

How complex HPC workflows are driving the architecture of the NERSC-10 system



NERSC

Debbie Bard
Data Department Head, NERSC
March 13th, 2024

As the Mission HPC Center, NERSC is Highly Connected to the Office of Science

NERSC USERS ACROSS US AND WORLD

50

States,
Washington D.C.
& Puerto Rico

53

Countries

~10,000 Annual Users from ~800 Institutions + National Labs



32%

Graduate Students



19%

Postdoctoral Fellows



15%

Staff Scientists



13%

University Faculty



8%

Undergraduate Students



5%

Professional Staff



60%

Universities



29%

DOE Labs



5%

Other Government Labs



4%

Industry



1%

Small Businesses



<1%

Private Labs

SBIR

0.5%

NP

12.4%

HEP

16.7%

FES

13.9%

ASCR

2.9%

BER

15.2%

BES

38.3%

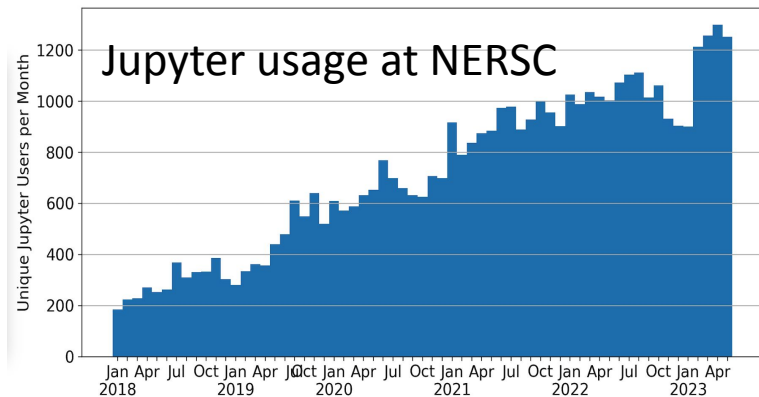


incl. 551 users from 34 HBCU+MSI

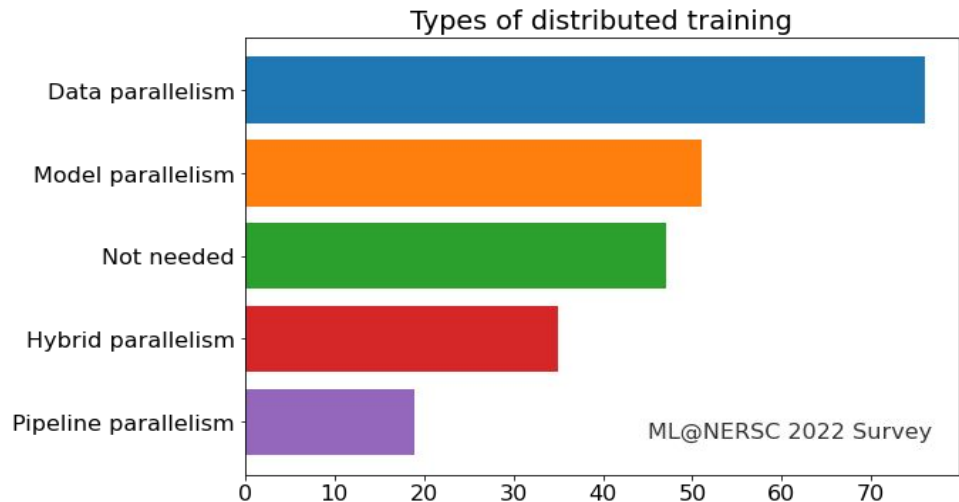
>80% of time at NERSC is allocated by DOE program managers

User community needs are diversifying

- Users interact with the system in new ways
 - > 2.5k Jupyter users - as popular as ssh
 - > 4.2k Python users - majority of active users
 - NERSC's Top 500 result run in Shifter container
 - Superfacility API: 1 request logged every 2 sec

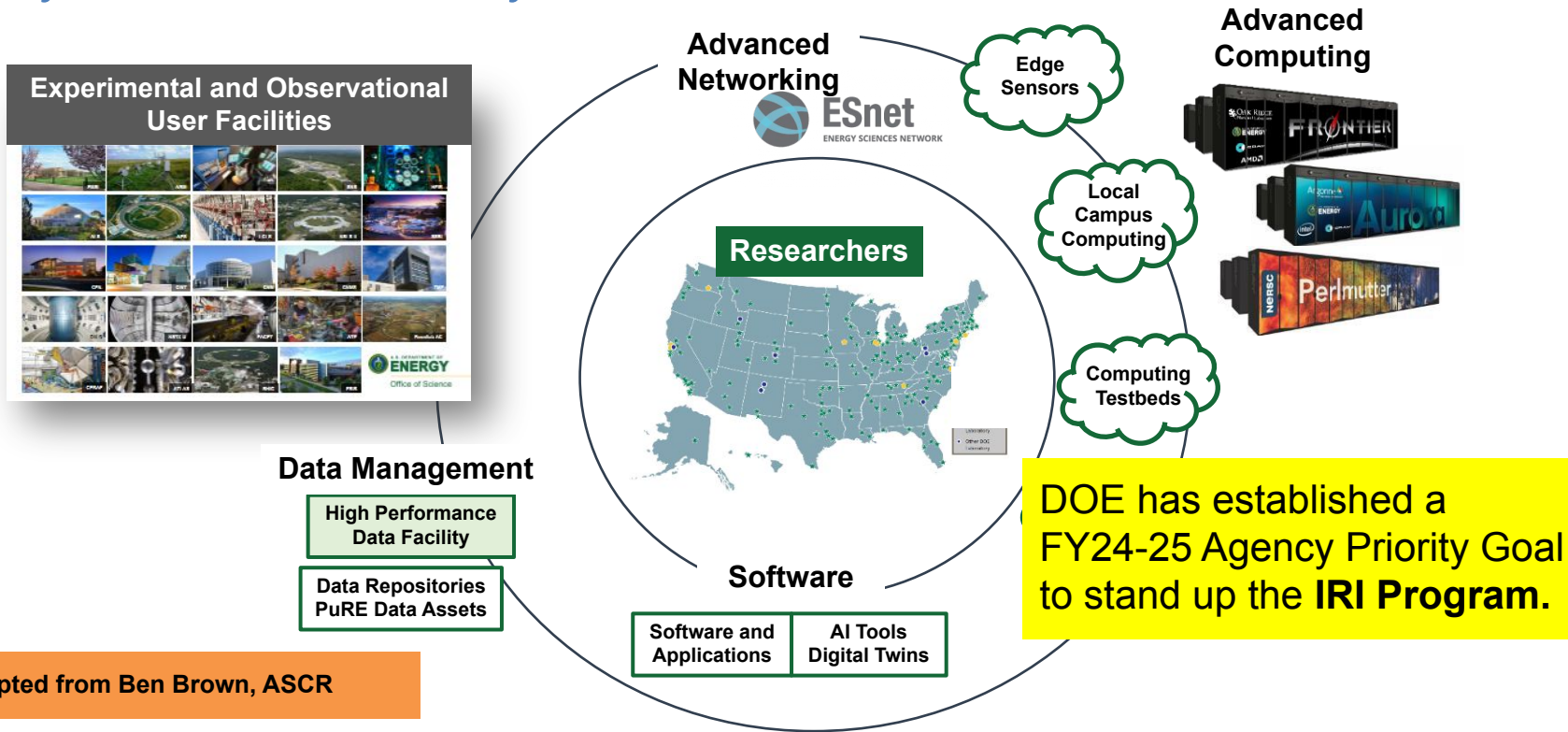


- Demand and capability for HPC-scale AI is increasing
 - >20x increase in AI users from 2017
 - Multiple AI applications run at full scale on Perlmutter



DOE's Integrated Research Infrastructure (IRI) Vision:

To empower researchers to meld DOE's world-class research tools, infrastructure, and user facilities seamlessly and securely in novel ways to radically accelerate discovery and innovation

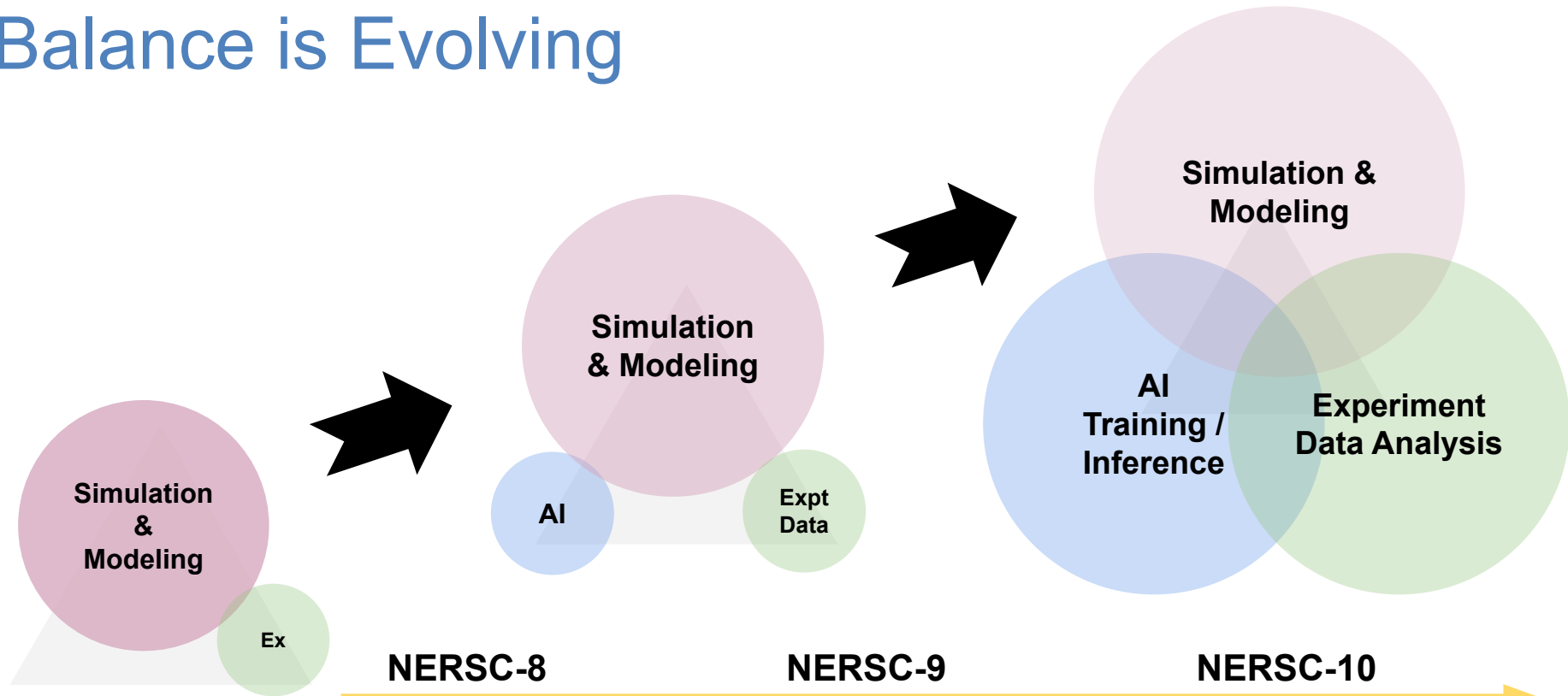


IRI will enable seamless workflows via close collaboration between ASCR facilities & the DOE Scientific Community

ASCR is implementing IRI through these major elements

- 
- 1 Invest in IRI foundational infrastructure
 - 2 Stand up the IRI Program governance and FY24 workstreams
 - 3 Bring IRI projects into formal coordination
 - 4 Deploy an IRI Pathfinding Testbed across the four ASCR Facilities

The HPC Facility Workload Balance is Evolving

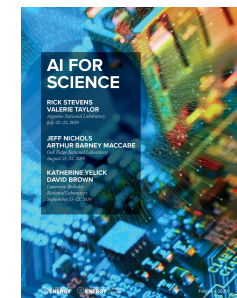
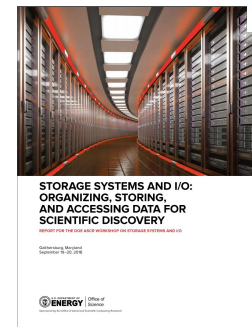
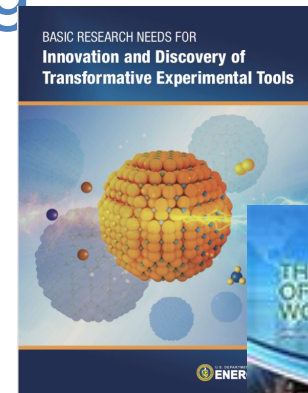


N10 User Requirements are Evolving

Users require support for new paradigms for data analysis with **real-time interactive feedback between experiments and simulations.**

Users need the ability to search, analyze, reuse, and combine data from different sources into **large scale simulations and AI models.**

NERSC-10 Mission Need Statement:
*The NERSC-10 system will **accelerate end-to-end DOE SC workflows** and enable new modes of scientific discovery through the integration of experiment, data analysis, and simulation.*



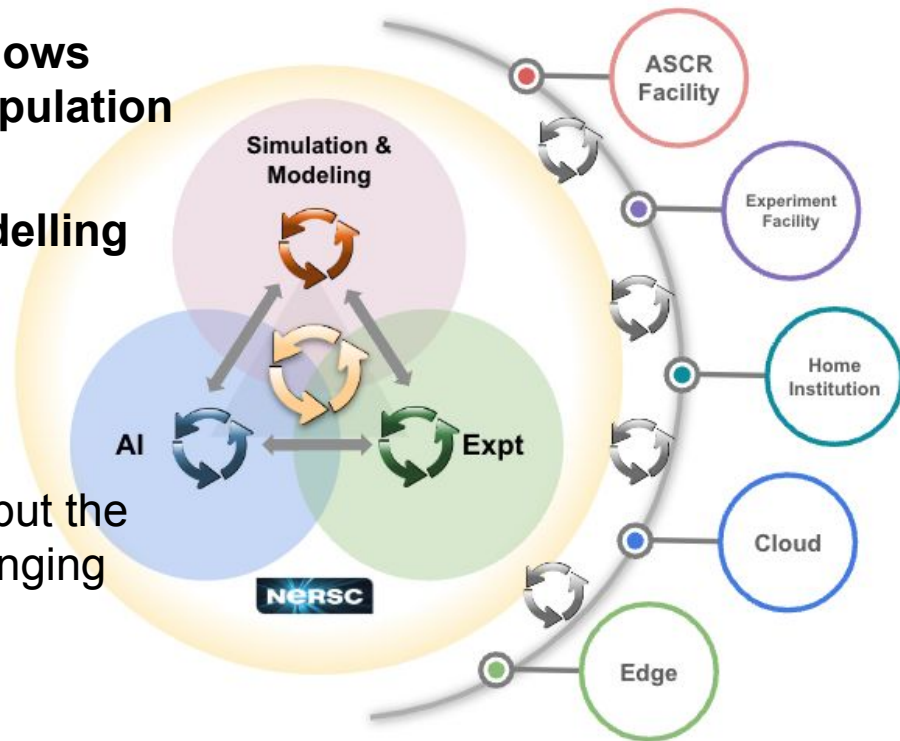
What is an HPC Workflow?

Workflows are interconnected computational and dataflow tasks with data products. They have task coupling (control flow) and/or data movement between tasks (data flow).

High performance computing (HPC) workflows interconnect computational and data manipulation steps across one/some/all of:

- High performance simulation and modelling
- High performance AI workflows
- High performance data analytics

We've been running workflows for decades - but the complexity and timeliness of workflows is changing which motivates a new approach with N10.



We identified 6 workflows archetypes to help define our vision for N10

1. High-performance simulation & modeling workflow	large-scale multi-physics applications with checkpoint/restart, data post-processing, visualization
2. High-performance AI (HPAI) workflow	data integration-intensive science patterns such as training, inference, hyperparameter optimization
3. Cross-facility workflow: Rapid data analysis and real time steering	time-sensitive science patterns such as superfacility, edge, and hybrid cloud
4. Hybrid HPC-HPAI-HPDA workflow	long-term campaign science patterns, AI-in-the-loop, AI-around-the-loop
5. Scientific data lifecycle workflow: Interactive, data-analytics and viz	data integration-intensive science patterns such as Jupyter, scientific databases, VSCode
6. External event-triggered and API-driven workflow	time-sensitive science patterns such as function-as-a-service, microservices

We identified 6 workflows archetypes to help define our vision for N10

1. High-performance simulation & modeling workflow	large-scale multi-physics applications with visualization
2. High-performance	s such as ization
3. Cross-facility analysis and	superfacility,
4. Hybrid HP	n-the-loop,
5. Scientific data Interactive, data-analytics and viz	search for “NERSC workflows white paper” erns such as Jupyter, scientific databases, VSCode
6. External event-triggered and API-driven workflow	time-sensitive science patterns such as function-as-a-service, microservices

HPC Workflows Drive N10 Technology Capabilities

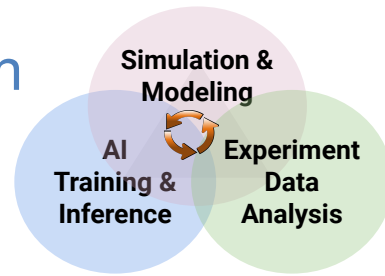
	Cloud native/ containers	QoS storage system (QSS)	End-to-end API	Network/ scheduling QoS	IRI/ Multi-site workflows	Smart networking	Prog. Env	Workflow Enablement Nodes (WEN, fka Spin)
1.Simulation & modeling		X	X			X	X	
2.AI	X	X	X	X	X	X	X	X
3.Cross-facility	X	X	X	X	X	X		X
4.Hybrid HPC-HPAI-HPDA	X	X	X	X	X	X	X	X
5.Scientific data lifecycle	X	X	X	X			X	X
6.Event-triggered & API-driven	X	X	X	X		X	X	X

HPC Workflows Drive N10 Technology Capabilities

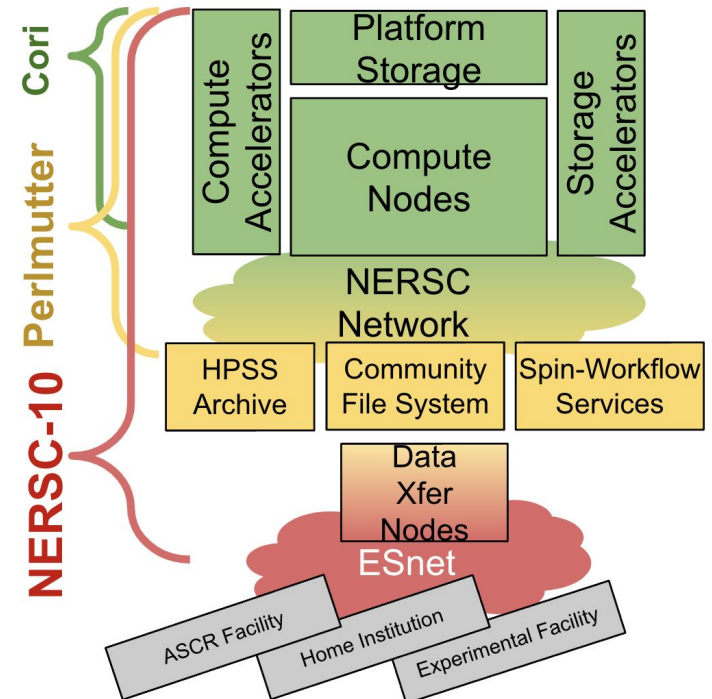
	Cloud native/ containers	QoS storage system (QSS)	End-to-end API	Network/scheduling QoS	IRI/ Multi-site workflows	Smart networking	Prog. Env	Workflow Enablement Nodes (WEN, fka Spin)
1.Simulation & modeling		X	X			X	X	
2.AI	X	X	X	X	X	X	X	X
3.Cross-facility	X	X	X	X	X	X		X
4.Hybrid HPC-HPAI-HPDA	X	X	X	X	X	X	X	X
5.Scientific data lifecycle	X	X	X	X			X	X
6.Event-triggered & API-driven	X	X	X	X		X	X	X

Pink: cannot be done today
Orange: can be done only with extraordinary effort
Green: can be done today in limited way

NERSC-10 is Designed to Support Complex Simulation and Data Analysis Workflows at High Performance



- **Quality of Service:** computation, storage and networking enables response-time plus utilization.
- **Seamlessness:** tight integration of system components enables high performance workflows.
- **Programmability:** APIs manage data, execute code, and interact with system resources.
- **Orchestration:** coordinates resource management across domains.
- **Portability:** Modular workflow execution across IRI sites.
- **Security:** authentication, authorization and auditing.



Innovation in software is key to enabling complex workflows

New capabilities:
FaaS/serverless,
specialized HW, AI
deployment, data
lifecycle, quantum...

Support usage of both
ssh and Jupyter

Meet federal security
requirements

User software/ workflows

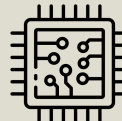


Workflow environment

System software



System hardware



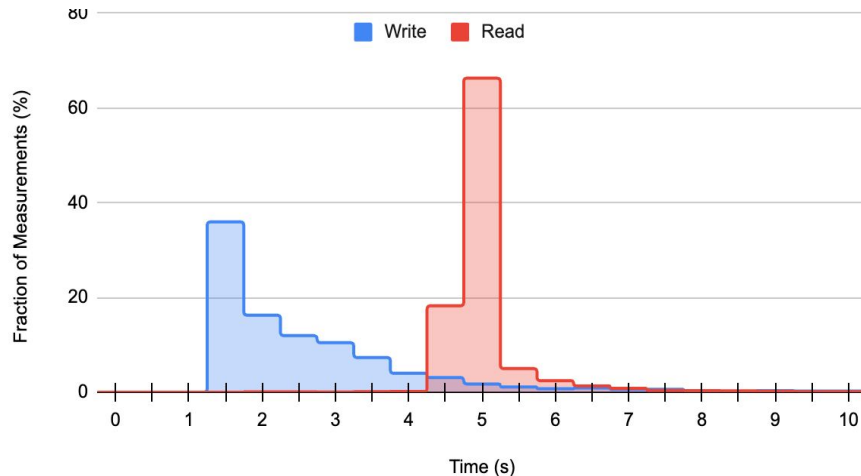
RESTful user-facing
APIs support
automation

System-side APIs for
workflow observability,
administration and
reconfigurability

Containerize the user
environment

The NERSC workload requires capabilities that are hard to reconcile in a single file system

IOR performance on Perlmutter



- 21% of all write tests took more than twice as long as the mode (1.5 sec)
- 2% of all write tests took at least **five times longer** than the mode

For instrument-driven and time-dependent workflows such variance could be catastrophic

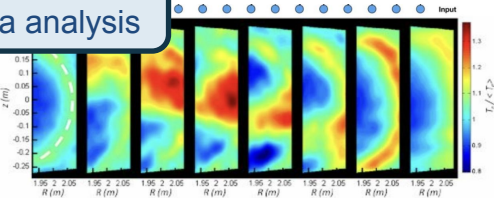
- **Quality of Service Storage System (QSS)** will provide controllable, guaranteed IOPs / bandwidth to meet the needs of time-sensitive workflows
- Platform Storage System (PSS) is a more traditional FS that will meet the needs of much of the NERSC workload

Cross-facility workflow example: Fusion science with DIII-D, preparing for ITER

Data readout from tokamak sensors sent to NERSC



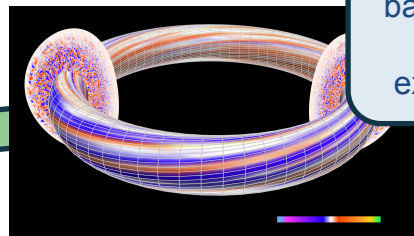
AI-driven data analysis



Feedback to scientist in minutes

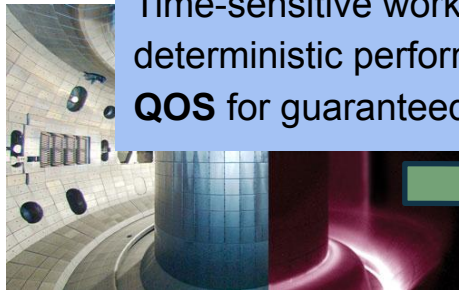


Digital Twin simulation based on analyzed data: recommends new experiment parameters



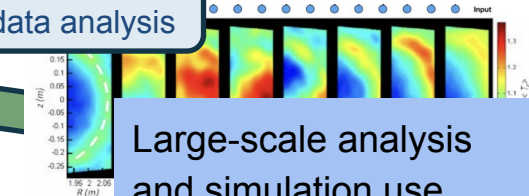
Cross-facility workflow example: Fusion science with DIII-D, preparing for ITER

Time-sensitive workflow requires **QSS** for deterministic performance and **network QOS** for guaranteed response in $O(\text{min})$



Data movement and compute progress tracked using **APIs** by automated workflow orchestrator

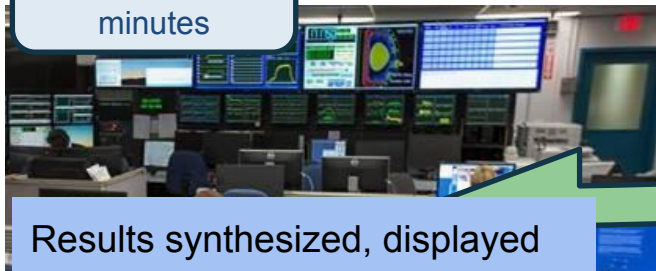
AI-driven data analysis



Large-scale analysis and simulation use **containerized apps** and **accelerated nodes**.

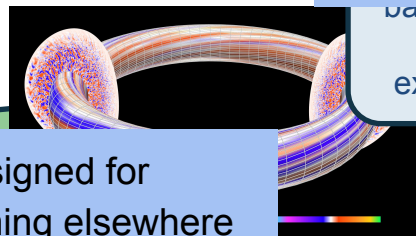
Based on analyzed data, recommends new experiment parameters

Feedback to scientist in minutes



Results synthesized, displayed and shared via **Jupyter** and **python** ready for the next shot

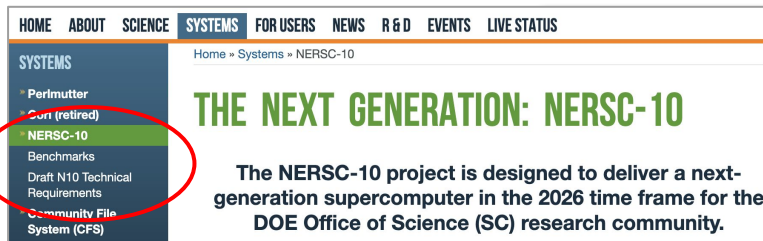
Portable workflows designed for resiliency, possibly running elsewhere if NERSC is unavailable (IRI)



NERSC-10 RFP: Technical Requirements

Technical Summary:

- No peak flops requirement
 - 10x on workflow component benchmarks
- CPU + GPU nodes
- Two kinds of storage
 - PSS - 120 PB, 20 TB/s
 - QSS - 80 PB, performance guarantees
- Workflow Environment (beyond the programming environment)
- Modular system software and management to support complex workflows



HOME ABOUT SCIENCE SYSTEMS FOR USERS NEWS R & D EVENTS LIVE STATUS

SYSTEMS

- Perlmutter
- Curt (retired)
- NERSC-10**
- Benchmarks
- Draft N10 Technical Requirements
- Community File System (CFS)

Home » Systems » NERSC-10

THE NEXT GENERATION: NERSC-10

The NERSC-10 project is designed to deliver a next-generation supercomputer in the 2026 time frame for the DOE Office of Science (SC) research community.

search "nersc rfp"

September 15, 2023

RFP

Technical Requirements
Document

for

NERSC-10 System

Version 3.0

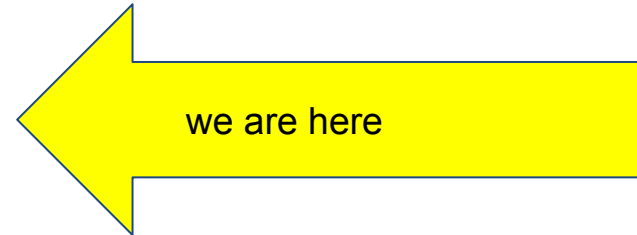
Lawrence Berkeley National Laboratory is operated by the University of California for the U.S. Department of Energy under contract NO. DE-AC02-05CH11231.

1

RFP Technical Requirements Document for NERSC-10 System, Version 3.0, September 15, 2023

NERSC-10 Timeline

- Project Authorized by DOE (CD-0) - Sept 2021
- Advanced Acquisition Plan approved by DOE - March 2023
- **Draft RFP Release - 20 April 2023**
- Technical Design Review - August 2023
- Berkeley Lab Director's Review (Red Team) - Fall 2023
- **CD-1 - December 2023**
- **RFP Release - March 2024**
- **Contract signed (CD-2) - Late CY 2024**
- (Potential) Phase I or Pilot System- mid 2025
- Technical Decision Point - Late 2025
- **Main System Delivery - Late 2026**
- **User access 2027**

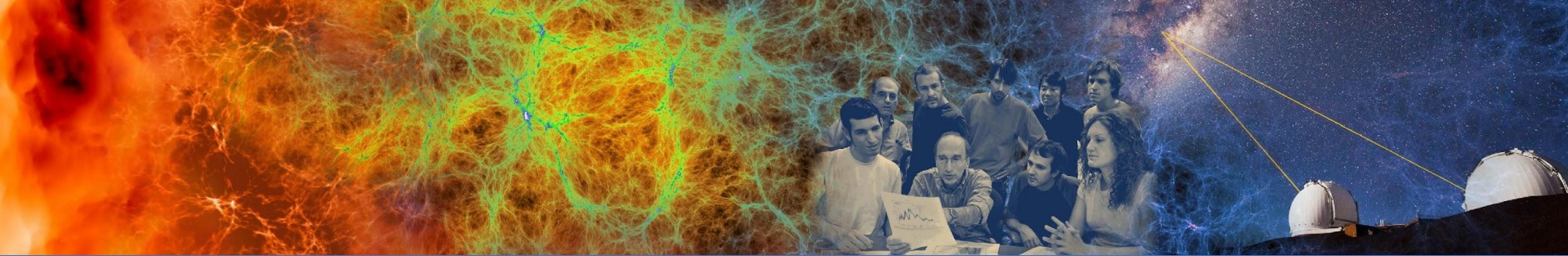


The NERSC-10 system will accelerate end-to-end DOE SC workflows and enable new modes of scientific discovery through the integration of simulation, data analysis and experiment.

Our technology choices for NERSC-10 are informed by the work we've done over the past 5 years to develop, operationalize and support Perlmutter and our users - including lessons learned from the Superfacility project and IRI.

We're building an engagement model to coordinate a complex set of requirements and stakeholders in a changing technology landscape.

- *N10 will deliver 10x Perlmutter performance on HPC workflows.*
- *N10 is designed to be IRI-ready.*
- *GPU-enabled applications should have minimal issues in porting/running their applications.*
- *The N10 RFP will be released any moment now, with system delivery in 2026.*

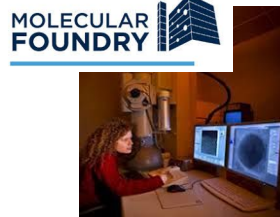
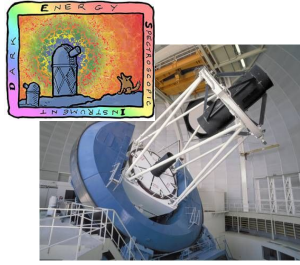
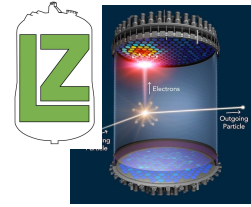
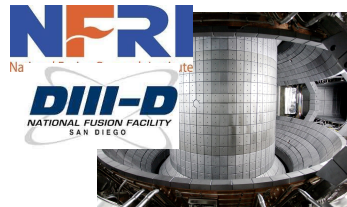


Thanks!

Multiple science teams are using NERSC for superfacility-enabled science, in production

The 3 year Superfacility project kick-started this work, building the base infrastructure and services. We now support **multiple science teams using automated pipelines to analyze data from remote facilities at large scale**, without routine human intervention, using:

- **Real-time** computing support
- Dynamic, high-performance **networking**
- Data management and movement tools, incl. **Globus**
- **API-driven** automation
- HPC-scale notebooks via **Jupyter**
- Authentication using **Federated Identity**
- Container-based edge services supported via **Spin**



Multiple science teams are using NERSC for superfacility-enabled science, in production

A set of 8 initial close science engagements drove this work, but the impact has scaled to benefit all NERSC users

- **Real-time** computing support
- Dynamic, high-performance **networking**
- Data management+movement tools, incl. **Globus**
- Interactive HPC via **Jupyter**
- Container-based edge services supported via **Spin**
- **API** interfaces
- **Federated Identity**/auth
- Collaboration accounts for **automated “robot” access**

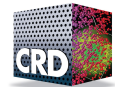
>20 science teams use the **realtime** qos to process urgent data

>1500 unique **Jupyter** users per month, similar to number of users who ssh into our systems

>250 users, >85 projects use **Spin**

>40 projects use the NERSC **API**, ~19M logged requests since May 2022 = one request every 2 sec

>1400 users are now logging in with a home lab identity



AMCR
SciData

