# HPDF
## High Performance Data Facility

# High Performance Data Facility Project Status and Plan

June 2024

**This material is based upon work supported by the U.S, Department of Energy, Office of Science,
Office of Advanced Scientific Computing Research under contracts
DE-AC05-06OR23177 and DE-AC02-05CH11231**

# Table of Contents

## List of Figures

# High Performance Data Facility Project: Status and Plan

## 1.0     Introduction

On October 16, 2023, the U.S. Department of Energy (DOE) launched the High Performance Data Facility (HPDF) Project to create a new scientific user facility specializing in advanced infrastructure for data-intensive science. DOE's Office of Science (SC) named Thomas Jefferson National Accelerator Facility as the HPDF lead and sited the HPDF Hub infrastructure at the lab's campus in Newport News, Virginia. The HPDF Project, sponsored and supported by the Advanced Scientific Computing Research (ASCR) program, will be a partnership between Jefferson Lab and Lawrence Berkeley National Laboratory (LBNL). The two labs moved immediately to form an integrated team led by Jefferson Lab.

When completed, HPDF will be a first-of-its-kind SC user facility that will expand and contribute to the world-class capabilities of the ASCR and SC data and computing infrastructure ecosystem. The facility's mission will be to enable and accelerate scientific discovery by delivering state-of-the-art data management infrastructure, capabilities, and tools. As a cornerstone of the DOE's Integrated Research Infrastructure (IRI) initiative, a successful, fully realized HPDF will be widely recognized as a national and international leader in uplifting data science and high-performance data infrastructure.[1]

HPDF is envisioned as a hub-and-spoke model, in which the Hub will host centralized resources and enable high-priority DOE mission applications at Spoke sites by deploying and orchestrating distributed infrastructure at the Spokes or other locations. The number and variety of Spokes are expected to grow and evolve with mission requirements, consistent with SC's deep experience with the Cooperative Stewardship model [1] for major research infrastructure.

The project team is tasked with designing and delivering a geographically resilient and innovative HPDF capable of meeting the needs of diverse users, institutions, and use cases. The joint project will itself provide the template for the first Spoke partnerships and blaze new paths in institutional engagement and outreach in the emerging era of artificial intelligence (AI)-enabled integrated science.

This document provides a summary of the HPDF Project's status and immediate plans to advance its design before the delivery of its FY 2024 funding.

## 2.0     Mission Needs and Capability Gaps

Every year, tens of thousands of scientists from universities, industry, and national laboratories rely on DOE user facilities to conduct research that expands what humans understand about the world around us. Whether this work involves small groups of researchers or collaborations across many domains, transformational research demands disruptive technology that can provide unprecedented data analysis, networking, and storage resources for the nation's science enterprises.

The SC stewards the world's most powerful constellation of major research infrastructure and scientific tools in its national laboratories and other sites across the country. Research at these facilities produces staggering amounts of experimental data that must be analyzed, preserved, and

---

[1] 2024 Advanced Scientific Computing Advisory Committee (ASCR) Facilities Subcommittee Recommendations (Technical Report) | OSTI.GOV

made accessible to the wider scientific community. The DOE envisions a revolutionary ecosystem – the IRI – that can deliver seamless, secure interoperability across these facilities and other national laboratory capabilities. From a human-centered perspective, the IRI vision empowers researchers using DOE's world-class tools, infrastructure, and user facilities in novel ways to radically accelerate discovery and innovation. The ASCR facilities – Energy Sciences Network (ESnet), Argonne Leadership Computing Facility (ALCF), Oak Ridge Leadership Computing Facility (OLCF), National Energy Research Scientific Computing Center (NERSC), and HPDF – will collaboratively forge the foundational infrastructure to enable IRI.  HPDF will be designed to overcome identified gaps in the current DOE computing and data infrastructure ecosystem, providing paradigm-shifting data-focused capabilities [2-4].

HPDF will lead the stewardship of the scientific data life cycle, advancing DOE's commitment to findable, accessible, interoperable, and reusable (FAIR) data principles. In January 2023, the White House Office of Science and Technology Policy launched the Year of Open Science, taking steps throughout the federal government to advance national open science policy, provide access to the results of the nation's taxpayer-supported research, accelerate discovery and innovation, promote public trust, and drive more equitable outcomes.

The 2023 IRI Architecture Blueprint Activity[2] identified three broad science patterns that demand research infrastructure interoperability:

- **Time-sensitive patterns**. Science cases that require end-to-end urgency, such as streaming data for real-time analysis, real-time experiment steering, real-time event detection, AI-model inference, or deadline scheduling to avoid falling behind.

- **Data-integration-intensive patterns**. Science cases that demand the combination and analysis of data from multiple sites and sources, such as AI federated learning, experiments, and/or simulations [2].

- **Long-term campaign patterns**. Science cases that need sustained access to resources over a long period to accomplish well-defined objectives, such as sustained simulation production, large data processing, extended experimental campaigns, and archiving for collaborative use and reuse.

These IRI patterns point to broad requirements for the DOE computational and data infrastructure ecosystem:

- **Data management.** A dynamic and scalable data management infrastructure that integrates with the DOE computing ecosystem (i.e., "networked data infrastructure at scale").

- **Data capture.** Availability of dynamically allocatable data storage and edge computing at the point of generation.

- **Data staging.** Dynamic placement of data in proximity to appropriate computing for reduction, analysis, and processing.

- **Data archiving.** Availability of extreme scale distributed archiving and cataloging of data with FAIR data stewardship principles.

---

[2] Integrated Research Infrastructure Architecture Blueprint Activity (Final Report 2023) (Technical Report) | OSTI.GOV

- **Data processing.** A diversified computing ecosystem that melds distributed and centralized computing assets (i.e., "the right compute resource at the right time") and combines the federated data catalog with a central workflow and automation system [5-8].

## 3.0 HPDF Project Design and Scope

The HPDF Project team is focused on meeting the mission needs and fill the capability gaps for providing the innovative capabilities and expertise that will enable scientists to integrate data services with diverse computing resources at scale through a distributed, resilient infrastructure. The HPDF Hub and Spokes will support the full data life cycle to facilitate global scientific discovery and advance the practice of scientific data stewardship across scientific communities. HPDF will play a critical role in developing and training the next-generation scientific data workforce.

HPDF will combine centralized and distributed infrastructure, incorporating resources at several locations. The design will deliver three key elements:

1. **Hub.** Centralized data-centric infrastructure with high availability and performance, as well as geographically and operationally resilient active-active failover.
2. **Spokes.** Distributed data-centric infrastructure to enhance HPDF access and support for science users by integrating distributed computing or storage resources.
3. **Integration and services.** Orchestration hardware, software, and services for data movement, storage and retrieval, and science workflow automation.

The HPDF Project scope will comprise the following:

- Data center site preparation, power, and cooling infrastructure at Jefferson Lab and Berkeley Lab.
- Design, acquisition, delivery, and commissioning of the Hub infrastructure at Jefferson Lab and Berkeley Lab.
- Design, acquisition, delivery, and commissioning of infrastructure for a selected set of initial Spokes.
- Integration of HPDF infrastructure with ESnet and the ASCR high-performance computing (HPC) facilities: NERSC, ALCF, and OLCF.
- Software development for core HPDF services and development of an operations team that will support the infrastructure and scientific users.

In conjunction with the HPDF Project, Jefferson Lab will design, construct, and commission the Jefferson Lab Data Center (JLDC) building using Commonwealth of Virginia funds. The JLDC scope will comprise construction of the building and infrastructure required for the HPDF Project to take occupancy. Design and construction of the JLDC will adhere to the Commonwealth of Virginia construction project management process and the guiding principles of DOE Order 413.3B. The JLDC requirements and schedule will be driven by and formally coordinated with the HPDF Project. The JLDC will report progress to the HPDF Project, as well as tracking performance for accountability to the Commonwealth.

### 3.1 Design Overview

Designing the HPDF requires careful consideration of multiple factors, including resilience, data locality, and the needs of streaming and real-time data. The vision for the HPDF distributed

architecture is illustrated in Figure 1. Jefferson Lab and Berkeley Lab will host the Hub core infrastructure. The infrastructure is designed to maximize planned availability, with a second site providing active-active failover and geographic resilience. The centralized Hub resources are supplemented by Spokes, long-term partnerships that manage distributed infrastructure and components. A particular Spoke might be established to support a particular SC user facility, the overarching needs of a specific SC program area, the needs of an entire cross-cutting DOE institution such as a national laboratory, or to support a specific DOE-funded project.
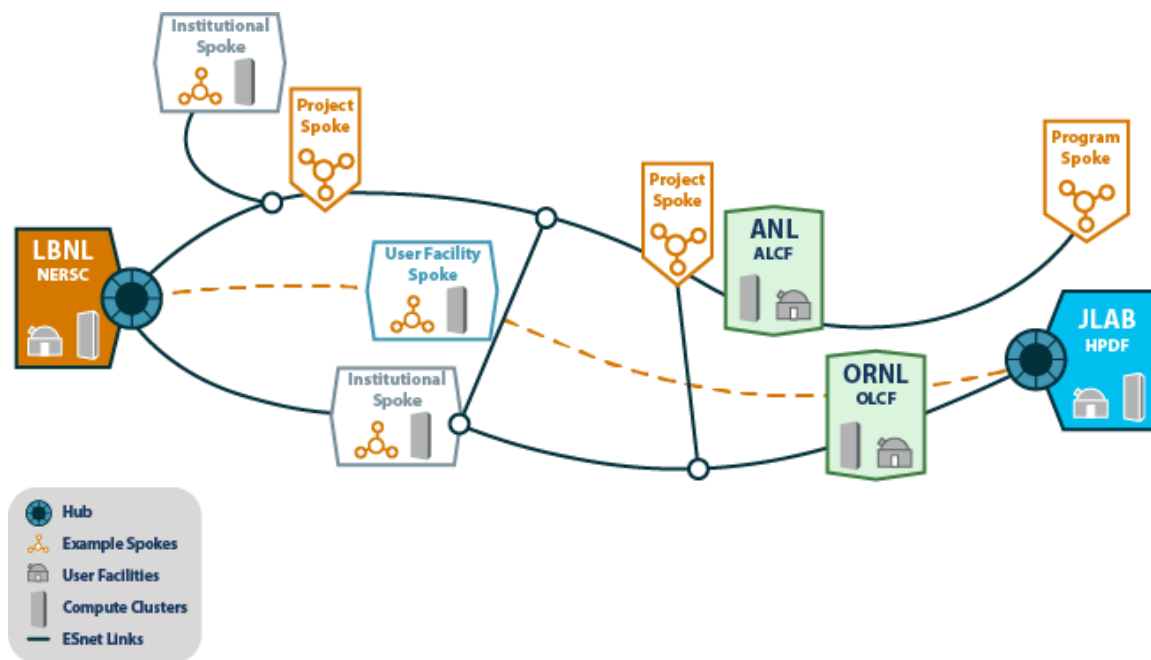


***Figure 1. Concept of the HPDF Hub and Spoke Model***
*The hub-and-spoke model will integrate the essential Hub infrastructure at Jefferson Lab and Berkeley Lab and the managed distributed infrastructure at defined Spokes locations. This model will ensure seamless resilience, performance, engagement, and customization for users.*

## 3.2 Hub Computing and Data Infrastructure Design

HPDF will provide operationally critical capabilities for data-intensive science. An extended gap in Hub services could potentially harm many science programs, and it is therefore important that the Hub maintain high availability and performance. The data and compute infrastructure at Jefferson Lab and Berkeley Lab will have the same architecture, running common software and providing consistent software and data services to enable portability and resilience.

HPDF will provide a data-centric infrastructure design that includes compute resources optimized for data processing and analysis. Broadly, this infrastructure will provide data management, data processing, and orchestration services and resources, as described in the following sections.

The Hub architecture will achieve these goals through the following design tenets:

- **Flexibility.** The computing infrastructure will feature modular standard units that can be configured to address the needs of scientific use cases.
- **Resilience.** The hub-and-spoke system will avoid single failure points. The infrastructure architecture will allow incremental maintenance and enable high availability through use of AI/machine learning (ML) for data center optimization, orchestration, and automation.

- **Consistency.** Data and compute infrastructure at both labs will feature common hardware and software to enable portable and resilient services.
- **Diversity.** HPDF will combine high-performance storage for active data, archive storage where necessary, and a distributed federated data catalog within the HPDF ecosystem.

## 3.3     Data-Centric Orchestration of Hardware, Software, and Services

HPDF will offer innovative production data services and software tools to support the entire data life cycle (Figure 2) in concert with the ASCR Facilities ecosystem and the Office of Scientific and Technical Information (OSTI). A key goal is facilitating data management and interoperability, i.e., making data available to a broad scientific community, providing for new technologies and user access patterns, and preserving the data for future use.
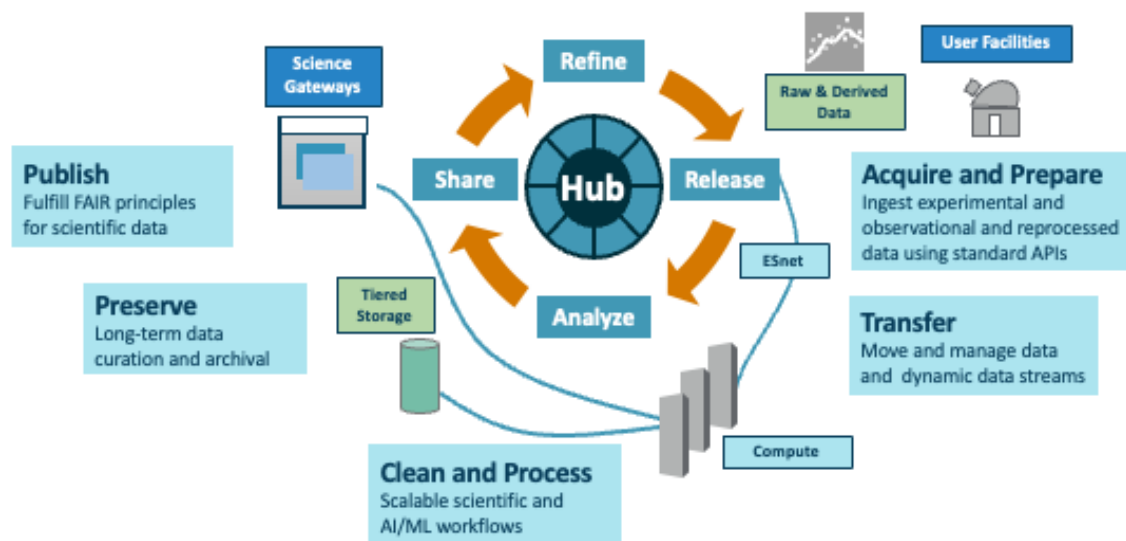


***Figure 2. HPDF in the ASCR Facilities Ecosystem***
*The HPDF Hub will provide core capabilities for management and curation of the full data life cycle.*

- **Data storage, access, and discovery.** HPDF will incorporate core capabilities to spearhead scientific data-driven discovery and provide access to data through various services (e.g., search tools) and interfaces (e.g., web and application programming interfaces).
- **Data life cycle tools and services.** HPDF will offer support for the entire data life cycle, providing tools and services to move, clean, process, analyze, share, and collaborate while supporting new technologies, such as AI and novel hardware or software technologies. User support and data stewards will help users and scientific communities advance data analysis.
- **Data preservation**. HPDF will deliver storage and access through a data repository framework that ensures FAIR data support and future availability via federated archives and catalogs while providing digital object identifiers in partnership with OSTI.
- **Seamless Data and Compute Infrastructure.** The free flow of data and workloads among HPDF resources, including institutional clusters, cloud, and HPC using IRI.

## 3.4     Infrastructure Design: Spokes

The hub-and-spoke infrastructure design and management framework will allow the scaling of HDPF resources and services for various strategic priorities, partnership contexts, and scientific communities. The Spokes are the most sophisticated and deepest HPDF partnerships, expressing

a mutual commitment to align the teams, resources, and policies users depend on for maximum productivity. They will be the first level of service for the communities they serve, providing an enhanced user experience and customizations that co-evolve with the Hub infrastructure as computing and data technology and solutions advance. The Spokes will add value by providing community-specific software and services, user support, and possibly hardware resources that mirror, supplement, or complement Hub resources.

A robust, multi-tiered science engagement process between the HPDF Project and its scientific communities will further refine the Hub and Spokes concepts. Ongoing user outreach will foster the dialogue essential for shared governance and a coherent, seamless user experience. The success of the model will depend on the human, policy, and technical interfaces between the Hub infrastructure and the Spokes, allowing scientific communities to achieve interoperability while developing specialized solutions. Early engagement between HPDF and its users will lead to creative thinking about how Spokes can help accelerate scientific discovery.

## 3.5    Integration and Services Design

Innovation and partnership will figure heavily in integration and development efforts. Our integration and services design is grounded in a cycle that identifies existing data services, evaluates these solutions, and integrates these solutions into the HPDF software stack. We anticipate leveraging products, best practices, and lessons learned from major scientific software integration projects, including the Exascale Computing Project and other ASCR projects, as well as earlier investments in research. Figure 3 illustrates the strategy for building and sustaining these integrated components contemporaneously with the development of the wider IRI ecosystem.

Design of the software infrastructure for integration of the HPDF Hub and Spokes and integration of HPDF into the IRI ecosystem will feature:

- Common application programming interfaces and data services to facilitate portability.
- Distributed orchestration and execution layers.
- Data transport, caching, communication, and monitoring built on ESnet6 capabilities.
- Dynamic virtualized compute and storage ensuring portability between sites.
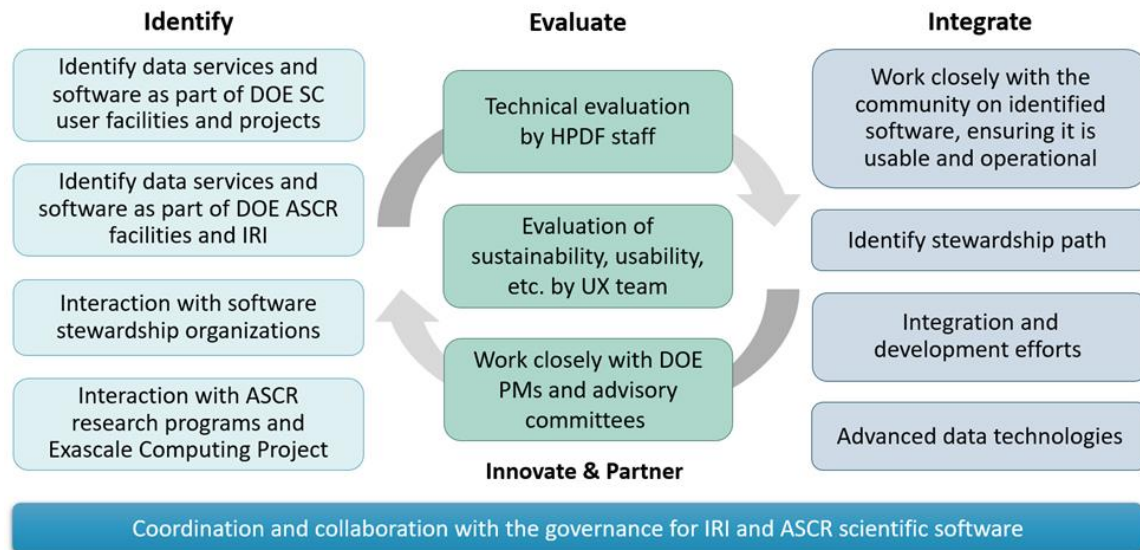- Cross-cutting components for security and monitoring.

*Figure 3. Strategy for HPDF Infrastructure Integration and Services*
*The HPDF Project will create a strategy for innovation and stewardship of the software and data services that span the Hub, Spokes, and, with the ASCR facilities, the broader scientific software ecosystem.*

## 3.6    Hub Physical Infrastructure Design

The design of the physical infrastructure for HPDF at Jefferson Lab and Berkeley Lab will be based on core requirements articulated above. It will be undertaken through iterative evaluation of those requirements against various risks, opportunities, and limitations.

### 3.6.1    Jefferson Lab Data Center

Jefferson Lab has previously conducted an analysis of alternatives to evaluate the concept of expanding the existing Jefferson Lab data center infrastructure vs. the concept of a new data center building; this analysis revealed that limits on space and access for utilities will make a new building necessary.

The Commonwealth of Virginia has appropriated $3M for preliminary design and $43M for design and construction of the Jefferson Lab Data Center building in the FY23 state budget cycle. Jefferson Lab anticipates these funds will be available to the Lab in FY25 ($33M) and FY26 ($10M). Planning is in progress with Dominion Energy for a new campus substation to accommodate additional power loads for HPDF.

The JLDC scope will comprise construction of the building and infrastructure required for the HPDF Project to take occupancy. The HPDF Project team has provided preliminary power, cooling, and space requirements to the JLDC Facilities team to begin design activities for the facility that will house HPDF systems. These requirements, based on estimates for compute, storage, and service hardware, have been reviewed and validated by experts across the lab complex and are being incorporated into a final design report that will inform the JLDC design contract. Additionally, JLDC and HPDF subject matter experts have visited several leading laboratories for data center architecture best practices and lessons learned.

Once the design of the JLDC is complete, construction and commissioning will take an estimated five years. Existing data center space will be adequate for near-term needs, including HPDF pilot activities and development of Hub software and hardware technologies.

### 3.6.2 Berkeley Lab Facilities

The HPDF Hub infrastructure at Berkeley Lab will be located in Shyh Wang Hall. This existing state-of-the-art LEED gold-certified facility uses free-air cooling and is also home to NERSC and ESnet. The building design allows for highly efficient operations, with a power usage effectiveness of <1.1 for its current systems.

Wang Hall has 20,000 square feet of existing machine room space and an additional 10,000 square feet of unfinished shell space to accommodate future growth. This unfinished space provides HPDF with space to meet the immediate needs of the project as well as future growth space. The existing finished space also offers possibilities to meet near-term needs such as hosting evaluation and early testbed systems. LBNL is actively studying potential long-term power and cooling needs for Wang Hall, incorporating consideration of HPDF, NERSC, and LBNL research computing requirements.

## 4.0 HPDF Project Plan of Effort

The near-term plan of effort for the HPDF Project is focused on the following activities:

- Establish a road map and timeline to achieve critical decision (CD)-1 for critical aspects of the technical design (hardware, software and services), including a notational staffing plan.
- Define cybersecurity considerations and requirements for the HPDF conceptual design.
- Define the design elements for distributed infrastructure.
- Define the design elements for Spoke integration.
- Establish the technical specifications that might affect facility needs, e.g., power, space, cooling, and other considerations, including and especially requirements for the JLDC.
- Conceptualize the goals and design for an early access system.
- Establish the HPDF Project risk management framework and processes, emphasizing a risk-based approach that considers trade-offs between ease of use and requirements.
- Create a preliminary community Spokes document that includes at least three distinct Spoke examples, refining the template as needed.
- Establish project advisory capacities, especially an interface to the IRI Management Council and other ASCR Facilities upgrade projects.
- Commence community engagement activities with a focus on core user requirements.

## 5.0    Appendices

## 5.1    Governance

### 5.1.1    IRI Governance

In FY 2024-25, the DOE is standing up the IRI Program governance, which will provide the structure to build and sustain resource interoperability across a broad set of stakeholders. The IRI governance mirrors the DOE management and operating partnership model: a steering committee of federal program offices that provides oversight of resources and requirements. Operationally, the governance structure centers around the IRI Management Council (Figure 4).

The Management Council's Executive Committee will provide direct oversight and is composed of executive-level leadership representing user facilities and IRI infrastructure contributors [9]. Day-to-day program management will be conducted by the IRI Leadership Group, in the future headed by a program spokesperson and deputy, positions that may ultimately be elected by the IRI community. Technical implementation and continual assessment of user requirements and satisfaction will be the purview of the Technical Subcommittees.



***Figure 4. Status of Startup of IRI Program Governance, March 2024***
*The ASCR Facilities (ALCF, ESnet, NERSC, and OLCF) provide the initial nucleus
of the IRI Management Council at startup.*

### 5.1.2    HPDF Project Governance and Execution

The governance and execution of the HPDF Project is framed by DOE Order 413.3B, with well-defined CD points [10]. HPDF will play a foundational role, with the existing ASCR facilities, in advancing the IRI. The IRI Management Council will serve in an advisory capacity to HPDF Project management at the executive and technical levels. The HPDF Project will work in conjunction with IRI technical subcommittees, as well as contributing functionality to the IRI (Figure 3).

The HPDF Project Executive Steering Committee, composed of executive leadership from the partner laboratories, will provide oversight and adjudication of any project-related concerns. This structure (Figure 5) provides a scalable model, with the Steering Committee expanding to include

Spoke representation as the project progresses. The role of the project leadership team will be to develop and execute the HPDF conceptual design scope, cost, and schedule. As Spokes are selected, their activities will be aligned with the HPDF Project through a joint management plan.



*Figure 5. Governance of the HPDF Project*

*The HPDF Project will have an integrated project team led by Jefferson Lab in partnership with Berkeley Lab, with oversight and support from a project Executive Steering Committee. The IRI Management Council will provide scientific and technical advice to the HPDF Hub Project. ASCR, Thomas Jefferson Site Office, and Berkeley Site Office will provide federal oversight and support.*

### 5.1.3    HPDF Leadership Team

**Amber Boehnlein** is the HPDF Project Director and Jefferson Lab's Chief Information Officer (CIO). She has been working on scientific computing for experimental workloads for 30 years. Dr. Boehnlein designed the software trigger for one of the Fermi National Accelerator Laboratory Tevatron experiments. She led the first large-scale data reconstruction campaign on distributed computing resources. She has worked in scientific computing with a focus on experimental data management at four DOE Laboratories covering three SC program offices.

**Lavanya Ramakrishnan** is the HPDF Deputy Project Director. She is Division Deputy in the Scientific Data Division at Lawrence Berkeley National Lab. Her work focuses on methods and tools to manage workflows and data while working closely with scientific groups on the design next-generation HPC systems. Dr. Ramakrishnan established and leads a scientific user research program that studies how scientists and communities use data and workflows to build usable tools for science. She serves on the High Performance Distributed Computing Conference Steering Committee, as Director of Training and User Engagement for the Workflows Community Initiative, iHARP National Science Foundation Harnessing the Data Revolution Institute's Advisory board, and the Computing Research Association's Widening Participation Board of Directors.

**Graham Heyes** is the HPDF Technical Director. He has been working on data acquisition systems and scientific computing at Jefferson Lab for more than 30 years. He is the lead architect of CODA, the common DAQ software and hardware toolkit used across the Jefferson Lab experimental halls. He served for three years as Computer Center Director. More recently, he served as the liaison between the experimental physics program and the Scientific Computing group before transitioning to lead Scientific Computing at Jefferson Lab.

**Richard Shane Canon** is the HPDF Deputy Technical Director. He is a Senior Engineer at Lawrence Berkeley Lab and has nearly 25 years of experience focused on supporting breakthrough science through large-scale computing systems, including some of the fastest computers and

storage systems in the world. Dr. Canon has deep experience with HPC, large-scale storage systems, data portals and platforms for science, and HPC containers. He leads the Data Science Engagement Group at NERSC, as well as Technology Integration efforts for the NERSC-10 Project.

**Theresa Bamrick** is the HPDF Project Manager. She brings 30 years of project and program management experience spanning several industries, including high-tech software development and commercial construction management. She established the IT Project Office at Stanford Linear Accelerator Center (SLAC) in 2012 and introduced IT Service Management. Ms. Bamrick served at SLAC as Associate CIO of IT Services, Deputy CIO for five years, and CIO from 2017-2020. She is Stanford Advanced Project Management and UC Berkeley Project Management certified, and an Agile Certified Scrum Master. She is a U.S. Marine Corps veteran.

**Becci Totzke** is the HPDF Deputy Project Manager. She is the Business Operations and Services Group Lead in the NERSC Division at Lawrence Berkeley National Lab and the NERSC-10 Project Manager. NERSC-10, a DOE O 413.3b Project, will deploy a high-performance computing system for the Department of Energy's Office of Science (DOE SC). Previously, as the NERSC-9 Project Manager, Ms. Totzke successfully closed out this DOE O 413.3b Project, which achieved CD-4 approval in 2022. She holds a Project Management Professional certification along with a UC Berkeley Project Management certification and is a U.S. Army veteran.

## Executive Steering Committee

**David Dean** is chairman of the HPDF Executive Steering Committee and Deputy Director for Science and Chief Research Officer at Jefferson Lab. He is responsible for the overall development, strategic planning, and oversight of the lab's scientific enterprise, ensuring alignment with DOE mission priorities. He provides leadership for strategic planning, mission diversification, and user community engagement and serves as the principal scientific contact with funding agencies, university partners, and the scientific community. Dr. Dean came to the lab in January 2022 from Oak Ridge National Laboratory, where he had served as associate laboratory director for the Physical Science Directorate since 2011. As ALD, he led four research divisions: Materials Science and Technology, Chemical Sciences, Physics, and the Center for Nanophase Materials Science.

**Jonathan Carter**, co-chairman of the HPDF Executive Steering Committee, is Berkeley Lab's Associate Laboratory Director for Computing Sciences, which includes the SC User Facilities NERSC and ESnet, as well as research and development divisions focused on applied mathematics, computer and data science. In prior roles, he served as deputy director for the DOE Exascale Computing Project Software Technologies Focus Area, participated in and managed the growth of the Quantum Information Science portfolio in the Computing Sciences Area, and led the project to procure the NERSC-6 Hopper HPC system. He joined Berkeley Lab in 1996.

## 5.2 Facilities and Other Resources

### 5.2.1 Thomas Jefferson National Accelerator Facility

Jefferson Lab, in Newport News, Va., is the home of the Continuous Electron Beam Accelerator Facility (CEBAF), a world-leading nuclear physics scientific resource in operation since 1995, supporting an 1,800-person user community. HPDF will leverage Jefferson Lab's experience and expertise in supporting a large scientific user community.

The Jefferson Lab Computational Science and Technology Division (CST) supports enterprise, scientific, and innovative computing activities. This leads to organizational alignment for scientific computing, networking, cybersecurity, and systems support. Tight integration between the IT and scientific computing teams enables the implementation of technical and security solutions that meet contractual requirements while maintaining a focus on service delivery.

Enabling science is the cornerstone of scientific computing at Jefferson Lab. CST supports multiple use cases in nuclear physics, with extensive real-life experience processing data intensive workloads from data acquisition through accepted publications. CST is integral to CEBAF science operations and accountable to the user community for efficient 24/7 operations to support acquisition of raw data at high rates.

Jefferson Lab's leadership in scientific computing and data science in the nuclear physics community provides a conduit to connect nuclear physics experimental use cases and HPC, analogous to the role Berkeley Lab plays in connecting high-energy physics use cases to HPC facilities.

Jefferson Lab continues to innovate in computing and software for experimental workloads. This work has taken place in collaboration with the Lab's Experimental Nuclear Physics and Theory Divisions as part of a commitment to deliver working solutions for scientific discovery. Jefferson Lab's track record in innovation was recognized with a core competency in Advanced Computational Science in 2022.

### 5.2.2  Lawrence Berkeley National Laboratory

Berkeley Lab is a multi-program science and technology laboratory managed for the DOE SC by the University of California. Scientists and engineers at Berkeley Lab conduct basic and applied research to create scientific knowledge and technological solutions that strengthen the nation's leadership in key areas of science, including renewable, clean and efficient energy, energy storage, Earth systems, materials science, chemistry, AI, fusion energy, quantum information systems, and high-energy and nuclear physics. Berkeley Lab had a leading role in creating the modern-day model of team science. This is exemplified by the numerous national user facilities it operates, including the Advanced Light Source, ESnet, the Joint Genome Institute, and NERSC.

**NERSC**

As the mission high-performance computing and data facility for the SC, NERSC provides computing and data analysis resources for scientists and will celebrate its 50th anniversary in 2024. The center hosts supercomputers, large data storage systems, and edge systems, providing expert technical consulting and support services to more than 10,000 users associated with more than 1,000 projects. In addition to supporting simulation-based computing, it partners with a number of large experimental facilities to provide for their extreme scale data processing and analysis needs.

**ESnet**

Berkeley Lab also manages and operates the Energy Sciences Network, the high-performance national networking facility of the DOE SC. ESnet is a high-performance unclassified network built to support scientific research. It provides services to more than 50 DOE research sites, including the entire National Laboratory system, its supercomputing facilities, and its major scientific instruments. ESnet also connects to 140 research and commercial networks, permitting DOE-funded scientists to collaborate productively with global partners.

**Science User Facilities**

In addition to its ASCR facilities, Berkeley Lab is home to three science user facilities:

- The **Advanced Light Source** is a specialized particle accelerator, known as a synchrotron light source, that generates bright beams of X-ray, infrared, and extreme ultraviolet light useful for research. It supports the research of 2,000 users annually.

- The **Joint Genome Institute** advances genomics science in support of the DOE's clean energy and environmental missions by providing scientific users from around the world access to integrated, high-throughput gene sequencing, DNA design and synthesis, metabolomics, and computational analysis.

- The **Molecular Foundry** is a DOE-funded nanoscience research facility that gives scientists from around the world access to world-class expertise and instrumentation in a collaborative, multidisciplinary environment.

## 5.3    Acronyms

| | |
|---|---|
| AI | Artificial Intelligence |
| ALCF | Argonne Leadership Computing Facility |
| ASCR | Advanced Scientific Computing Research |
| CD | Critical Decision |
| CEBAF | Continuous Electron Beam Accelerator Facility |
| CIO | Chief Information Officer |
| DOE | U.S. Department of Energy |
| ESnet | Energy Sciences Network |
| FAIR | Findability, Accessibility, Interoperability, and Reusability |
| HPC | High Performance Computing |
| HPDF | High Performance Data Facility |
| iHARP | Institute for Harnessing Data and Model Revolution in the Polar Regions |
| IRI | Integrated Research Infrastructure |
| JLDC | Jefferson Lab Data Center |
| LBNL | Lawrence Berkeley National Laboratory (Berkeley Lab) |
| LEED | Leadership in Energy and Environmental Design |
| ML | Machine Learning |
| NERSC | National Energy Research Scientific Computing Center |
| OLCF | Oak Ridge Leadership Computing Facility |
| OSTI | Office of Scientific and Technical Information |
| SC | Office of Science |
| SLAC | Stanford Linear Accelerator Center |
| TJNAF | Thomas Jefferson National Accelerator Facility (Jefferson Lab) |
| UX | User Experience |

## 5.4    Glossary

**Advanced Scientific Computing Research**. The ASCR mission is to discover, develop, and deploy computational and networking capabilities to analyze, model, simulate, and predict complex phenomena important to the DOE.

**Critical Decision**. Formal transition points during a project's life cycle where a set of required deliverables are evaluated by approvers to ensure they were properly completed and accepted.

- CD-0, Approve Mission Need
- CD-1, Approve Alternative Selection and Cost Range
- CD-2, Approve Performance Baseline
- CD-3, Approve Start of Construction/Execution
- CD-4, Approve Start of Operations or Project Completion

**Cybersecurity and Federated Access**. Users require a distributed research infrastructure with seamless access and consistent services, while the infrastructure must be operated according to cybersecurity requirements and policies set at the federal level. Operators of user facilities also have different missions, and thus different requirements, across the lab complex. Balancing these constraints can also lead to sources of impedance. Novel secure design patterns and architectures will be required to support open science-integrated architecture for seamless scientific collaboration [2].

**Data Stewardship**. Managing the full life cycle of scientific data requires the adoption and promotion of best practices in data curation, repository management, discoverability and reuse, FAIR/open science, and trusted data. HPDF will play a unique role, providing functionality that enables scientific data to be leveraged by all users and AI/ML techniques to be applied optimally. HPDF will provide an interface to data as well as the monitoring infrastructure to track how data are queried, accessed, and moved through the system, providing provenance and security checks to ensure integrity. HPDF will also work with users to ensure data privacy and ethical use.

**Data Transport and HPDF Integration**. HPDF integration technology will build upon ESnet's existing hardware and programmatic infrastructure to develop a data distribution substrate that leverages ESnet6 infrastructure and its presence at national laboratories and research institutions. The reliable, low-latency, high-bandwidth data path will simplify data transfers, provide dynamic computational load-balancing capabilities, and improve efficiency of scientific workflows.

**Distributed Architecture**. Aligns with the IRI's goal to "empower researchers to seamlessly and securely meld DOE's world-class research tools, infrastructure, and user facilities in novel ways to radically accelerate discovery and innovation." Jefferson Lab and Berkeley Lab will provide management and operational support for HPDF, focusing on deployment and operation of Hubs, coordination of national federated architecture operations, and development of strategic directions for HPDF and IRI efforts in partnership with DOE ASCR and other stakeholders.

**Distributed Data Storage**. HPDF will coordinate and add resources to HPDF distributed data storage. A federated meta catalog and object store will provide access to the rich metadata needed for FAIR data management. Tiered storage will support local data processing and archiving.

**Early Access Facility Testbed**. Using existing data center space and infrastructure capacity to implement an initial hardware platform. As the HPDF Project matures, we will implement a single standard unit of Hub design with a subset of capabilities to support select use cases (e.g., data-

intensive and time-critical processing). Our partnerships will inform the design and support development of prototypical Spoke integration concepts for science programs. These partners will also assist in evaluation of potential approaches for Spoke facilities.

**Exascale Computing**. Exascale supercomputers are the current generation of supercomputers, allowing scientists to better simulate the complex processes involved in stockpile stewardship, medicine, biotechnology, advanced manufacturing, energy, material design and the physics of the universe, more quickly and with higher definition. The DOE's National Labs have some of the most significant HPC resources available, including some of the world's fastest supercomputers.

**High Performance Computing** (HPC). The most powerful and largest-scale computing systems, which enable research that would otherwise be impractical or impossible in the real world. Within the federal government, the DOE leads the effort of pushing the boundary of what is possible, with the nation's fastest and most capable supercomputers housed at National Laboratories.

**High Performance Data Facility**. HPDF is envisioned as a national resource that will serve as the foundation for advancing DOE's IRI program, enabling and accelerating scientific discovery by delivering state-of-the-art data management infrastructure, capabilities, and tools.

**HPDF Hub**. The central component of HPDF, providing services for data access, data stewardship, data manipulation, and user community support. The Hub will encompass the hardware, software, and staff needed to organize, orchestrate, and provide resources for HPDF.

**Integrated Research Infrastructure** (IRI). A new DOE program intended to give researchers frictionless access to a seamlessly integrated ecosystem of computing, networking, instruments, and experimental facilities.

**SC User Facility**. An Office of Science user facility is a federally sponsored research facility available to advance scientific or technical knowledge.

**Spokes.** Satellite facilities that operate in concert with the HPDF Hub. As HPDF components, they will enable delivery of support and data services in a way that reduces latency and improves user support and overall system resilience. HPDF participants, including Spokes, computational Hubs, users, and other sites, will form the integrated data ecosystem within DOE.

**Standard Units**. Individual tightly coupled clusters of CPU, GPU, memory, and storage, meshed by a high-bandwidth, dynamically steered network. Highly flexible; can be designed to operate individually or collectively in high-throughput and high-performance computing workloads, and streaming data processing modes.

## 5.5  Bibliography

1. National Research Council. 1999. Cooperative Stewardship: Managing the Nation's Multidisciplinary User Facilities for Research with Synchrotron Radiation, Neutrons, and High Magnetic Fields. Washington, DC: The National Academies Press. https://doi.org/10.17226/9705.

2. Miller, William L.; Bard, Deborah; Boehnlein, Amber; Fagnan, Kjiersten; Guok, Chin; Lançon, Eric; Ramprakash, Sreeranjani; Shankar, Mallikarjun; Schwarz, Nicholas; and Brown, Benjamin L. 2023. Integrated Research Infrastructure Architecture Blueprint Activity (Final Report 2003). United States. https://doi.org/10.2172/1984466.

3. ASCR Integrated Research Infrastructure Task Force. Toward a Seamless Integration of Computing, Experimental, and Observational Science Facilities: A Blueprint to Accelerate Discovery. 2021. https://doi.org/10. 2172/1863562.

4. Helland, Barbara. Future of Computational Infrastructures: Exascale Computing and an Integrated Research Infrastructure. [Online] 2022. https://science.osti.gov/-/media/bes/besac/pdf/202212/7-Helland-BESAC-Panel.pdf.

5. Alexander, Francis J., et al. Co-design Center for Exascale Machine Learning Technologies (ExaLearn). 2021, The International Journal of High Performance Computing Applications, Vol. 35, pp. 598-616.

6. An Implementation Plan for a National Artificial Intelligence Research Resource. [Online] https://www.ai.gov/wp-content/uploads/2023/01/NAIRR-TF-Final-Report-2023.pdf.

7. Wilkinson, MD, Dumontier, M, and Aalbersberg, IJ. FAIR Guiding Principles for Scientific Data Management and Stewardship. Sci. Data 3 (2016). 2016, Scientific data, Vol. 3.

8. Nelson, Alondra. Ensuring Free, Immediate, and Equitable Access to Federally Funded Research. [Online] 2022. https://www.whitehouse.gov/wp-content/uploads/2022/08/08-2022-OSTP-Public-Access-Memo.pdf.

9. IRI: The Integrated Research Infrastructure Update to ASCAC. May 20, 2024.

10. Program and Project Management for the Acquisition of Capital Assets – DOE Directives, Guidance, and Delegations.