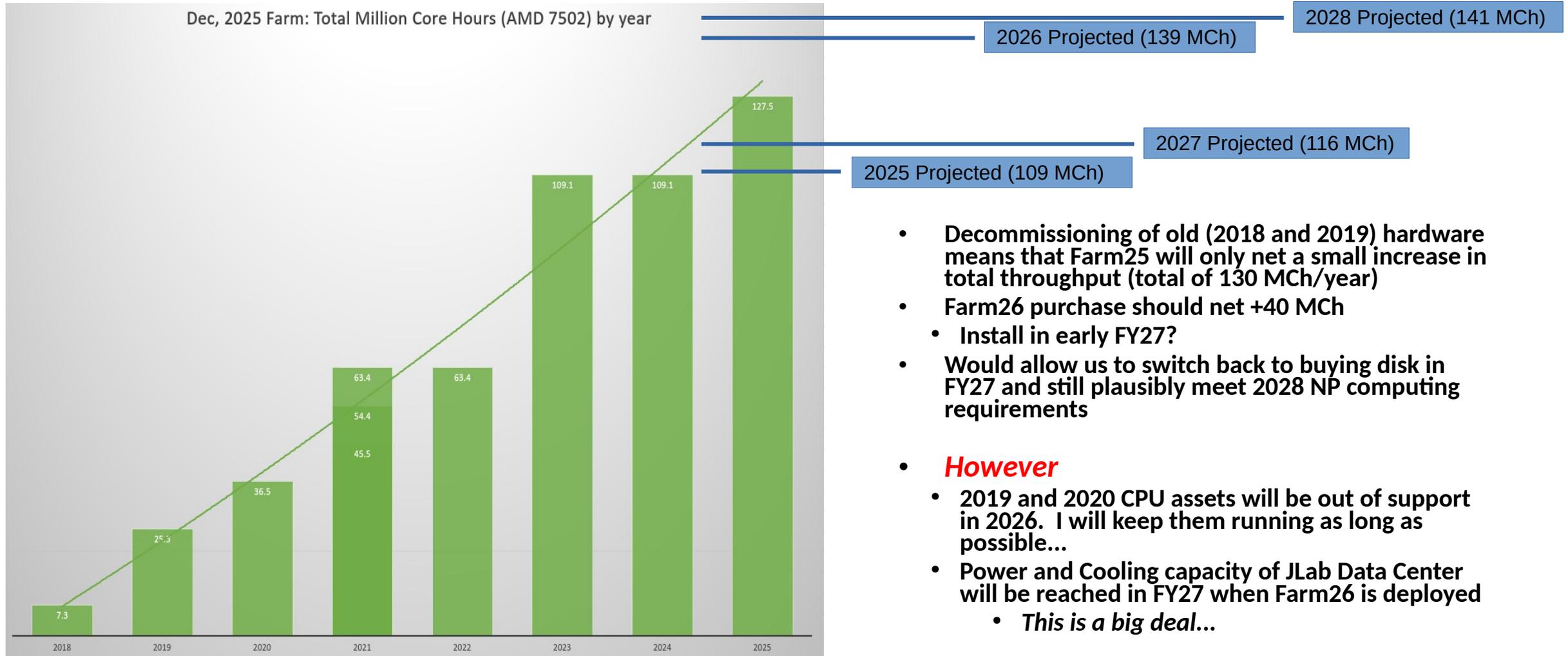


Steps to getting started using the Physics Farm system at Jefferson Lab are listed below:

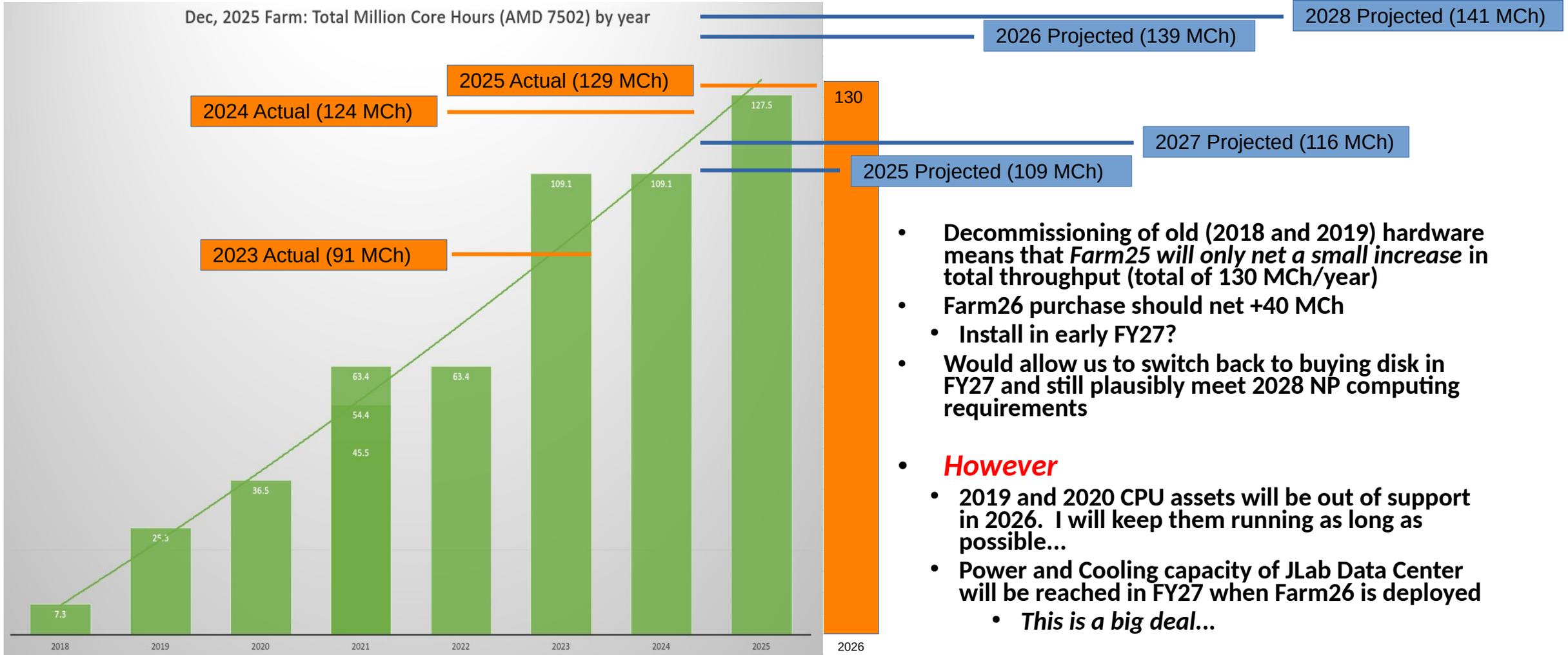
1. Obtain [Farm/ifarm access](#). Note that you will require a [2-factor authentication](#) token to login to the Farm and QCD nodes.
 2. [Generate a SciComp user certificate](#) for your account if you wish to interact with the JLab tape system and other infrastructure.
 3. Please [join the email list](#) `jlabs-scicomp-briefs@jlab.org` to get email about the status and updates of resources. This is **highly** recommended.
- Considerable information about the status of the computing resources is available at the [SciComp Portal: `scicomp.jlab.org/scicomp/`](#). On the entry page you can find information about the status of the various systems, as well as important announcements about new capabilities, current problems or planned outages. From the menu on the left you can get more detailed information about status and utilization, including reports by user or project for any arbitrary time interval.
5. For all other account questions please contact your [Hall Compute Coordinator](#), the Physics Computing Coordinator [Brad Sawatzky](#), or the [Computer Center Help Desk](#).
 6. New users of Farm resources (ifarm, farm, swif2, slurm, etc) are also encouraged to review the following resources (in addition to the documentation at the bottom of this page).

- [Basics of Farm Computing \(Tips & Tricks\)](#)
- [SWIF2 Farm Submission Example](#)
- [Recent Workshops \(with How-to's and Walkthroughs\)](#)
 - [GSPDA Computing Bootcamp \(May 2025\)](#)
 - [GSPDA Mini-Software Workshop \(Part 2 -- Sept 2024\)](#)
 - [GSPDA Mini-Software Workshop \(Part 1 -- May 2024\)](#)
- [JLab Knowledge Base Search](#)
- See also: [data.jlab.org](#)

JLab NP Computing Situation



JLab NP Computing Situation



- Decommissioning of old (2018 and 2019) hardware means that *Farm25 will only net a small increase in total throughput (total of 130 MCh/year)*
- Farm26 purchase should net +40 MCh
 - Install in early FY27?
- Would allow us to switch back to buying disk in FY27 and still plausibly meet 2028 NP computing requirements
- **However**
 - 2019 and 2020 CPU assets will be out of support in 2026. I will keep them running as long as possible...
 - Power and Cooling capacity of JLab Data Center will be reached in FY27 when Farm26 is deployed
 - *This is a big deal...*

JLab NP Data Center Power/Cooling Crisis

HPC DC Projected Growth (kW)--Updated on 8/14/25

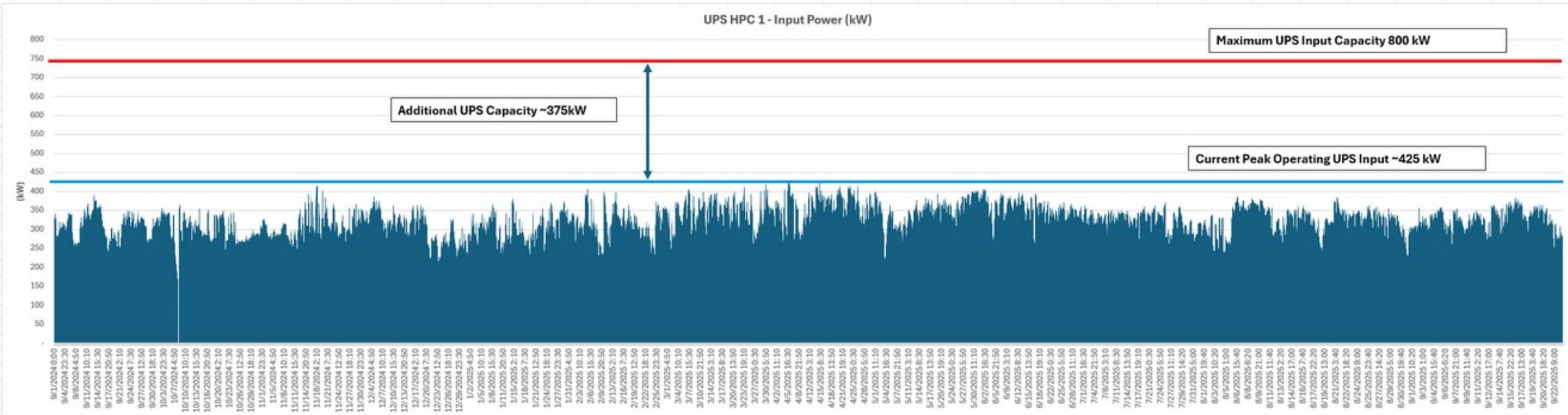
	FY23	FY24	FY25	FY26	FY27	FY28	FY29	FY30
Average Actual DC Load @ End of FY (Lower Bound Growth)	291	311	312	522	577	632	562	617
Max. Actual DC Load in FY (Upper Bound Growth)	355	404	418	628	683	738	668	723
Projected Growth During FY	89	175	50	310	55	55	105	55
-Farm CPU	48		0	100	0	50	0	50
-Farm GPU (AI/ML)	8	8	5	30	0	0	0	0
-HPDF	0	0	0	175	0	0	0	0
-LQCD	0	125	0	0	50	0	100	0
Ad hoc potential customer requirements (in kW load)	0	0	26	0	0	0	0	0
HACS 4 Testbed	33	42	19	5	5	5	5	5
Possible Decoms. During FY			192	100			175	

- Increases to projected Farm CPU and GPU loads **not** captured (~70kW ea. for CPU and ~40kW for GPU)
 - HPDF anticipated growth as of Jan 26 is closer to ~50kW (vs. 175 kW)
 - The increase in Farm power/cooling is offset by the lag in HPDF early-access projections*
- EIC Echelon-1 deployment timeline is unknown and not captured on this chart

JLab NP Data Center Cooling Capacity/Usage



JLab NP Data Center Power Usage/Capacity



Wrapping Up with Random Bits...

Helpdesk changes (from JLab Weekly announcement)

NEW HELP DESK HOURS & PROCESSES TO BEGIN - MARCH 9

- As the CST Division implements more efficient processes and procedures, the Help Desk will introduce new operating hours and support workflows beginning Monday, March 9, at 8 a.m.

- The new Help Desk hours will be Monday-Friday from:

- 8 a.m.-noon: Open for normal Help Desk operations, including walk-ups and urgent issues
- Noon-1 p.m.: Closed
- 1-4:30 p.m.: Open for appointments and urgent issues only

- To make an appointment or report an urgent issue, call x7155. For questions, contact the Help Desk.

- Calls will be routed to a live person during working hours (8am—5pm)

- Triaged and passed to a human immediately if urgent, otherwise a ticket will be created

- A kiosk has been set up at the helpdesk to ease ticket entry

Information Resources

- scicomp.jlab.org
 - [SciComp web page](#)
- [scicomp-briefs](#)
 - [mailing list for JLab Scientific Computing](#)

- [Documentation links](#)
 - [Getting Started](#)
 - [SciComp Knowledge Base](#)
 - [CST User Portal](#)
 - [JLab Helpdesk](#)

- » helpdesk@jlab.org
- » [Incident Request](#)

Jlab Scientific Computing

Welcome to the Jefferson Lab Scientific computing home page. [New users start here.](#)

Feb-27-24 Software Environment and Filesystem Changes The use of /apps is deprecated and is not available on farm AlmaLinux 9 machines. CVMFS is now used to distribute software. It is rooted under OASIS and can be used with [modulefiles](#) as before. For questions about software package availability, please submit a ServiceNow incident. For hall-specific software distribution questions, contact your computing coordinator. The legacy /site area has been removed. The path to Jasmine (tape) and cache tools is changed from /site/bin to /usr/local/bin. The CUE / u/scratch area has also been removed.

Feb-26-24 Farm Upgrade Schedule and Worker Node Selection The farm is being upgraded in a series of steps. Between now and June, the farm composition will change from majority CentOS 7 to predominantly AlmaLinux 9. At the time of this writing, CentOS 7 is the default. This default will change at a later step in the conversion process. Users may currently select which nodes run their jobs using slurm features/constraints. [This article](#) provides details on feature-based node selection. SWIF can pass features through to Slurm. See the SWIF introduction and [SWIF command line](#) reference for details. The interactive (farm) nodes currently run CentOS 7. A new machine, ifarm9.jlab.org is available for AlmaLinux 9 use now. Two new ifarm machines that will run AlmaLinux 9 are on order. They will replace the existing ifarm machines and include more per-core memory and temporary disk space.

Slurm Job (Outstanding jobs)					Slurm Job (past 24 Hrs finished jobs)						
Running	Pending	Held	Other		Success	Failed	Cancelled	Timeout	OverMemory	NodeFail	
5,427	12,010	1	5		45,322	13,485	66	222	35	19	

Cluster Node Status

farm16	farm18	farm19	farm23	scim2
~25	~75	~100	~25	~5

Datamover Status

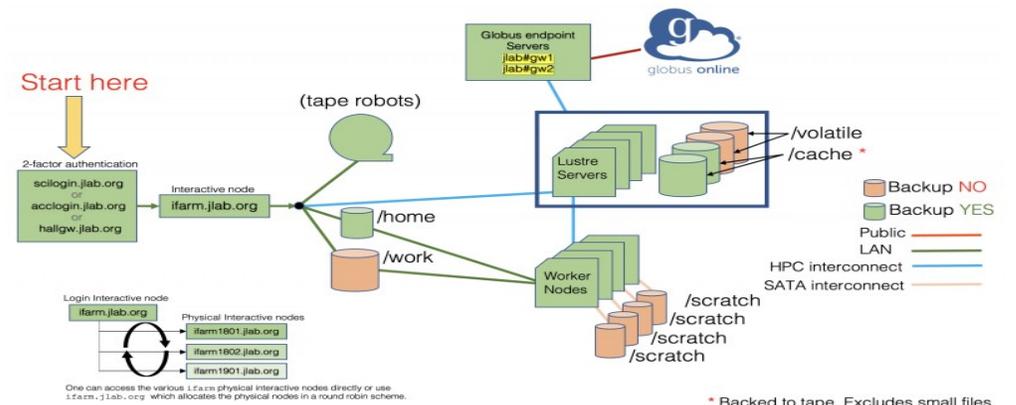
LTO7	LTO8
~5	~15

File System Status

Lustre	Cache	Volatile	work
~4k	~3k	~1k	~1k

Job Info Last 24 Hrs

Graph showing job counts over time for various farms.



JLab NP Computing Situation

Core Problem

- We do not have the compute necessary to meet your projected requirements
 - Driven by procurement and budget delays
 - Farm25 nodes sat on the shelf for 8 months waiting on Facilities, power, and procurement/budget...
 - Farm26 procurement has been 'on the bubble'
 - Seems we have the budget restored (for now), but AI boom has driven costs up a lot – I'm expecting long lead and less compute/\$.
 - New Disk is going to have to wait until 2027. It will also have to compete with a tape library upgrade to LTO10
 - OSG future is complicated – a loss of those resources would be a huge problem
 - Seems unlikely, but I expect the situation to get harder before it gets better...
 - ***Elephant in the room is Data Center cooling and power (next slides)***
- **Upshot is we must use the capacity we have efficiently**
 - Greater oversight and training on the part of Hall Compute Coordinators to support effective use of the Farm could help
 - Low efficiency jobs: <https://scicomp.jlab.org/scicomp/slurmJob/cpuEff>
 - Identify inefficient workflows / code
 - Worth looking into profiling and optimization for core code
 - User error / Insufficient testing before large campaign
 - Users submitting O(100k) jobs that burn allocation and then notice a problem