



Strategic Plan

FY 2024 – 2034

National Energy Research
Scientific Computing Center



BERKELEY LAB



U.S. DEPARTMENT OF
ENERGY

Office of
Science

Lawrence Berkeley National Laboratory | 1 Cyclotron Road, Berkeley, CA 94720-8148

Table of Contents

| | |
|--|-----------|
| 1. Executive Summary | 3 |
| <hr/> | |
| 2. Overview | 4 |
| 2.1 Facility Overview | 4 |
| 2.2 The Current HPC Landscape | 4 |
| 2.2.1 Future Users and Workloads | 5 |
| 2.3 Crosscutting Themes for the Next 10 Years | 6 |
| 2.3.1 Breadth and Depth of Science Impact | 6 |
| 2.3.2 Collaboration | 7 |
| 2.3.3 Operational Excellence and Innovation | 8 |
| <hr/> | |
| 3. Focus Area: Impact on Science Through Partnerships with Users | 9 |
| 3.1 Integrated Research Infrastructure and Preparing Users for Future Systems | 9 |
| 3.2 Pervasive AI for Science | 11 |
| 3.3 User Engagement | 13 |
| <hr/> | |
| 4. Focus Area: System and Data Center Architecture | 15 |
| 4.1 Future Systems | 15 |
| 4.1.1 NERSC-II and Beyond | 15 |
| 4.1.2 Quantum Computing | 16 |
| 4.2 Software Ecosystem | 18 |
| 4.3 Data Center Design and the NERSC Center Roadmap | 20 |
| <hr/> | |
| 5. Focus Area: Smart Green Facility | 21 |
| 5.1 Sustainability | 21 |
| 5.2 Smart Facility | 22 |
| 5.3 Security in an Open Science Environment | 23 |
| <hr/> | |
| 6. Focus Area: Workforce Development | 25 |
| 6.1 A Catalyst for Workforce Development | 25 |
| 6.2 Creating a Hybrid Work Environment for the Future | 26 |
| <hr/> | |
| References | 28 |

1. Executive Summary

The National Energy Research Scientific Computing Center (NERSC) is the U.S. Department of Energy (DOE)'s mission high performance computing (HPC) facility for the Office of Science (SC). NERSC serves the largest and most diverse research community of any HPC facility in the DOE complex, and one of the largest in the world. NERSC has over 10,000 users and provides services to unclassified research programs in high-energy physics, biological and environmental sciences, basic energy sciences, nuclear physics, fusion energy sciences, mathematics, and computational and computer science.

NERSC's mission is to accelerate scientific discovery at the DOE SC through high performance computing and data analysis. Our vision is to enhance the productivity of our users by enabling workflows incorporating simulation, pervasive AI, quantum computing, and data analysis to operate seamlessly in the DOE Advanced Scientific Computing Research (ASCR) Integrated Research Infrastructure.

NERSC is driven by a set of core values: maximizing the breadth and depth of science impact through collaboration with users, industry, ASCR facilities, DOE, and the HPC community, and delivering operational excellence.

NERSC has been accelerating scientific achievement for fifty years – and is laying the groundwork for the next fifty. For the next ten years, NERSC's strategic areas of focus are:

- Maximizing impact on science through partnerships with users, developers, and industry
- Developing and deploying first-of-a-kind systems and data center architectures
- Innovating the data center to become a world-leading, next-generation smart green facility
- Continuing to serve as a catalyst for HPC workforce development

The following describes these actions in more detail and outlines our plan to put them into action.

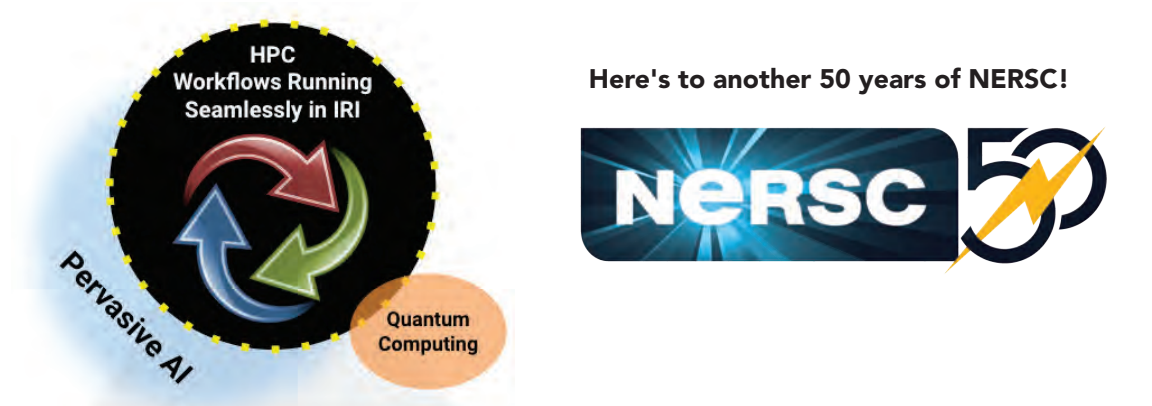


Figure 1-1. The future of NERSC is one of seamless connection through integrated workflows, supplemented by pervasive AI and quantum computing.

2. Overview

2.1 Facility Overview

The National Energy Research Scientific Computing Center (NERSC) is the U.S. Department of Energy (DOE)'s mission high performance computing (HPC) facility for the Office of Science (SC). Founded at Lawrence Livermore National Laboratory in 1974 and located at Lawrence Berkeley National Laboratory (Berkeley Lab) since 1996, NERSC has been accelerating scientific achievement for fifty years—and is laying the groundwork for the next fifty.

NERSC serves the largest and most diverse research community of any HPC facility in the DOE complex, and one of the largest in the world. NERSC has over 10,000 users and provides services to unclassified research programs in high-energy physics, biological and environmental sciences, basic energy sciences, nuclear physics, fusion energy sciences, mathematics, and computational and computer science. DOE SC program managers directly allocate more than 90% of the available compute time at NERSC.

NERSC collaborates with its user community to accelerate scientific achievement. We architect, deploy, and operate large-scale, state-of-the-art computing, data storage systems, and edge services to users at universities and national laboratories and in industry. NERSC also offers an extensive and comprehensive set of services, including substantial user documentation, a broad training program, and access to its staff's expertise.

Key to the mission of NERSC is enabling computational science at scale, in which large, interdisciplinary teams of scientists attack fundamental problems in science and engineering that require massive calculations and have broad scientific and economic impacts. In addition to deploying first-of-a-kind HPC systems, NERSC pursues deep engagements with SC users, developers, and industry partners (such as NESAP, the NERSC Science Acceleration Program) to help users cross the chasm of disruptive technologies and, eventually, scale the impact of these technologies to thousands of users through world-leading support, training, and documentation.

2.2 The Current HPC Landscape

The HPC landscape has changed significantly in the last ten years, and the rate of change is likely to increase even more rapidly over the next ten. As described in the [2024 ASCAC Facilities Subcommittee report](#) [1], science and engineering research is addressing increasingly complex questions that require a deeper integration of disciplines and methodologies, and accordingly need significant and growing computational and data analysis resources. There are numerous challenges we anticipate NERSC having to respond to:

- The NERSC workload is growing larger and more diverse, with an increasing emphasis on integrated research workflows that require time-sensitive transport and analysis of unprecedented quantities of data.
- User workloads are shifting in such a way that the three main capability thrusts of supercomputing—simulation and modeling, data analysis, and AI training and inference—form a tightly integrated workflow. This increases the complexity of users' needs and the systems required to meet them.
- AI has become pervasive, transforming science through innovations like automation, fast experiment design, unsupervised detection of novel science, inference with full science models, AI knowledge discovery assistants, and extreme-scale surrogate models.

- As Moore’s Law comes to an end, disruptive technologies such as quantum computing will become more prominent, and there will be increased focus on energy-efficient, and potentially specialized, solutions.
- Industry AI trends will increasingly drive hardware design choices, offering the opportunity to leverage them to meet the varied and complex needs of NERSC users.

2.2.1 Future Users and Workloads

Today, in addition to the SC community’s insatiable demand for increased computing power in the twilight of Moore’s Law, NERSC must enable scientists to leverage the immense potential of new synergies by integrating simulation, data, and AI (in multiple fields of research) in complex workflows. The need for enhanced capabilities supporting new workflow-centric computing paradigms is highlighted in analysis of the NERSC workload and allocation requests, and through deep engagements with users across initiatives like the [Superfacility Project](#) and the [Integrated Research Infrastructure \(IRI\) Architecture Blueprint Activity](#).

IRI will tie together experimental and observational facilities and computing centers through DOE’s ESnet network. As IRI enables increasingly collaborative research workflows and larger teams gain access to HPC resources, the number of NERSC users is expected to double—or potentially triple—in the next ten years.

These users require new workflow capabilities, including real-time interactive feedback between experiments and simulations; persistent services supporting on-demand data analysis; and new tools to search, analyze, reuse, and combine data from different sources into large-scale simulations and AI models. They need not only more powerful supercomputers, but ones that are more dynamic, feature-rich, and programmable to enable more complex end-to-end workflows.

Looking toward the end of the decade, NERSC anticipates a convergence of three trends that will influence the design and deployment of NERSC-II and the surrounding data center ecosystem.

- The IRI ecosystem will be fully realized, greatly increasing the connectivity of resources across DOE including full interactivity with HPC workflows running at scale, massively connected to clients (users, sensors, or experiments) potentially numbering in the millions.
- AI will be pervasive in nearly every aspect of the scientific process, from how users interface and interact with systems, to the application of foundation models across multiple domains, to embedding AI throughout end-to-end workflows, and to automated and self-driving compute systems and experimental facilities.
- The post-Moore’s Law era will impact both the design of classical computing systems and, potentially, lead to the broader availability of new energy-efficient technologies such as quantum, neuromorphic, and optical computing.

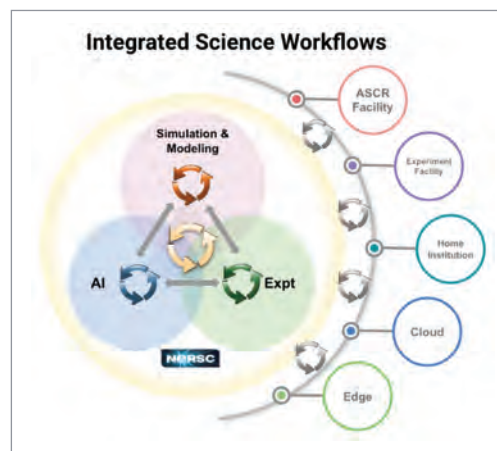
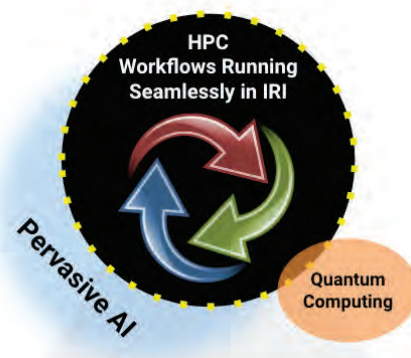


Figure 2-1. Increasingly, NERSC must enable scientists to leverage the immense potential of new synergies by integrating simulation, data, and AI (in multiple fields of research) in complex workflows.



These increased capabilities will come with added complexity for NERSC and its users. This highlights the importance of our focus on enabling scientific discovery by deeply engaging with our user community.

2.3 Crosscutting Themes for the Next 10 Years

This report describes the NERSC strategy roadmap for the next decade, which will enable us to navigate the changing landscape and enable researchers to drive new scientific discovery and leverage opportunities in technology advancements. NERSC identified four focus areas to prioritize the use of limited resources and create actionable goals. The four focus areas are:

- Impact on science through partnership with users
- System and data center architecture
- Smart green facility
- Workforce development

Across these focus areas, three fundamental cross-cutting themes carry NERSC's core values. These cross-cutting themes include:

- Breadth and depth of science impact
- Collaboration (with users, industry, ASCR facilities, DOE and the broad HPC community)
- Operational excellence and innovation

The following four sections describe the key focus areas for NERSC's 2024–2034 Strategic Plan and provide short- and long-term goals incorporating the cross-cutting themes.

2.3.1 Breadth and Depth of Science Impact

One way of measuring success at NERSC is breadth and depth of scientific impact. NERSC resources and expertise have been used to support seven Nobel Prize-winning scientific achievements, and its role is acknowledged in over 2,000 peer-reviewed scientific publications yearly. Since 2020, NERSC has been referenced in 8,790 refereed journal articles, including in science's most prominent publications. This vast number of citations and the quality of the results NERSC enables makes it an extremely productive computing center.

Due to the size and diversity of its user community and the broad spectrum of scientific work it supports, NERSC will need to accommodate an increasingly wide range of workflows, codes, and algorithms to maintain that breadth and depth. In the coming decade, considerable focus will be on scaling impact to that broad user community.

To build a tech-savvy scientific workforce for the future, NERSC is actively building its programs and services aimed at developing users' skills. In addition to a full slate of training events and workshops, NERSC is improving its already-robust documentation to include a new online training platform and documentation specifically aimed at orienting new users. Further improvements are in the pipeline to ensure that future users are prepared to work efficiently in the changing scientific-computing landscape.

2.3.2 Collaboration

In addition to providing computing and storage resources, NERSC also helps prepare users and their work for new technologies and workflows. NERSC has laid the groundwork for addressing these needs in the form of the NERSC Scientific Acceleration Program (NESAP), a collaborative effort with code teams, vendors, and library and tools developers to prepare for advanced architectures and new systems, as well as other initiatives. As we look toward NERSC-10 and the future, NESAP is evolving to focus on holistic end-to-end workflows including new aspects of performance optimization which include data-transfer and advanced workflow capabilities like the use of quality-of-service guaranteed storage.

NERSC has supported experiment science teams for decades, but the past ten years have shown that the needs of user teams are evolving as their workflows become more complex. The [Superfacility Project](#) (2019–2022) kick-started NERSC’s work to support these evolving needs, co-developing tools and services like Spin, Jupyter, Federated ID and the API. This work was performed in close collaboration with a set of eight science partners (shown below) from a broad range of science areas who helped steer the development process.

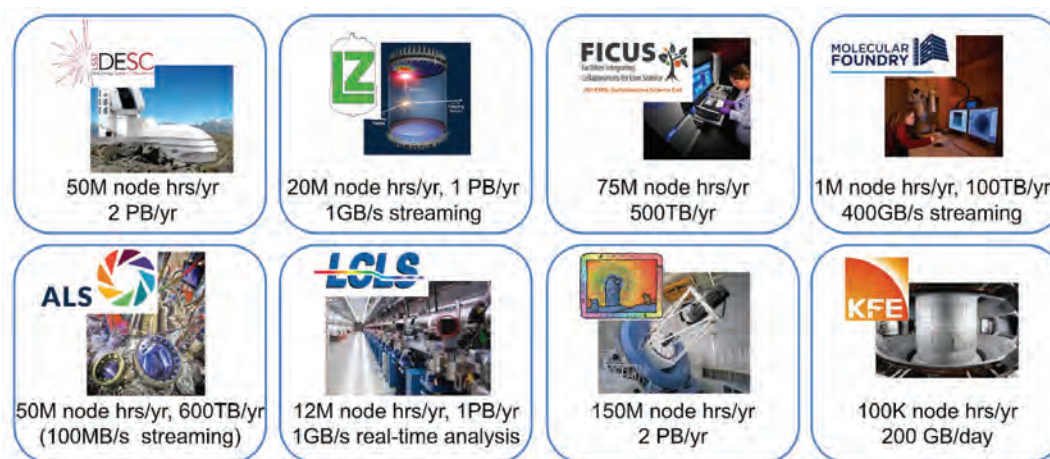


Figure 2-2. Eight [science teams](#) from across the scientific landscape partnered with NERSC to develop tools and services that support complex integrated scientific workflows.

In the coming years NERSC will build on and expand these collaborations in partnership with the ASCR community and partner ASCR facilities to support the DOE’s broader IRI initiative, making it easier for scientists to access and leverage all DOE resources. NERSC’s experience directing the Superfacility Project will be an invaluable resource for that work.

Collaboration with HPC vendors has been a part of NERSC’s strategy for many years, and will continue in the next decade. In an era where the vendor community is increasingly focused upon generative AI, it will be essential to provide quantitative data about NERSC’s workload and engage in co-design to ensure that the needs of NERSC’s users are addressed, especially in situations where they diverge from the needs of AI.

In AI, NERSC has participated in numerous collaborations over the last ten years, including with academia, industry, and domain scientists. This includes collaborations for optimizing AI software on NERSC systems (such as with NVIDIA and HPE, developing the NVIDIA Collective Communications Library [NCCL] for Slingshot II); benchmarking (such as for [MLPerf HPC](#)); development of AI tools and technologies; application of novel AI methods into science workflows (numerous [recent projects](#)); and training of the wider community (e.g. for the Deep Learning at Scale tutorial from SCI18 to [SC23](#)).

Moving forward, we plan to build on this foundation to enter an era of pervasive AI: increasingly, wherever AI can be used to accelerate science and HPC, NERSC will enable that capability. This will require activity across the systems, software, and applications at NERSC, in addition to the ways in which users interact with the center. Through cross-cutting activity described in the sections below, AI will become fully integrated with scientific simulations and data pipelines as well as the heart of NERSC operations, including systems, workflow orchestration, user support, and security.

2.3.3 Operational Excellence and Innovation

As a DOE National User Facility [2], NERSC is committed to operational excellence to maximize scientific impact and productivity. NERSC staff work tirelessly to deliver and maintain state-of-the-art supercomputers, large data storage systems, and edge services that enable and advance cutting-edge scientific research. NERSC ensures high resource availability and incorporates best practices to minimize downtimes.

NERSC also offers a range of services, including thorough and dynamic user documentation, an active training program, and access to its staff's expertise. An annual survey allows staff to collect feedback on user satisfaction with NERSC's services and HPC resources, helping to maintain operational excellence. And to align with the changing workload and demands of the growing user community, NERSC staff are continuously implementing innovations to improve the facility's operations across technical, management, and workforce development dimensions.

3. Focus Area: Impact on Science Through Partnerships with Users

3.1 Integrated Research Infrastructure and Preparing Users for Future Systems

The future of scientific research is one of collaboration and integration, utilizing the power of team science and seamless technology to accelerate the pace of innovation. NERSC is already leading in these areas and is continuing to prepare for this integrated future.

The Superfacility model at Berkeley Lab connects experiment and compute facilities with the expertise and community they need for success. Close collaboration between NERSC, ESnet, and the Berkeley Lab Computing Sciences research divisions supports multiple DOE science teams using automated pipelines to analyze data from remote facilities at large scale (see Figure I-1). These services have since been deployed to the full NERSC user base, where they are widely used, inside and outside the experiment science community.

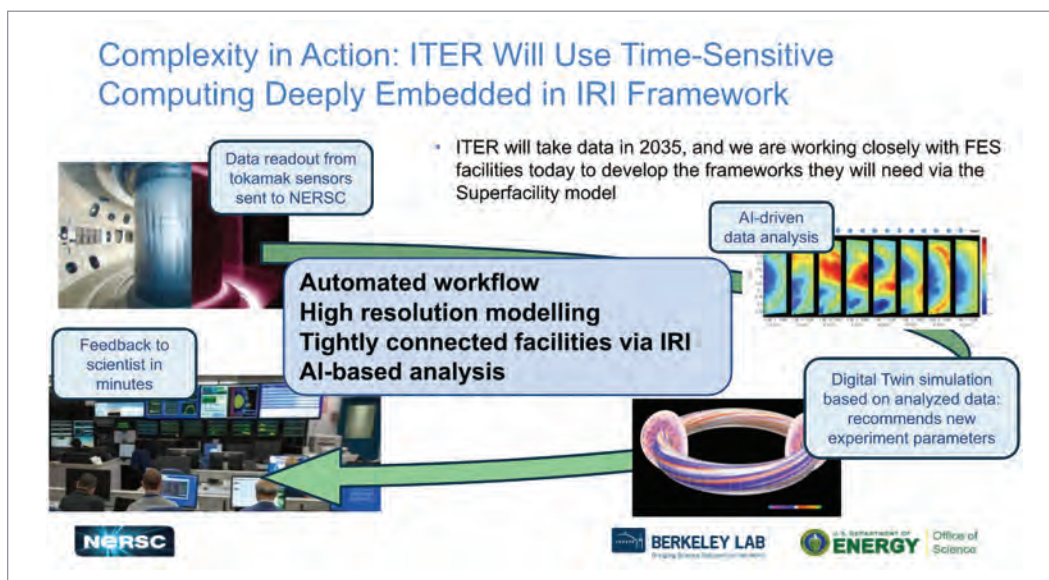


Figure 3-1. Connecting facilities like the ITER fusion facility to HPC resources at NERSC via ESnet is a complex workflow that will yield data to scientists across the globe in a matter of minutes.

At the DOE level, the IRI initiative aims to enable seamless workflows via close collaboration between ASCR facilities and the DOE scientific community. The initiative has kick-started new collaborative work across ASCR facilities and highlighted the importance of coordinated cross-facility projects. NERSC is deeply involved with IRI collaborations across SC and NERSC staff are helping to lead the initiative. The experience gained in developing, deploying, and supporting these services has positioned NERSC well for future work with NERSC-10 and IRI, particularly in designing tools that can be supported as they scale out.

NERSC is preparing for a future in which integrated workflows are the norm by architecting systems designed for these workflows from the start. NERSC-10, NERSC-II, and the surrounding infrastructure will be designed to support complex simulation and data analysis workflows at high performance through the following:

- **Quality of service:** Computation, storage, and networking capabilities enable quick response time and utilization.
- **Seamlessness:** Tight integration of system components enables high-performance workflows.
- **Programmability and automation:** APIs are used to manage data, execute code, and interact with system resources to enable automation and integration.
- **Orchestration:** Resource management is coordinated across domains.
- **Portability:** Modular workflows are executed across IRI sites.
- **Security:** Authentication, authorization, and auditing are essential components in a more connected ecosystem.

In preparing for the future, sociology is as important as technology—so in addition to fielding systems built to support complex distributed workflows, NERSC is preparing the workforce of the future now, getting users and staff ready to use the latest technologies and those still on the horizon. The NESAP program, in which NERSC staff help prepare science teams and workflows, is one major tool for ensuring that users are ready. The NESAP strategy is to:

- Form deep partnerships with science teams, working to directly improve workflows for NERSC-10 (and beyond) and gain expertise that can be applied moving forward
- Take lessons learned from these deep-dive partnerships and share them with the NERSC community at large

This combination of approaches is positioning NERSC to continue its leadership in the field of IRI and preparing the larger scientific community for future workflows, team by team and project by project.

Goals

1–2 years

- Fully engage with DOE’s IRI program, including taking lead roles in the Leadership Group and the technical subcommittees
- Launch NESAP for NERSC-10, focusing on improving the performance of end-to-end workflows and enabling them to leverage all the features of the NERSC-10 system

3–5 years

- Stand up services, policies, and frameworks identified by the IRI program, including APIs, data management, and policies to enable users to operate across multiple ASCR facilities
- Demonstrate the first IRI workflows running on NERSC-10
- Define how NERSC will interface with the High-Performance Data Facility (HPDF)

6–10 years

- Complete Integrated Research Infrastructure deployment
- See scientists using IRI interfaces to seamlessly integrate their workflows with NERSC systems and access data stored at NERSC
- Offer multi-site data management and orchestration, so that users can move data across IRI sites seamlessly and data from NERSC is available wherever a job runs

3.2 Pervasive AI for Science

In the last decade, NERSC has helped drive the emergence of modern AI applied to science via HPC. Over time, this progress has supported innovation and discovery across DOE SC.

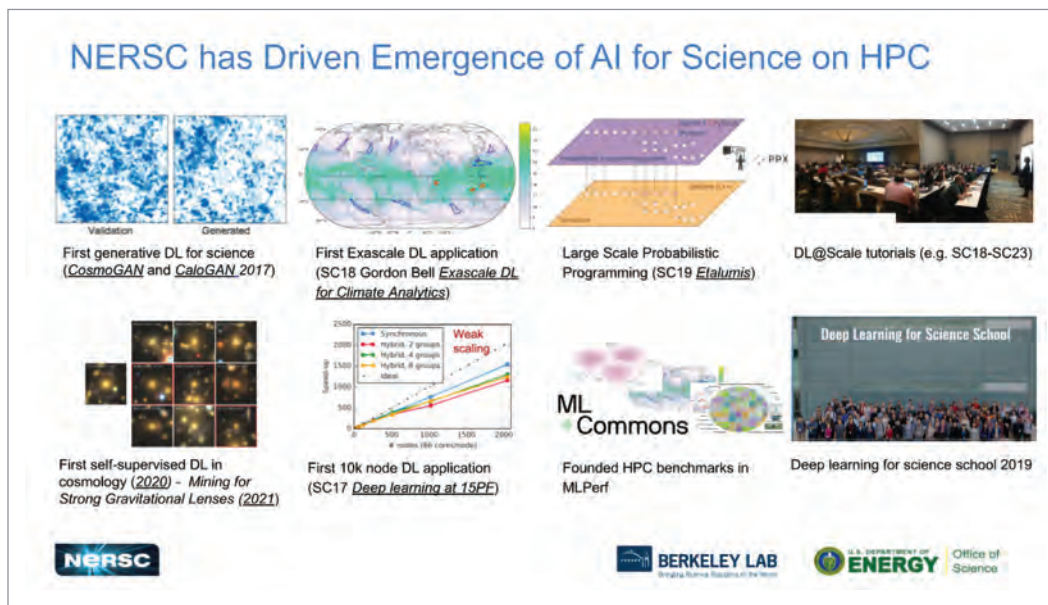


Figure 3-2. Highlighted examples of projects at NERSC that have driven the recent emergence of modern AI for science across different areas of applications, HPC, and empowerment of the community.

Perlmutter and the NESAP program are now enabling adoption of large-scale and groundbreaking AI with models trained on thousands of GPUs. These workflows are beginning to achieve the precision of traditional scientific simulation, with incredible speedups. This work is opening new doors of discovery — but to unlock the transformative potential of AI, we will need to go beyond these achievements to build an ecosystem for pervasive AI. With this ecosystem in place, NERSC can enable users to use AI to accelerate scientific discovery wherever and whenever possible.

To be successful, such an ecosystem must provide cutting-edge HPC systems for AI together with standardized frameworks and adaptable tools for use in AI training in addition to leveraging across experiments, applications, and science domains. The fast pace of AI development motivates the continuing of NERSC’s approach of deep engagements with domain scientists and AI experts to develop technology and expertise that is then brought to the wider community.

NERSC’s strategy for building a pervasive AI ecosystem will require activity across a variety of key areas including:

- **System** hardware and software that liberates scientists to apply large AI models
 - Accelerators for pervasive AI, as well as for workflow and data management, reflected in the architecture of NERSC-10 and beyond
 - Highly-instrumented, “self-driving” systems
- A **service platform** for seamless experimentation and integration of AI with simulation and data

- Host foundational AI models and datasets
- Intelligent AI-driven interfaces to compute
- **Applications** for science with large-scale, science-informed, robust, transferable models
- **Ecosystem** to empower scientists to use pervasive AI with human and AI-driven expertise

Beyond NERSC-10 and toward the ten-year horizon, by exploiting innovations in the other areas of NERSC’s strategy, AI will become seamlessly integrated into all science workflows, opening up new potential for scientific discovery.

Goals

1–2 years

- Determine AI capability requirements for NERSC science and HPC system operations
- Produce initial design and prototypes of components for a next-generation NERSC AI services platform

3–5 years

- Deploy initial AI services platform supporting full development and deployment life cycle for AI workflows
- Demonstrate impact on several (>3) NERSC-10 application workflows

6–10 years

- Provide world-class DOE mission platform for pervasive AI
 - Leading in rapidly evolving AI4Sci methods and systems landscape
 - Enabling AI in all aspects of scientific workflows and in HPC system operations
 - Driving impactful scientific discoveries

3.3 User Engagement

The breadth and depth of NERSC’s work and impact are constantly growing. The center currently serves nearly 11,000 active users and is cited in over 2,000 papers each year – both figures that continue to climb. The NERSC user base has grown by over 20% in the past five years, and we expect it to double in the next 10 years to roughly 20,000.

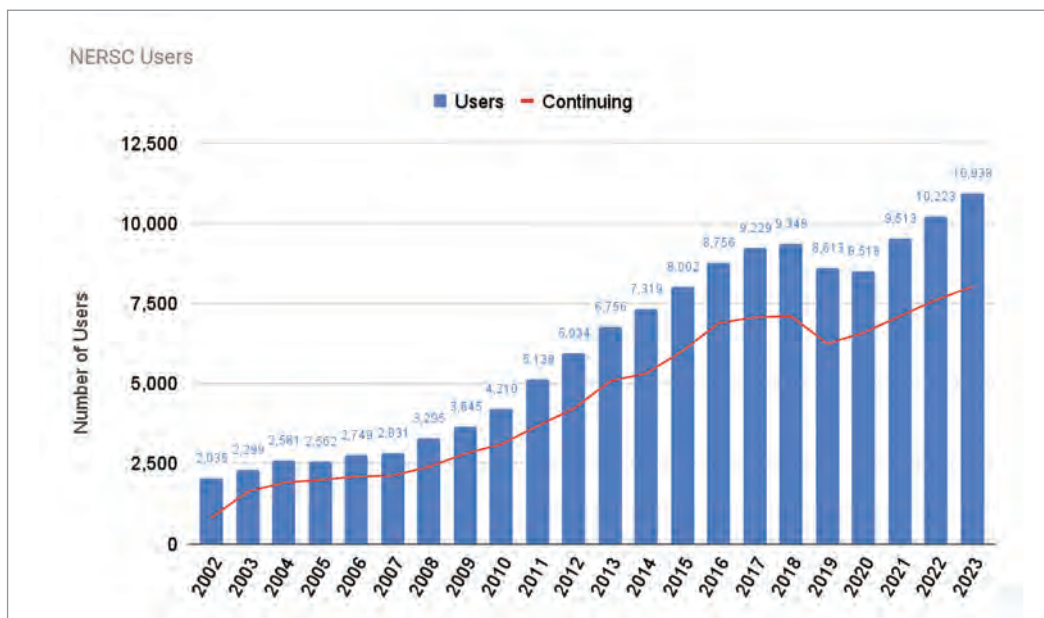


Figure 3-3. The NERSC user base has grown steadily in recent years. The dip in 2019 represents the point where PIs were required to confirm continuing users, and we accordingly dropped many inactive users.

The diversity of workflows users pursue at NERSC is also steadily increasing as users interact with NERSC systems in a greater variety of ways. Evolving the NERSC user engagement model can help NERSC scale with its growing community to support larger, more diverse, and more complex workflows.

NERSC currently has one of the best-reviewed engagement experiences across major supercomputing facilities, including a wide range of user education, technical training, and outreach opportunities. In user surveys and through direct engagement, users consistently praise the variety of supports NERSC offers, including:

- The NESAP program
- Globally referenced documentation
- Vendor interaction and technology development
- Broad user training at all levels
- Community-wide tutorials and knowledge sharing seminars
- Early allocation support for new growth areas (such as AI and quantum)

In ten years, the NERSC user base will have doubled and their workflows will have expanded and integrated across multiple facilities, including more AI and potentially moving beyond classical computing architectures. Expanding the world-class NERSC user experience to meet these challenges (without significantly increasing the NERSC staff effort) will be challenging. In order to scale, NERSC will need to

find effective and efficient ways to leverage expertise and knowledge from other sources, including users, vendors, in-house research, and other fields. This will include:

- Creating a NERSC user community of practice to enable users to learn directly from each other
- Expanding engagements with vendors, the open-source community, and other areas to leverage community-driven expertise
- Working with IRI and other ASCR facilities to build an HPC-expert community across the DOE
- Taking advantage of advances in technology to scale NERSC support and increase the impact of staff expertise, including automation and using AI-enabled tools where appropriate

Goals

1–2 years

- Improve, reorganize, and add to NERSC user support resources to better serve current users and advanced workflows and reflect changes in HPC education and knowledge
- Expand and enhance NERSC community engagement via diverse forums, including the NERSC User Group (NUG) & NERSC User Group Executive Committee (NUGeX), IRI, and community requirements reviews

3–5 years

- Investigate, and if beneficial deploy, an automated and interactive onboarding experience, including a learning management system with courses, short videos, and documentation.
- Ensure the 25% of new users each year can learn how to operate at NERSC without direct NERSC staff support
- Identify users with unique workflows, methodologies, and preferences, and better design the NERSC user experience to support and enhance these users' projects

6–10 years

- Develop NERSC users into a user community of practice, an active community helping and collaborating with each other and with NERSC in new and effective ways
- Ensure all NERSC users can productively use NERSC systems with minimal staff intervention

4. Focus Area: System and Data Center Architecture

4.1 Future Systems

NERSC will need to complete two facility upgrade projects (NERSC-10, NERSC-11) in the next ten years (2024–2034) to continue to support the DOE SC mission. Facility upgrade projects are driven by the approximately five-year life cycle of high-end supercomputing systems.

NERSC deployed *Perlmutter* (NERSC-9) in 2021 and will operate it through 2027. The NERSC-10 system anticipates user access in 2027. NERSC-10 is [designed](#) to deliver a large-scale HPC system with a rich workflow-enabling environment capable of running hundreds of diverse workflows simultaneously. NERSC-10 will support IRI services for complex multi-facility integrated science, automated AI-driven workflows, and real-time interactive computing in collaboration with HPDF, ESnet, Argonne Leadership Computing Facility (ALCF), Oak Ridge Leadership Computing Facility (OLCF), and other SC user facilities.

The [NERSC-10 Workflow Archetypes White Paper](#) outlines key workflow scenarios that the NERSC-10 system is expected to support. The NERSC-10 system is central to meeting SC goals and national initiatives [1] and avoiding the emergence of a gap between SC programmatic needs and HPC capabilities.

4.1.1 NERSC-11 and Beyond

Even as NERSC procures the NERSC-10 system, we are looking to the horizon and considering future systems and the technologies they might incorporate to support users in the 2030+ time frame, beginning with NERSC-11. These technologies include:

- Chiplets & specialization
 - Custom silicon for selected kernels
 - Estimated O(10x) performance AND O(1/100x) silicon
- Neuromorphic computing (spiking neural networks)
 - Novel non-von Neumann processing model on standard CMOS
 - Highly parallel, temporally sparse activity to allow for extremely efficient computation
 - Potential applications in image processing
- Analogue computing
- Optical computing
 - Based on optical interference instead of transistor switching
 - High-performance and low power for some workloads
 - Technology exists today (e.g. Lightmatter.com), but is not mature and currently limited to narrow workloads (matrix multiplication)
- Quantum Computing (See section 4.1.2)

Overall, as NERSC looks to the future and prepares for new paradigms and technologies, different approaches and new collaborations will be required.

The introduction of post-Moore's law technologies in NERSC-11 will unlock the potential to tackle [problems](#) that are infeasible on current HPC systems. In order to enable NERSC users to realize this potential, NERSC

will continue to partner with other laboratories and the vendor community to assess the applicability of these technologies for the NERSC workload; refine and use its benchmarking and assessment strategy to determine whether these technologies provide a significant enough advantage to warrant deployment in NERSC-II; assess programming methods and algorithms for these new technologies; and develop user support and training materials to accelerate their adoption. NERSC has initiated an early investigation of quantum computing, [partnering with hardware vendors](#) to get early access to the technology and assess its potential for accelerating science applications of interest to the DOE.

Goals

1–2 years

- Prepare for a beyond-Moore’s future, exploring power and resource efficiency and initiating collaborations with other HPC centers to chart the future

3–5 years

- Complete the conceptual design for NERSC-II system, including technology evaluation, extension of QOS capabilities, and NRE
- Complete the conceptual design for the NERSC-II facility

6–10 years

- Acceptance of the NERSC-II system with the following capabilities:
 - Integrated with IRI to execute workflows
 - Pervasive AI within user workflows and for system operations
 - Differentiated from commercial offerings
 - Incorporating post-Moore technologies as appropriate
 - World-leading resource efficiency (energy, power, water)
 - A substantial increase over NERSC-10 in performance

4.1.2 Quantum Computing

NERSC has been actively investigating and preparing users for potential technology disruptions such as quantum computing. Quantum computing may offer a solution for currently intractable problems, as well as improved time-to-solution for certain problems compared with classical computing; in the meantime, over 50% of computing cycles at NERSC solve quantum mechanical problems, suggesting that quantum computing could be a powerful tool for many users.

NERSC is preparing for the advent of performant quantum computing by developing an understanding of how it will interplay with future workloads and enable new scientific applications. Along with technological advances, NERSC is also preparing users and the workforce by providing training and tutorials, as well as access to both classical and quantum computing technologies. Currently, NERSC is working with a small fraction of its user base, with plans to scale more broadly as capabilities advance. Our goal is to ensure that users have access to state-of-the-art quantum computing resources including hardware, software, algorithms, and support.

However, the future of quantum computing is unclear; numerous quantum companies are developing a variety of available technologies, and it’s yet unknown which will be the “winner(s)”. It’s also unclear whether one technology will apply to all relevant parts of the NERSC workload or whether NERSC’s needs will be more diffuse.

Today's quantum computers are not yet viable for user scientific applications:

- Programming languages are not currently developed enough to enable widespread use
- Hardware is not mature, facing issues with qubit lifetimes, gate fidelities, scalability, etc.
- Algorithms need to be co-designed with hardware characteristics
- User projects require deep engagements by domain experts

In the midst of this uncertainty, NERSC is investing in quantum information science (QIS), investigating what's possible with hybrid quantum-classical technologies. The NERSC quantum team is planning for the quantum future over the next ten years by assessing the maturity of quantum hardware for the potential installation on premise in the 2030s.

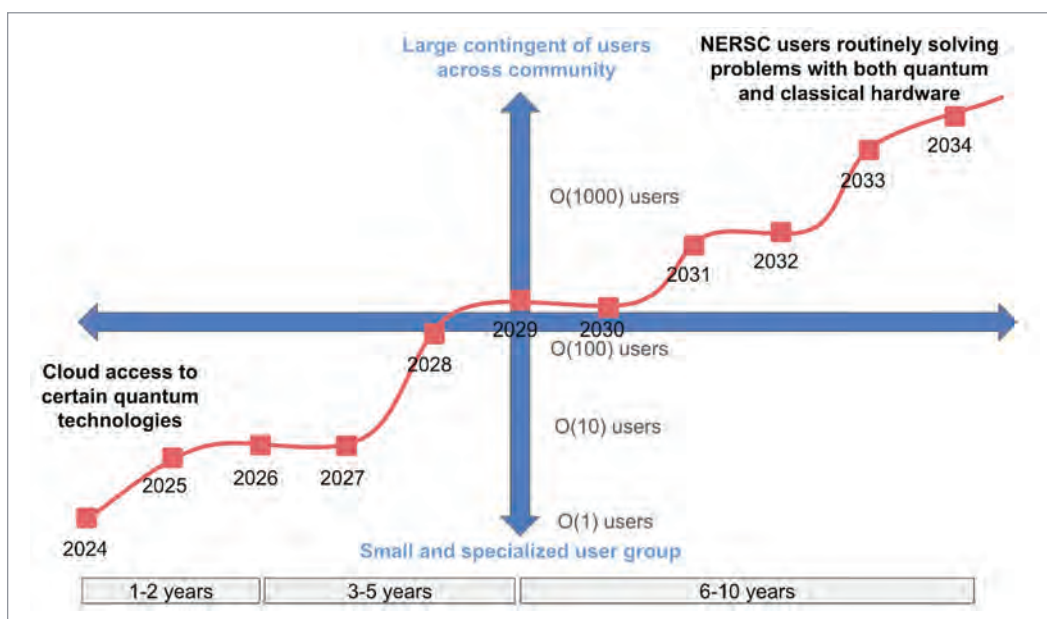


Figure 4-1. The path of NERSC quantum strategy looks to systems still on the horizon.

| Next 10 years | | | |
|--|---|--|---|
| 2022-2024 | 2024-2026 | 2026-2029 | 2030-2034 |
| <ul style="list-style-type: none"> • Ramp up engagement with QIS community • Director's Discretionary Reserve Call for quantum information science (QIS) on Perlmutter | <ul style="list-style-type: none"> • Enable user access to quantum hardware • Engage with relevant quantum vendors • Develop/evaluate quantum and hybrid quantum-classical algorithms • Identify opportunities for quantum-accelerated HPC codes • Benchmark quantum hardware • Perform resource analysis for executing useful quantum algorithms | <ul style="list-style-type: none"> • Availability of near-term quantum hardware becoming standard • Users request both classical and quantum resources • Workforce development through training / tutorials / quantum day • Evaluate the need and requirements for quantum hardware on premise | <ul style="list-style-type: none"> • High-performing quantum hardware becoming available • Potential integration with traditional HPC • Users routinely solve problems using both quantum and classical hardware |

Figure 4-2. The NERSC ten-year quantum computing roadmap, which is already in motion, is preparing NERSC staff, users, and systems for robust and tightly integrated quantum computing for science.

Goals

1–2 years:

- Develop a small-scale quantum user program, in the spirit of NESAP, with tens of quantum users
- Enact Quantum@NERSC outreach and engagements in the broader QIS community
- Expand the expertise of the NERSC Quantum Team to include software / hardware and algorithms, and increase staff in proportion to quantum user growth

3–5 years

- Grow the NERSC quantum user program to hundreds of users with a proportional increase in dedicated quantum staff
- Evaluate the needs and requirements (building, staffing, etc.) for procuring quantum hardware on premise
- Provide leading technology QC resources to hundreds of users through partnerships with relevant quantum hardware vendors

6–10 years

- Thousands of NERSC users solve problems using quantum hardware, enabled by appropriate staffing levels and expertise
- NERSC is positioned, assuming proper maturity of the technology, to procure and successfully manage post-Moore HPC systems that incorporate state-of-the-art quantum computers, opening a pathway to the future

4.2 Software Ecosystem

To meet the demands of increasingly complex scientific workflows and fill in gaps in vendor software offerings, NERSC, DOE and the HPC community must make critical investments in the scientific HPC software ecosystem, from the system management layer up to the user programming environment.

The NERSC HPC software ecosystem will be the cornerstone of a dynamic high-performance workflow environment. This hybrid environment will provide the complete HPC experience and enable seamless integration of cloud execution models needed for complex workflows that span within and beyond a single HPC facility.

This vision will be achieved through the Open HPC software environment, designed with a modular and vendor-agnostic infrastructure with standard interfaces for interoperability. This modularity starts at the system management layer with a control plane able to incorporate new hardware targets and flexibility to adapt to new technology trends without increased management burden. User and development software environments will be containerized and instantiated on systems, increasing portability and reproducibility from the base user software layer. In this modular environment, cloud-like services can be exposed and deployed alongside traditional HPC capabilities to create a complete workflow environment.

Since vendors are no longer providing the needed full HPC software environment, NERSC and DOE will need to develop this HPC software ecosystem. This will require key investments in workforce development to build the necessary skills. Collaboration will be essential to realize the modern HPC ecosystem; this will be a substantial effort with many components that no single HPC facility can provide. The HPC community will need to work together to implement a standards-based approach to enable modularity that retains flexibility for the use of site-specific and third-party software components.

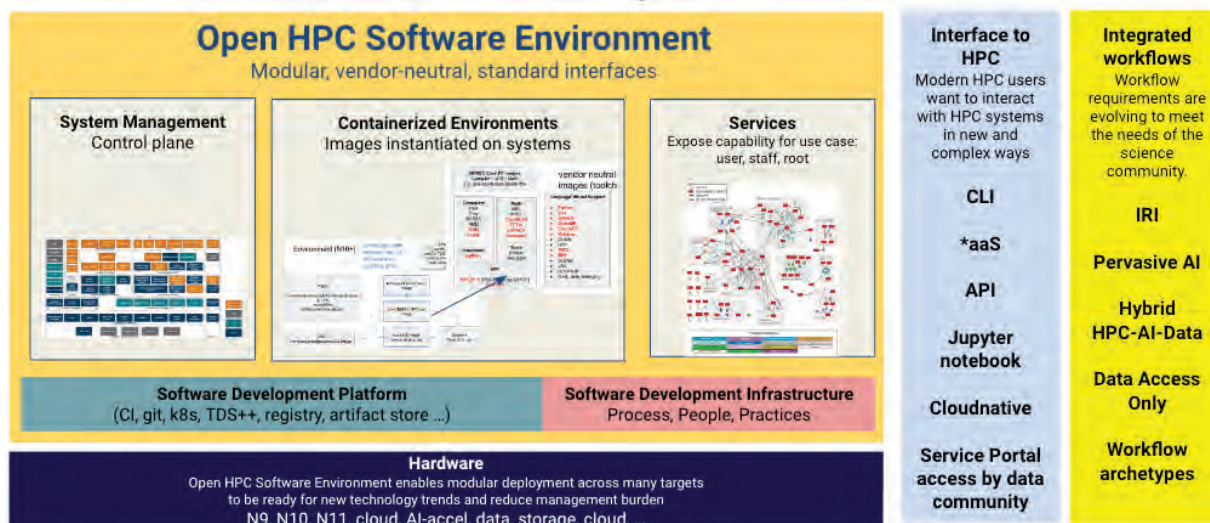


Figure 4-3. In the modern HPC software ecosystem, the NERSC software strategy builds on modular, flexible systems that can support a range of technologies and user needs.

Goals

1–2 years

- Develop a management plan for an open and modular software ecosystem
- Develop a software and services roadmap tied with resource (people and infrastructure) analysis

3–5 years

- Deploy the Open HPC software environment on the NERSC-10 system and demonstrate new capabilities for scientific workflows without sacrificing HPC performance
- Develop collaboration partners across the ASCR, DOE, and HPC communities

6–10 years

- Seamlessly accommodate new hardware and user requirements with vendor-neutral operational continuity

4.3 Data Center Design and the NERSC Center Roadmap

A primary goal for the NERSC Data Center Design is to evolve in alignment with large systems and seamlessly interoperate all systems to support future integrated workflows. Effective planning keeps the NERSC data center modern and enables updates that are largely non-disruptive and transparent to users. This allows continuity as storage, network, monitoring, and edge services remain constant, reliable, and available during the transition to the next large system—for example, portals remain available and serve data from off-platform storage, allowing user productivity to remain steady.

In addition to user productivity, effective data center planning allows NERSC staff to work more efficiently and have more time for strategic projects, enabling daily operational work to mesh seamlessly with forward-looking strategic planning.

Goals

1–2 years

- Establish a Cross-Team Center Roadmap Committee as a structure for project review, scheduling, and resource planning. Gather and prioritize requirements from NERSC-10 and IRI workflow analysis efforts.
- Develop a Center Roadmap Tool to collect and visualize NERSC projects

3–5 years

- Simplify and automate NERSC-10 interoperability with other data center resources
 - Increase ability to support emerging technologies
 - Implement continuous improvement to data center without major disruptions to users
- Mesh operational efforts seamlessly with forward-looking strategies, including NERSC-II
 - Train staff in using modern hardware, software, and techniques
 - Ensure staff development in critical areas

6–10 years

- A well-designed data center effectively interoperates to support evolving integrated workflows. The data center utilizes automation of operations and leverages IRI and other integrated resources to allow for seamless user productivity.

5. Focus Area: Smart Green Facility

5.1 Sustainability

NERSC's goal is to be a recognized leader in energy-efficient data centers. Although defining and measuring energy efficiency is complex, NERSC will strive to meet this goal by:

- Maximizing the science delivered per unit of energy, aiming for the smallest carbon footprint and the largest carbon handprint possible
- Influencing best practices in governance, policy, and operations
- Improving energy consideration in system operation (e.g., tuned approximations, job scheduling, node frequency adjustments, etc.) in tandem with conventional efficiency efforts of computer hardware and mechanical/electrical plants

In addition to optimizing NERSC's individual energy costs, the center's sustainability efforts are poised to have a broad impact across policy and operational methods within the scientific computing and general HPC industries. Leadership in this area supports Berkeley Lab, UC, and DOE sustainability goals and establishes a best-practices standard for other facilities.

Currently, NERSC uses Power Usage Effectiveness (PUE) as its major metric for energy efficiency, as is common across the industry. However, this measure is self-implemented and susceptible to inaccurate and inconsistent measurement methods as well as manipulative uses and interpretations. PUE is also not a great measure of comprehensive energy efficiency comparisons from site to site because it does not consider the following:

- Regional climate diversity
- Efficiency and carbon content of power sources
- Water consumption metrics and water sources
- Efficiency of IT equipment
- Efficiency of HPC system operations (SysOps)

In the coming years, NERSC plans continued adoption and implementation of comprehensive efficiency measures that consider energy and water sources use, computing throughput efficiency, and the identification of cost-effective waste heat activities on campus.

The challenges for this plan include a lack of standardization and regulatory requirements, especially for existing data centers. NERSC will also need to develop an accurate method of measurement for real-time computing efficiency (such as FLOPS per watt). However, beginning to address these issues through policy and action is one way in which NERSC can influence HPC best practices writ large and bring the field into greater alignment with the DOE mission.

Goals

1–2 years

- Develop and hone metrics for IT-Power Usage Efficiency (iTUE) and Water Usage Efficiency (WUE)
- Identify options for additional water and power savings through reuse of resources
- Continue to identify and develop opportunities for AI to provide additional levels of refinement in operational efficiency

3–5 years

- Develop plans to implement resource reuse for water or power
- Begin to assemble a comprehensive energy efficiency metric system for use at NERSC and extension to any data center
- Design the NERSC-II system to minimize the use of electrical equipment and power distribution

6–10 years

- Achieve recognition as the leader in energy-efficient and sustainable HPC data centers
- Leverage our expertise to assist other HPC centers and continue to incorporate continuous improvement initiatives learned from other HPC centers

5.2 Smart Facility

To provide reliable performance to users, swiftly address issues, and improve NERSC's efficiency, NERSC staff must have a clear view of how the full ecosystem operates and how users interact with it. NERSC is working toward a fully automated monitoring environment.

Within ten years, NERSC's vision is to reach the capabilities of a self-driving smart facility, offering the following:

- Automatic detection and correction of most data center issues
- Dynamic power management throughout the data center and HPC system, including smart scheduling that intelligently co-schedules multiple jobs and multiple hardware characteristics of the system, as well as across the wider center and IRI
- Identification and suggestion of optimal job parameters including topology

An additional goal is to gather and provide access to actionable user data that scales to full system usage and full community access. When this vision becomes a reality, both users and staff will be able to access detailed information about individual and aggregate job execution over time, as well as actionable insights, with the click of a button. This includes detailed use of hardware over time, IO performance, failure modes, and more.

This process will build on the foundations established by the NERSC OMNI team, which has already assembled cutting-edge monitoring infrastructure. From these efforts NERSC has developed expertise in:

- Collecting and centrally storing large amounts of center, system, and job data
- Understanding speeds and feeds characteristics
- Organizing data from multiple sources, which may also change over time, in one place with consistent schema

One early success in this area is the Lightweight Distributed Metric Service (LDMS) API. LDMS captures data from all nodes including detailed performance data from NVIDIA's Data Center GPU Manager. This allows staff to connect data to user jobs, visualizing job characteristics such as GPU usage, CPU usage, and memory leak detection, across all nodes used. By scaling this type of concrete progress center-wide, NERSC can make real steps toward its long-term vision of becoming a self-driving smart facility.

Goals

1–2 years

- Develop a centralized monitoring team at NERSC
- Develop a plan including milestones, software roadmap, hardware architecture, cost, and FTE count for center-wide monitoring, in the context of NERSC-10

3–5 years

- Create tools and dashboards to prepare for automated monitoring, including how to leverage AI
- Deploy automated monitoring on the NERSC-10 system.

6–10 years

- Automate analysis of monitoring data so that > 90% of software-related issues with systems including the facility may be resolved automatically
- Create the ability to simulate some of the complex environments associated with the center, system and jobs to help facilitate changes or mitigation of environmental change

5.3 Security in an Open Science Environment

NERSC Security is advancing NERSC's mission to perform open science by leveraging new technologies and processes to protect the NERSC technical ecosystem. NERSC strives to enable users to conduct their science with the least friction possible while balancing the need for managing risks and maintaining compliance. NERSC has a track record of developing innovative approaches to achieve usability without sacrificing security. However, the evolution in workloads and user interaction with the system, including integrating across facilities, creates new challenges. These, coupled with the increasing sophistication of cybersecurity attacks and new compliance requirements, underscore the need for continued attention, investment, and innovation.

The security landscape at NERSC is broad and diverse. At the network level, NERSC experiences several million scans – interactions from external IP addresses that are exploratory in nature for security vulnerabilities – every month. The HPC, storage, and building systems also require a broad set of security controls. In addition, NERSC provides several application-level services like containers and APIs, which are publicly accessible and could be used as attack vectors. In addition to standard user accounts, NERSC supports non-human/programmatic access use cases where services can connect and use center resources. For these reasons, NERSC Security is focused on four main areas:

- **Enhanced threat detection and monitoring:** Ensuring full security visibility within the NERSC environment, enabling NERSC to quickly detect and respond to malicious activity and gain insight into any attack vectors that are being utilized. Identifying any gaps and continuously evolving our monitoring techniques to match emerging technologies and new methods for user access to NERSC systems.
- **Risk assessment and management:** Performing in-depth security risk self-assessments on a regular basis. Implementing a robust third-party risk management (TPRM) program to identify, track, and reduce risks that are introduced to NERSC through our relationships with vendors, service providers, and other third parties.

- **Ensuring compliance** with regulatory requirements (Zero Trust, Export Controls, NIST 800-53, Berkeley Lab, DOE, OFAC, etc.)
- **Integrating security assessment and review** throughout NERSC projects and efforts to ensure that security is considered throughout the project life cycle. This includes secure architecture design, code reviews, CI/CD pipelines, vulnerability management, and sharing best practices for developing secure software and services.

In the future, NERSC's goal is to use emerging technologies to its own advantage, placing the center in a proactive position with regards to possible threats.

Goals

1–2 years

- Develop and deploy new security monitoring capabilities for different modes of access to NERSC
- Ensure that NERSC's Zero Trust Infrastructure (ZTI) is in compliance with Berkeley Lab
- Improve and streamline the vulnerability management process within NERSC

3–5 years

- Improve Security Orchestration Automation and Response (SOAR) via new capabilities and automation
- Implement tooling to enable NERSC to automatically search publicly accessible resources for exposure of sensitive information
- Integrated Security Processes: Ensure that all teams collaborate to mature the security processes such that:
 - Security reviews are done from the design phase throughout the life cycle of software and services
 - Vulnerabilities are patched before systems, services, or code are put into production
 - Tooling is implemented in the SDLC and generates less than 3% false positives

6–10 years

- Use AI and automation to proactively identify and remediate vulnerabilities across the entire systems stack

6. Focus Area: Workforce Development

6.1 A Catalyst for Workforce Development

NERSC is a catalyst for HPC workforce development and holds influence across the HPC landscape. The skills staff and postdocs develop at NERSC benefit more than just NERSC; in the past decade, the NESAP postdoc program has prepared 35 professionals for careers in HPC, including alumni currently holding important roles at vendors and partners such as Microsoft, Intel, NVIDIA, and AMD. NERSC staff have also gone on to serve and lead across the DOE lab complex and at important labs internationally.

In addition to developing the skills and gifts of its staff, NERSC also invests in its users through world-leading training and documentation. Approximately 50% of NERSC's 10,000+ users are early-career (students/postdocs), and the new-user experience is a major thrust of user engagement efforts; in 2023, over 4,000 users took part in some kind of training event. Additionally, NERSC leads all other institutions in GPUHackathon.org mentors, supporting the HPC community at large.

NERSC also actively supports those without HPC experience who may be interested in joining the field. In 2023, NERSC hosted HPC Bootcamp in collaboration with ALCF and OLCF. The program exposed non-major students to HPC through the lens of sociological studies. NERSC has also been an early supporter of the Sustainable Research Pathways program matching emerging scientist students with research projects at DOE labs. In 2023, approximately 15 students were placed at NERSC through SRP.

Still, challenges to NERSC workforce development remain.

- **Internal workforce development:** Next-generation systems will require new (additional) skills due to changes in user use cases, technologies, and the vendor landscape.
- **NERSC user workforce development:** As NERSC and the HPC ecosystem evolves, we need to offer our users opportunities to grow as technologists.
- **NERSC pipeline development:** NERSC must continue to attract new people with the relevant skill set to HPC.

To address these challenges, NERSC is setting a new course for preparing and developing the NERSC workforce.

Goals

1–2 years

- Examine NERSC's workforce and develop a plan to strategically address any skill gaps. This also allows for employee career advancement and hiring new employees when necessary, and accommodates attrition through succession planning.

3–5 years

- Evaluate our performance management models and identify any changes that we would benefit from, e.g., appropriate ways to assess management success as well as peer assessments
- Build on NERSC's outreach strategies to bring more people into the HPC ecosystem
- Leverage the NERSC-10 NESAP program to develop our users and the broader workforce via best-in-class documentation, training, hackathons, and workshops.

6–10 years

- Establish NERSC as the preeminent HPC facility for enabling and catalyzing a range of HPC careers: within NERSC, as a NERSC user, or in the greater HPC ecosystem.

6.2 Creating a Hybrid Work Environment for the Future

Currently, NERSC's approximately 120 staff and postdocs work in a variety of modes. The majority work one to three days per week onsite, with about 15% teleworking or fully remote and a smaller fraction working full-time onsite. As hybrid work has become the new normal, NERSC is striving to adjust and build a new culture of work that remains positive, equitable, and productive under these conditions.

A positive hybrid work environment that accommodates many modes of work will:

- Improve collaboration, particularly between those in different work modes
- Maintain high visibility of outstanding employees, whatever their work mode
- Keep staff morale high

NERSC is planning new ways to adjust its work culture to promote collaboration and productivity regardless of work mode and to use its existing space efficiently. This includes leveraging NERSC's space at Berkeley Lab – Wang Hall – as an attractive location to collaborate.

Goals

1–2 years

- Continue to train staff and managers to work successfully with a hybrid workforce
- Develop metrics to measure employee satisfaction

3–5 years:

- Remodel Wang Hall for optimal collaboration experience
- Promote an internal culture of gratitude, with awards for helping others
- Continue measuring employee satisfaction and see significant improvements

6–10 years

- Establish NERSC as the preeminent HPC center that supports employees working in a variety of work modes and has excellent collaboration practices, including an amazing collaboration space in Wang Hall

References

1. 2024 Advanced Scientific Computing Advisory Committee Facilities Subcommittee Recommendations Report, May 2024. <https://science.osti.gov/-/media/ascr/ascac/pdf/reports/2024/52224Draft-ASCAC-Facilities-Report-003.pdf>
2. US Department of Energy Office of Science User Facilities. August 2024. <https://science.osti.gov/User-Facilities>

For more information about NERSC, contact:

NERSC Communications
Lawrence Berkeley National Laboratory
1 Cyclotron Road
Berkeley, CA 94720-8148
Email: CSComms@lbl.gov

NERSC's Web Site

www.nersc.gov

NERSC Strategic Plan Editor

Elizabeth Ball

Design

Design, layout, illustration, photography, and printing coordination:
Berkeley Lab Creative Services

Cover Image

Berkeley Lab's Building 59 Shyh Wang Hall (CRT/NERSC), exterior photo at sunset,
Credit: Roy Kaltschmidt, Berkeley Lab.

DISCLAIMER

This work was supported by the Director, Office of Science, Office of Advanced Scientific Computing Research of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor The Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or The Regents of the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof, or The Regents of the University of California.

Ernest Orlando Lawrence Berkeley National Laboratory is an equal opportunity employer.

