

January 26, 2024

RFP Technical Requirements Document

for

NERSC-10 System

Version 6.0

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

TABLE OF CONTENTS

REQUIREMENTS DEFINITIONS	4
1.0 INTRODUCTION	5
1.1 NERSC-10 Mission Need	5
1.2 HPC Workflows	6
1.3 Schedule	7
2.0 HIGH-LEVEL SYSTEM REQUIREMENTS	8
2.1 System Description	8
2.2 Software Description	9
3.0 BENCHMARKS	10
4.0 WORKFLOW ENVIRONMENT	12
4.1 Scalable and Reliable Workflow Services	12
4.2 Compilers and Libraries	13
4.3 Python and Artificial Intelligence/Machine Learning (AI/ML)	14
4.4 Profiling and Debugging	14
4.5 Containers	15
4.6 Workflow Readiness Support	15
4.7 Programming the Data Center	16
5.0 SYSTEM SOFTWARE & MANAGEMENT	17
5.1 Infrastructure Services	17
5.2 Operating System	18
5.3 Platform Management	18
5.4 System Software Deployment	19
5.5 Data Collection and Monitoring	19
6.0 SYSTEM NETWORKS	20
7.0 STORAGE SYSTEMS	22
7.1 Platform Storage System (PSS)	22
7.2 QoS Storage System (QSS)	23
8.0 SYSTEM OPERATION	24
8.1 Resilience, Reliability and Availability Metrics	24

8.2 System Security	25
8.3 Power and Energy	26
8.4 Maintenance and Support (Hardware/Software)	26
8.5 Documentation	27
9.0 FACILITIES AND SITE INTEGRATION	27
10.0 NON-RECURRING ENGINEERING (NRE)	28
11.0 TECHNICAL OPTIONS	28
11.1 Upgrades, Expansions, and Additions	28
11.2 Early Access Systems	30
11.3 Test Systems (TS)	31
12.0 DELIVERY AND ACCEPTANCE	31
13.0 PROJECT AND RISK MANAGEMENT	31
APPENDIX A: Sample Acceptance Test Plan	35
APPENDIX B: Project and Risk Management - Planning Deliverables Descriptions	40
APPENDIX C: Facility and Site Integration Specifications	47
DEFINITIONS AND GLOSSARY	55

REQUIREMENTS DEFINITIONS

Technical requirements have priority designations, which are defined as follows:

(a) Target Requirements designated as (TR-1, TR-2, or TR-3)

Target Requirements (designated TR-1, TR-2, or TR-3) are features, components, performance characteristics, or other properties that are important to the University. Target Requirements are prioritized by dash number. TR-1 is most desirable to the University and forms the baseline system, while TR-2 is desirable and adds additional capabilities or increases productivity. TR-3s are stretch goals. Target Requirement responses will be evaluated as part of the proposal evaluation process.

(b) Technical Option Requirements designated as (TO-1, TO-2, or TO-3)

Technical Option Requirements (designated TO-1, TO-2, or TO-3) are features, components, performance characteristics, or upgrades that are important to the University. Technical Options add value to a proposal. Technical Options are prioritized by dash number. TO-1 is most desirable to the Laboratory, while TO-2 is more desirable than TO-3. Technical Option responses will be evaluated as part of the proposal evaluation process; however, the University may or may not elect to include Technical Options in the resulting subcontract(s). Each proposed TO should appear as a separately identifiable item in an Offeror's proposal response.

Note: There are no mandatory requirements or mandatory options in the NERSC-10 technical requirements document.

1.0 INTRODUCTION

The Regents of the University of California (the “University”), which operates the National Energy Research Scientific Computing (“NERSC”) Center residing within Lawrence Berkeley National Laboratory (“LBNL”), is releasing a Request for Proposal (RFP) for the next-generation high performance computing (HPC) system, NERSC-10, to be delivered in the 2026-time frame.

The successful NERSC-10 Offeror will be responsible for delivering, installing, supporting, and maintaining the NERSC-10 system.

Each proposed solution in response to this document should clearly describe the role of any lower-tier subcontractor(s) and the technology or technologies, both hardware and software, and value-added that the lower-tier subcontractor(s) provide(s).

The statement of work for any subcontracts resulting from this RFP will be based on this Technical Requirements Document and the successful Offeror’s responses/proposed solutions.

NERSC-10 has maximum funding limits over its system life, to include all design and development, site preparation, maintenance, support, and analysts.

Application performance and workflow efficiency are essential to NERSC. Success will be defined as meeting NERSC-10 mission needs as described in Section 1.1, below. The advanced workflows aspects of the NERSC-10 system will be pursued both by fielding first-of-a-kind workflow-enabling technologies as part of the system and by selecting and participating in strategic Non-Recurring Engineering (NRE) projects with the Offeror and applicable technology providers. A compelling set of NRE projects will be crucial for the success of NERSC-10 by enabling the deployment of first-of-a-kind technologies in such a way as to maximize their utility.

Supporting information can be found on the NERSC-10 website (<https://www.nersc.gov/systems/nersc-10/>).

Additional information on proposal preparation are provided in Section 2.4, the *Proposal Submittal Requirements* Section of the RFP.

1.1 NERSC-10 Mission Need

The DOE Office of Science (SC) is the largest supporter of basic and applied research programs in the areas of efficient energy use, reliable energy sources, improved environmental quality, and fundamental understanding of matter and energy. One of the principal thrusts within SC is the direct support of the development, construction, and operation of unique, open-access HPC scientific user facilities. These [HPC facilities](#) are critical to supporting the research programs that help accomplish the DOE’s mission.

The NERSC User Facility at LBNL, funded by the DOE SC’s [Advanced Scientific Computing Research](#) (ASCR) Office, is the mission HPC facility for SC, uniquely supporting the needs of science across the entire Office. ASCR’s mission is to discover, develop, and

deploy computational and networking capabilities to analyze, model, simulate, and predict complex phenomena important to the DOE. ASCR has a long history of supporting cutting-edge research in applied math, computational and computer science, and the deployment of advanced HPC and networking facilities.

For almost 50 years, NERSC has supported SC-funded researchers using cutting-edge supercomputers to develop materials and devices for clean energy generation and storage, study the environmental impact of a changing global ecosystem, investigate the fundamental properties and interactions of matter, and explore other science areas within the DOE science mission. Potential offerors are encouraged to review [NERSC's history](#) available on the NERSC website. Over that time, NERSC has helped guide the SC computational science community through many disruptive changes and evolving national scientific priorities.

The SC community is now witnessing another transition, with new science use-cases requiring more dynamic and programmable systems to accommodate increasingly complex workflows that may require running many interdependent simulations and/or analysis tasks and integrating simulations with artificial intelligence (AI). These workflows can involve moving vast amounts of complex data among storage hierarchies both within NERSC and externally, with requirements to provide time-sensitive feedback to experiments.

For SC to fulfill its mission to maintain and extend U.S. leadership in scientific discovery, the NERSC-10 system must leverage available new technologies and support the emerging needs in AI and experimental/observational science by accelerating end-to-end DOE SC workflows and enabling new modes of scientific discovery through the integration of experiment, data analysis, and simulation.

A NERSC-10 system upgrade is essential to meet SC goals and national initiatives and to avoid creating a gap between SC programmatic needs and HPC capabilities. To meet the Mission Need, the NERSC-10 system must accomplish the following objectives:

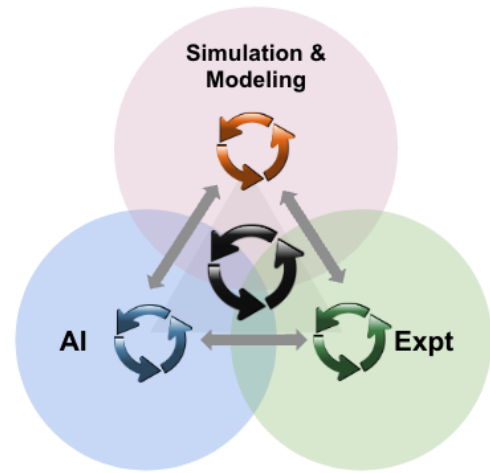
- Deliver significantly enhanced capabilities beyond the current generation [Perlmutter](#) system in the 2026 timeframe
- Maximize the performance of DOE SC software applications and next-generation workflows while supporting the needs of the diverse and broad NERSC user community
- Maximize the effective and efficient transition of the NERSC user community to advanced hardware architectures
- Continue to provide high-speed access to DOE SC community data stored in NERSC's Community File System and archival file system
- Fit within the space confines and the power and cooling envelope of the NERSC data center located at Shyh Wang Hall on the LBNL campus.

1.2 HPC Workflows

The NERSC workload is increasingly diverse, with a growing demand for complex high-performance workflow capabilities from its user community. NERSC supports and will continue to support HPC science campaigns through the following objectives:

- High-performance simulation and modeling workflows (e.g., large-scale multi-physics applications)
- High-performance AI workflows (e.g., training, inference, hyperparameter optimization)
- Cross-facility workflows: Rapid data analysis and real-time steering of experiments (Expt)

The new HPC workflows couple the computational tasks and data flow within, between, or across all three modes shown in the figure to unlock opportunities for new scientific discovery. Enabling these complex and novel workflows requires increases in high- performance computing, composable resources for targeted workload performance, modular workflow execution through orchestration frameworks and application programming interfaces (APIs). It also requires spanning resources into the commercial cloud, other computational facilities, and the edge, tight integration of system components to enable “seamless” execution, leveraging data stored/collected at other facilities, and more.



The NERSC-10 system must enable numerous simultaneous workflows with different resource and time requirements. NERSC is interested in technologies that can help support this vision. The [NERSC-10 Workflows Whitepaper](#) provides workflow scenarios that we expect to support with the NERSC-10 system.

1.3 Schedule

The following is the tentative schedule for the NERSC-10 system acquisition.

Table 1.1. NERSC-10 summary schedule.

Solicitation Schedule	
Request for Proposal Release	Q1 CY, 2024
Offeror Questions Due	Q1 CY, 2024
Proposals Due	Q1 CY, 2024
Award and Performance Schedule	
Early Access System (Pilot and/or Phase I)	CY 2025
NERSC-10 System Delivery	2H CY 2026
NERSC-10 System Acceptance	CY 2027

2.0 HIGH-LEVEL SYSTEM REQUIREMENTS

This section describes the high-level technical requirements for the NERSC-10 system proposals. In addition, per Section 2.4.2.2-Technical Volume Enclosures of the RFP, the

Offeror should complete and submit the Technical Data Summary Sheet and the NERSC-10 Benchmark Results Worksheet.

2.1 System Description

2.1.1 The NERSC-10 system should be sited at the NERSC data center in Building 59 (Shyh Wang Hall) on the LBNL campus in Berkeley, California. The Offeror should provide details of the physical footprint of the system and all of the supporting components to be sited at the NERSC data center. [TR-1]

2.1.2 The Offeror should provide a detailed full-system architectural description of the NERSC-10 system that will deliver at least a 10x performance improvement of the Workflow-Sustained System Improvement (Workflow-SSI) metric over [Perlmutter](#) using any combination of baseline, ported, or optimized performance results, as defined in Section 3.0, below. The Offeror should describe how the system fits into the Offeror's long-term product roadmap. The Offeror's description should include quantities and define any minimum scalable unit sizing to maintain optimal performance and productivity across the system. The Offeror's full-system description should include diagrams and text describing the following details as they pertain to the Offeror's proposed system architecture(s), plus detail any unique features in the design: [TR-1]

- **Component architecture** – details of all processor(s), memory technologies, storage technologies, network interconnect(s), and any other applicable components
- **Compute node architecture(s)** – details of how components are combined into the node architecture(s). Details should include bandwidth and latency specifications (or projections) between components. Details should be provided for each compute node type. The balance of CPU-only to GPU nodes should be chosen to optimize the Workflow-SSI described in Section 3.0.
- **Board and/or blade architecture(s)** – details of how the node architecture(s) is integrated at the board and/or blade level. Details should include all inter-node and inter-board/blade communication paths and any additional board/blade level components.
- **Rack and/or cabinet architecture(s)** – details of how board and/or blades are organized and integrated into racks and/or cabinets. Details should include all inter rack/cabinet communication paths and any additional rack/cabinet level components.
- **Interconnect** – details of the system's high-speed network topology and connectivity across all system components (compute nodes, workflow environment nodes, Platform Storage, QoS Storage, management system).
- **Storage systems** – details of how the Platform Storage System (PSS) and QoS Storage System (QSS) are integrated with the system, including an architectural diagram and gateway nodes if applicable.
- **System architecture** – details of how racks or cabinets are combined to produce system architecture, including the high-speed interconnects and network topologies (if multiple) and storage systems.

- **Management node(s)** – details of hardware to support management and services to operate the NERSC-10 system. Management node types can include, but are not limited to, administrative nodes for the orchestration of system services, worker nodes for the deployment of services (e.g., the Slurm resource manager), and storage nodes for system management and administration. Multiple node types may be needed to optimize for different uses described in Section 5.0.
- **Workflow environment node(s) (WEN)** – details of hardware to support user access and user-driven workflow activities. A pool of WENs will be needed to address the different requirements described in Section 4.0, below. NERSC-10 may require multiple WEN types optimized for different use cases. At least one WEN type should be accessible even if the compute system is unavailable. The workflow environment nodes will collectively support 800 interactive users, 80 data analysis front-end applications, 5000 batch jobs, and 320 simultaneous compilations of software of equivalent complexity to the latest GNU Compiler Suite.

2.1.3 The Offeror should provide an alternative processor vendor and technology. The response shall provide a concise architectural description, including hardware and software. The response should include a description of the PCI-e technology generation, BabelStream and py-GEMM microbenchmark results, and power estimates (see Section 3 for benchmark descriptions). [TR-1]

2.1.4 The Offeror should describe any other technologies that are part of their NERSC-10 offering that are high risk and propose mitigation strategies, including alternative technologies. [TR-1]

2.2 Software Description

2.2.1 The Offeror should provide a detailed description of the proposed system software and user programming environment, including a high-level software architecture diagram. The Offeror should describe the provenance of the software component (for example, open source or proprietary), support mechanisms and level of support, and licensing, including any limitations on the simultaneous use of software, if applicable (for the lifetime of the system, including updates). [TR-1]

2.2.2 The Offeror should describe the high-level roadmap for the following: [TR-1]

- System software and tools provided for management and operation of the NERSC-10 system
- Provided user programming environment, including the ability to utilize new hardware features

3.0 BENCHMARKS

Assuring that real applications and workflows perform well on the NERSC-10 system is key to the success of the system. The workflow component benchmarks listed in Table 3.1, below will be used to evaluate workflow performance as part of both the RFP response and system

acceptance. These benchmarks demonstrate the most computationally intensive components of the workflows they represent, but without the data-flow and control-flow complexities of an integrated science workflow. The workflow benchmarks are supplemented by a collection of microbenchmarks listed in Table 3.2.

Benchmark performance targets will be negotiated after a final system configuration is defined. Performance targets will be reevaluated and converted into requirements at the technical decision point (TDP) described in Section 13.0, below. All performance tests must continue to meet acceptance criteria throughout the lifetime of the system.

The benchmarks and supplemental materials can be found on the NERSC-10 benchmarks website: <https://www.nersc.gov/systems/nersc-10/benchmarks>. All benchmark results must be reported in the accompanying “NERSC-10 Benchmark Results” worksheet. All benchmark results should conform to the Workflow Component Benchmark and Microbenchmark Instructions and Run Rules document (<https://gitlab.com/NERSC/N10-benchmarks/run-rules-and-ssi>).

Benchmark results and Workflow-SSI can be submitted for three categories of workflow optimization (baseline, ported, and optimized), which are defined in the Run Rules document. Workflow-SSI is the calculation used for measuring improvement and is documented on the NERSC-10 benchmarks website.

Table 3.1. Workflow Component Benchmarks

Workflow Name	Description	Application Components
Lattice QCD	Lattice Quantum Chromodynamics (QCD)	MILC generation MILC analysis
Optical Properties of Materials	<i>Ab initio</i> Electronic Structure	BerkeleyGW Epsilon BerkeleyGW Sigma
Materials by Design	Molecular Dynamics	LAMMPS
Climate Simulation and Analysis	Deep Learning Training	DeepCAM
Metagenome Analysis	Genomic Data Analysis	HMMsearch
CMB-S4	Cosmology Data Analysis	TOAST-3

- 3.0.1 The Offeror should provide baseline performance results for the proposed system and platform storage system for all of the workflow component benchmarks listed in Table 3.1. [TR-1]
- 3.0.2 The Offeror should provide ported performance results for the proposed system and platform storage system for any of the workflow component benchmarks. [TR-2]

- 3.0.3 The Offeror should provide optimized performance results for the proposed system and platform storage system for any of the workflow component benchmarks. [TR-2]
- 3.0.4 The Offeror should state a minimum Workflow-SSI for the NERSC-10 system, to be measured using baseline versions of the workflow component benchmarks. If baseline results cannot be obtained, ported results may be provided in their place. [TR-1]
- 3.0.5 The Offeror should state a minimum Workflow-SSI for the NERSC-10 system, to be measured using any combination of baseline, ported, or optimized versions of the workflow component benchmarks. [TR-2]
- 3.0.6 The Offeror should provide baseline performance results for the proposed compute system and QoS Storage System for the DeepCAM workflow component benchmark. [TR-2]
- 3.0.7 The Offeror should provide performance results for the proposed system for microbenchmarks listed in Table 3.2. Results should be provided for each test configuration listed in the table. [TR-1]
- 3.0.8 The Offeror should provide licenses for the NERSC-10 system for all software required to achieve benchmark performance, including but not limited to compilers and libraries. [TR-1]

Table 3.2. Microbenchmarks

Name	Description	Test Configuration(s)
BabelStream	Memory bandwidth	Every functional combination of processor type and memory domain
py-DGEMM	Floating-point performance	Every processor type
Ziatest	Job startup	Full-system
OSU Micro-Benchmarks (OMB)	MPI performance	Every functional combination of processor type and memory domain. It is not necessary to test communication between heterogeneous processor/ memory combinations
iperf	External networking performance	Through each path into the HPC network (e.g., between compute nodes and the data center Ethernet network or between compute nodes and off-platform file systems)
IOR	Storage bandwidth performance	PSS and QSS - single node to minimum number of nodes to achieve maximum performance; POSIX/MPI-IO; strided/random; read/write QSS - minimum scalable unit QSS - multiple independent scalable units
MDTest	Storage metadata performance	PSS and QSS - multiple levels of concurrency (single MPI process, optimal MPI processes/compute node, minimum MPI process to achieve peak performance, maximum concurrency on

		the proposed system) QSS - minimum scalable unit QSS - multiple independent scalable units
--	--	--

4.0 WORKFLOW ENVIRONMENT

4.1 Scalable and Reliable Workflow Services

- 4.1.1 The system should support running jobs up to the full scale of the compute node resources. The Offeror should describe factors (such as executable size) that may affect application launch time. [TR-1]
- 4.1.2 SchedMD's Slurm resource job management scheduler will be the primary scheduler and policy engine of the system. The University will directly procure the necessary software licenses and ongoing maintenance support from SchedMD. The Offeror will cooperate and work with the University and/or SchedMD to resolve operational problems with Slurm that may be caused by the Offeror's products. The Offeror will provide the necessary integration interfaces to support scalable job launch, including node placement, topology-aware scheduling, rank reordering, power-aware scheduling, and node configuration and re-provisioning of nodes if supported by the hardware. The system design should not limit Slurm's ability to support thousands of concurrent users and more than 20,000 concurrent batch jobs. [TR-1]
- 4.1.3 The system should support a container orchestration platform such as Kubernetes or similar to provide staff-managed and users' self-supported services on the workflow environment nodes. It should be capable of operating with Slurm to provide a unified workflow environment where users can securely and performantly launch job tasks on the compute resources from the workflow environment nodes or services running on them. The Offeror should describe any specialized hardware or software that may enhance this unified workflow environment capability. [TR-2]
- 4.1.4 The system workflow environment nodes should support interactive user access modes, including the following: [TR-1]
- command-line interface (CLI) through ssh and web-based user access modes for login, code compilation (cross-compilation is not desirable), application development, container builds, job lifecycle management, small-scale data analysis, and data transfer
 - persistent user services and frameworks (e.g., JupyterHub, databases, API services, and message brokers) that are staff-managed or user self-supported

The Offeror should describe any provided mechanisms to enable these access models and how they are managed.

- 4.1.5 The system should provide correct numerical results and minimize runtime variability. The Offeror should describe strategies for minimizing runtime variability in production and the method used to measure variability (e.g., the measured coefficient

of variation (standard deviation divided by the mean) of runtimes from dedicated/non-dedicated runs is not greater than X%). [TR-1]

4.2 Compilers and Libraries

- 4.2.1 The system should support building and executing C17 code, C++20 code, and Fortran 2018 code, including code utilizing OpenMP directives 5.2 or later. The Offeror should describe all provided compilers, including any enhancements or limitations that can be expected in meeting full support of the standards and other native language features for expressing parallelism, including but not limited to support for C++ parallel STL, Fortran *do concurrent*, and Fortran coarrays. The Offeror should describe the level of support, if any, for all provided compilers. [TR-1]
- 4.2.2 The Offeror should describe any support for LLVM backends for each processing element (e.g., CPU, GPU, specialized accelerator), such that it can be utilized with both Offeror-provided frontends and the open clang/flang projects. [TR-2]
- 4.2.3 The Offeror should describe any provided capability for the system to compile and run applications using Kokkos, SYCL, and/or OpenACC 3.x. If applicable, the Offeror should describe any provided performance enhancements. [TR-3]
- 4.2.4 The Offeror should describe any provided capability for the system to compile and run CUDA-based applications. The Offeror should describe any tools or software to support the translation of CUDA applications to the native GPU programming language if CUDA is not directly supported. [TR-1]
- 4.2.5 The Offeror should provide a Message Passing Interface (MPI) library that supports MPI 4.0 or higher, is capable of running a job at full-system scale, and makes GPU-aware MPI available wherever this is supported by the GPU vendor. The Offeror should describe any extensions or limitations to the MPI standard in the available MPI libraries. [TR-1]
- 4.2.6 The Offeror should enable the system to support the Process Management Interface (PMI). The Offeror should describe the version and supported integrations. [TR-1]
- 4.2.7 The Offeror should describe any provided optimized BLAS, LAPACK, ScaLAPACK, and FFT libraries for CPU and GPU. [TR-1]
- 4.2.8 The Offeror should describe any provided communication libraries (e.g., PGAS and task-based programming libraries). [TR-3]
- 4.2.9 The Offeror should describe any provided optimized scientific I/O libraries for CPUs and GPUs (e.g., HDF5, NetCDF). [TR-3]
- 4.2.10 The Offeror should describe any provided mechanisms for direct data movement and access from the GPU to improve network bandwidth and latency (GPU-to-GPU) and I/O performance (GPU-to-Storage). [TR-3]

4.3 Python and AI/ML

- 4.3.1 The Offeror should ensure any provided libraries (used to expose high-performance CPU, GPU, I/O, or communication capabilities in 4.2.5-4.2.9) can be used through standard Python packages (e.g., NumPy, SciPy, h5py, mpi4py) or using standard Python packaging tools. The Offeror should describe any limitations. [TR-2]
- 4.3.2 The Offeror should describe any provided performance optimizations for distributed or parallel execution of Python programs through Python multiprocessing, Dask, Ray, or other similar non-MPI-based runtimes. [TR-2]
- 4.3.3 The Offeror should describe any provided optimized libraries for execution of machine learning and AI workloads for both CPUs and GPUs, such as those required for optimal execution of deep learning frameworks like PyTorch and TensorFlow. [TR-1]
- 4.3.4 The Offeror should describe any provided distributed deep learning libraries (e.g., Horovod, PyTorch DDP) that enable scaling of training workloads across the full system and any provided tools that accelerate AI/ML and Deep Learning-based workflows (e.g., hyperparameter optimization, tracking experiments, and integration with simulation and data pipelines). [TR-2]

4.4 Profiling and Debugging

- 4.4.1 The Offeror should describe all provided profiling tools and their scaling capabilities. The Offeror should describe support for MPI and OpenMP profiling for CPU and GPU for all provided compilers. [TR-1]
- 4.4.2 The Offeror should provide support for APIs, including Linux perf, that enable profilers and other performance optimization tools to access CPU and GPU performance counters on the system. The Offeror should describe any restrictions on perf_event_paranoid, required kernel modules, or other security considerations. [TR-1]
- 4.4.3 The Offeror should describe all provided debugging tools for applications running on all user-accessible hardware, such as gdb for CPU and equivalents for GPU. The Offeror should describe capabilities for all supported languages and any limitations to scaling up to 10% of the system. [TR-1]

4.5 Containers

- 4.5.1 The Offeror should provide a mechanism for users to build and run [Open Container Initiative \(OCI\)](#)-compliant containers on the system without requiring privileged access to the system or allowing a user to escalate privilege. [TR-1]
- 4.5.2 The Offeror should describe how the system software and hardware dependencies (e.g., device driver libraries, MPI libraries, and libfabric) can be accessed by containers, including dynamic mechanisms to maintain accessibility of these dependencies when software updates are made. [TR-1]

- 4.5.3 The Offeror should describe any provided container images for users (e.g., libraries, applications), the licensing model and how they can be distributed (e.g., can we distribute a container build on top of an Offeror-provided container), and the image registry where these container images may be published (e.g., an internal or public registry). [TR-2]
- 4.5.4 The Offeror should describe any slowdowns and scaling limitations that would be observed due to running an application in a container up to the full scale of the system. [TR-2]

4.6 Workflow Readiness Support

- 4.6.1 The Offeror should provide a separately priced Workflow Readiness Support (Center of Excellence) plan to assist in transitioning select NERSC applications to the system (e.g., [NESAP](#) focus for simulations, data, and learning). Support should be provided by the Offeror and the CPU and GPU vendor. The Offeror should provide access to experts in the areas of compilers and application performance in the form of staff training and deep-dive interactions with a set of teams. The proposed plan will be used to mutually develop a Center of Excellence (COE) for Workflow Readiness Plan as described in Appendix B.10. [TR-1]
- 4.6.2 The Offeror should include support for transitioning select workflows to the system. Support should be provided by the Offeror and/or key technology providers (e.g., the CPU and GPU vendors, storage, networking, third-party software) to address overall workflow performance. The Offeror should include how it will collaborate with third-party developers from open-source communities. [TR-2]
- 4.6.3 The Offeror should propose a separately priced Workflow Readiness Support (Training and Education Plan) available both during the COE and for the lifetime of the system. Activities should target effective use of the user environment, performance, and optimization. The description should include topics, frequency, and format (e.g., classroom training, online training, hackathons). Proposed training will be used to mutually develop a user Training and Education Plan as described in Appendix B.11. [TR-1]

4.7 Programming the Data Center

- 4.7.1 The system should support complex workflows, as described in the NERSC-10 Workflows Whitepaper, through REST API interfaces or other mechanisms that expose functionality to system administrators and/or users to automate services. The Offeror should describe the capabilities that its provided REST APIs expose and detail how they are documented and tested. This description may include but is not limited to the following capabilities: [TR-1]
- System and subsystem status and health
 - Data transfer, management, and archiving
 - Orchestration of workflows, persistent services, CI/CD workflows (including container deployment), and complex science workflows

- Dynamic reconfiguration of storage, compute, and networking hardware
 - Any authentication and authorization requirements/expectations for their REST APIs
- 4.7.2 The Offeror should describe any provided capabilities that enable or improve multi-tenancy support (on compute and/or WEN nodes) that goes beyond the status quo (shared node jobs). This description may include but is not limited to the following capabilities: [TR-2]
- Protecting users from one another through minimized privileges and other mitigation techniques to prevent escalations in privileges
 - Virtualization and container networking (e.g., SR-IOV, VXLAN), including details of hardware offload capabilities, the number of tenants supported, and guarantees of isolation between tenants
- 4.7.3 The system should support integrating with external cloud capabilities provided by the Offeror and/or commercial cloud service providers to enable resilient workflows, urgent computing, or specialized services (see the NERSC-10 Workflows Whitepaper). The Offeror should describe any provided support for hybrid-cloud capabilities, including bursting to cloud resources, portability of programming environments, data management and movement, and access to specialized hardware and services. [TR-2]
- 4.7.4 The Offeror should describe any provided capabilities to support “server-less” or Function-as-a-Service (see the NERSC-10 Workflows Whitepaper), including how it could integrate with the system and scheduler, security model, scaling, and performance. The Offeror should describe all capabilities to support an event-based message model that can be used to publish and subscribe to system events, job events, data-related events, and other event types that can be integrated into and used to support complex workflows. [TR-3]
- 4.7.5 The Offeror should provide specialized hardware or integrated technologies to enhance support for composable workflows. The hardware could be present in the WEN and/or compute partitions of NERSC-10. The description should outline the programmable capabilities, the intended scope (e.g., user programmable versus a secure platform managed by trusted identities), the available memory, processing hardware, interfaces, and hardware acceleration capabilities. [TR-3]
- 4.7.6 The Offeror should provide data lifecycle management capabilities, such as policy-driven data movement and capability to integrate with a site-wide scheduling resource. The Offeror should describe any methods for enabling persistent user-defined metadata to enable tracking and sharing data across the entire storage ecosystem or automatically attaching workflow-based metadata to files. [TR-3]
- 4.7.7 The Offeror should describe how its proposed system could interface with quantum computing hardware, including any provided software and library support. The Offeror should describe any software provided to support quantum computing simulation at scale that will be provided with the NERSC-10 system. [TR-3]

5.0 SYSTEM SOFTWARE & MANAGEMENT

5.0.1 The NERSC-10 system should include management capabilities that facilitate integration with the evolving HPC workflow-driven environment. The management system should include the following:

- Employ configuration management to ensure reproducibility and automation of critical tasks (e.g., continuous deployment of operating system images, container images, microservices on compute nodes, server nodes and any support devices, and reinstallation when necessary for operational reasons)
- Employ software components that should not restrict the evolution of the NERSC-10 workflow ecosystem, comply with open standards when available (e.g., Redfish), and provide documented programming interfaces. NERSC may choose to use open source or third-party software over the course of the production lifecycle for NERSC selected components

The Offeror should provide a high-level overview of its proposed system management solution and any limitations toward achieving a modular environment. [TR-1]

5.0.2 The Offeror should provide access to source code and necessary build environments for all software except for firmware, compilers, and third-party products. The Offeror should provide updates of source code and any necessary build environments for all software over the life of the system. [TR-1]

5.1 Infrastructure Services

5.1.1 The Offeror should provide capabilities and describe how the University would use them to remotely manage compute nodes, network switches, platform and QoS storage, power distribution units, and servers comprising the system, including power control and console access, firmware updates, zero-touch provisioning, diagnostics, event logs, and alert capabilities. These capabilities should be accessible via documented APIs, preferably based on open standards, and a user interface (e.g., graphical or command-line). [TR-1]

5.1.2 The Offeror should provide capabilities and describe how the University would use them to perform scalable full-platform management to automate the management of all hardware, see a comprehensive overview of system operations, and automate whole-system maintenance actions. Relevant features include, but are not limited to, sequenced power up and power down of the system; summarization of temperature, power, and other sensors; automating firmware and configuration updates; maintaining an inventory of field-replaceable units over the system lifetime; and collecting alert and error information from hardware. [TR-1]

5.2 Operating System

5.2.1 The Offeror should provide a commercially-supported, non-proprietary Linux operating system (OS) environment on all visible service partitions (e.g., front-end

nodes, service nodes, I/O nodes). The Offeror should describe the proposed Linux environment. [TR-1]

- 5.2.2 The Offeror should provide an optimized compute partition operating system such that all optimizations are limited to open-source Linux kernel modules that can be rebuilt onsite to provide an efficient execution environment for applications running up to full-system scale. The Offeror should describe any HPC-relevant optimizations made to the compute partition operating system. [TR-1]
- 5.2.3 The Offeror should explain how they will enable all provided device drivers or kernel modules to be rebuildable and manageable by the University with the operating system proposed in 5.2.1. [TR-1]

5.3 Platform Management

- 5.3.1 The Offeror should describe the provided system configuration management and diagnostic capabilities of the system that address the following details of system management: [TR-1]
- Any effect or overhead of software management tool components on the CPU or memory available on compute nodes
 - Support for multiple simultaneous or alternative system software configurations, including the estimated time and effort required to install both a major and a minor system software update
 - User activity tracking, such as audit logging and process accounting
 - Any restrictions to privileged access to all hardware components delivered with the system
- 5.3.2 The system should have no single points of failure that would cause a **system outage** (as defined in the Definitions and Glossary). The system should remain in an operational or degraded state after the unexpected failure or planned maintenance of any single field-replaceable unit (FRU), server, or switch and during any repair or other maintenance action. The Offeror should describe the provided reliability, availability and serviceability (RAS) capabilities to mitigate single points of failure (hardware or software) and the potential effect on running applications and system availability. [TR-1]
- 5.3.3 The Offeror should describe the provided resilience, reliability, and availability mechanisms as well as describe the capabilities of the system to mitigate any condition or event that can potentially cause a job interrupt. The Offeror should detail how a job maintains its resource allocation and is able to relaunch an application after an interrupt. [TR-2]

5.4 System Software Deployment

- 5.4.1 The Offeror should provide mechanisms for the system to perform rolling upgrades and rollbacks on a subset of the system while at least half of the system remains in

production operation. The Offeror should describe the mechanisms and any limitations of the continuous deployment framework to perform rolling upgrades. [TR-1]

- 5.4.2 The Offeror should describe the provided solution for scalable boot, reconfiguring, and rebooting of compute, server, and any other node types in the system. The description should include an overview of the node boot process (warmboot and coldboot), including secure boot, stateless/stateful node provisioning, and infrastructure automation for customization and configuration of a node, the coordination, ordering, and parallelism of the boot process, and techniques to provide rapid configuration and rebooting. The Offeror should include how the time required to reboot scales with the number of nodes being rebooted. [TR-2]
- 5.4.3 The Offeror should describe any provided system development tools to make deployments easier (e.g., container registry, container image management, automated testing, and version control) and describe how it integrates with the system management workflow. [TR-3]
- 5.4.4 The Offeror should describe how the University can add new third-party hardware to the system (e.g., an AI accelerator partition, fabric-attached memory, specialized storage, and a cluster of Linux servers). The Offeror should detail any requirements, interfaces, or standards that must be provided by this third-party hardware for addition to the system. The Offeror should describe any provided specialized networks for resources requiring higher bandwidth and lower latency than the HPC network would provide. [TR-3]

5.5 Data Collection and Monitoring

- 5.5.1 The Offeror should provide a secure mechanism whereby all system monitoring data and logs captured are available to the University. The Offeror will support an open monitoring API to facilitate lossless, scalable sampling and data collection to publish and subscribe to monitoring data. The mechanism should include a sampling and connection framework that allows the system manager to configure independent alternative parallel data streams to be directed off the system to site-configurable consumers. Any filtering that may need to occur will be at the option of the University. [TR-1]
- 5.5.2 The Offeror should provide mechanisms to collect, provide, store, and generate alerts to monitor the status, health, and performance of the system. These mechanisms and data should adhere to available open standards (when available), be open-source, or provide documented APIs and data definitions if only a proprietary solution is available. The Offeror should describe these capabilities, which should include at least the following: [TR-1]
- Environmental measurement capabilities for all systems and peripherals and their sub-systems and supporting infrastructure, including power, energy consumption, voltage, cooling, and temperature, including sampling frequency, accuracy of the data, and timestamps of the data for individual points of measurement
 - Metrics related to memory, network, and other error correction or faults

- Metrics of both HPC protocols and TCP/IP flows. This shall include switch and/or router data, load balancing, error counters, congestion state, throttling, throughput, and latency for select packets traversing the network(s).
- Resource utilization for memory, CPU, network, storage, and accelerator devices
- The system as a whole, including all levels of integrated and attached storage and their associated hardware performance counters, degraded components, and impending failure

5.5.3 The Offeror should describe the provided tools for the collection, analysis, integration, and visualization of metrics and logs produced by the system (e.g., peripherals, integrated and attached storage, and environmental data, including power and energy consumption). [TR-1]

6.0 SYSTEM NETWORKS

6.0.1 The Offeror should describe the high-speed network, including the following: [TR-1]

- High-level description of how traffic would be routed and protocol translations or bridges
- Scale of transfers and the number of connections may be established per network interface and in aggregate
- Aggregate bandwidth and transfer rates between compute nodes and the other networks on (the management and any storage networks) and off (the external data center network) the system
- The network topology, the number of network tiers, and the degree of bandwidth tapering
- The expandability of the network and options to increase the number of endpoints and switches as needed (e.g., how many endpoints can be added to the current network configuration, and how would this affect switch-to-switch connectivity?)
- The Offeror-provided lower-level communication (LLCA) API in the form of either UCX (<https://openucx.org/>) or libfabric (<https://ofiwg.github.io/libfabric/>). The Offeror should describe any enhancements or limitations that can be expected in meeting full support of the standards and latest version of the LLCA. [TR-1]

6.0.2 The Offeror should describe link failure resilience throughout the network (e.g., the number of links, network interfaces, and switch failures that can occur while maintaining connectivity and how performance degrades as links fail). [TR-1]

6.0.3 The Offeror should describe the out-of-band management network and mechanisms to securely extend management segments through an intermediary network, such as a data center network. [TR-1]

6.0.4 The Offeror should provide a management platform that enables dynamic addition/removal of network segments (e.g., subnets, host routes, IPv4, IPv6, and MAC addresses) with stateful tracking of failure and recovery, including removal from and addition to the network. [TR-3]

- 6.0.5 The Offeror should describe any provided mechanisms that ensure quality of service (QoS) for the interconnect (such as congestion control and traffic classes/virtual channels). The Offeror should describe how these could be configured by the University and the guarantees they provide. The offeror should describe mechanisms that enhance the NERSC-10 system's performance on the complex workflows described in the NERSC-10 Workflows Whitepaper. [TR-2]
- 6.0.6 The Offeror should provide a high-bandwidth and resilient solution that allows external connectivity to and from the system to the NERSC data center Ethernet network with support for thousands of simultaneous transfers and capability of at least 26 terabits per second aggregate throughput. NERSC will provide file systems configured for different purposes to users, including but not limited to a HOME file system. The Offeror should work with the University to ensure the NERSC-10 system mounts these file systems to achieve a high level of user satisfaction. [TR-1]
- 6.0.7 The Offeror should describe proposed support for the following: [TR-1]
- Jumbo frames, IPv6, IPv4, TCP/IP, UDP, and virtual networks support
 - Ability to control IP traffic using Access Control Lists
 - Ability to load balance and route traffic across multiple paths
 - Ability to dynamically configure routing and exchange routing information between the data center, storage, and other external networks
- 6.0.8 The Offeror should describe how the external connectivity solution will utilize Internet Engineering Task Force (IETF) standards-compliant technology for functionality that includes but is not limited to the following: [TR-2]
- Congestion control mechanisms across network boundaries
 - QoS required for reliable connections with configurable buffers
 - Traffic draining capability at a link or adjacency level, with flow-tracking ability using sampled flow or similar
 - Encryption and authentication
 - Integrated software-defined network (SDN) with job management and scheduling systems
- 6.0.9 The Offeror should describe support for CXL 2.0 (or greater) standard in its CPU-, GPU-, NIC, and storage offering. The Offeror should include the level of support and the intended usage models of CXL.io, CXL.mem, and CXL.cache. [TR-3]

7.0 STORAGE SYSTEMS

The NERSC-10 system requires two distinct and independent storage systems to accommodate diverse I/O use cases.

The platform storage system (PSS) should be designed to accommodate the I/O and storage needs of multiple workload types, be designed for scale, and be maximized for bandwidth.

The QoS storage system (QSS) should be designed to provide deterministic performance through the use of a configurable Quality of Service (QoS) mechanism. The QSS will be

designed for scale and have the ability to deliver guaranteed I/O to a subset of workloads. Of particular interest are mechanisms that seamlessly enable the time-sensitive workflows described in the NERSC-10 Workflows Whitepaper.

The PSS and QSS will both need to support jobs of varying sizes (from a single node to the whole system) and support workloads that require MPI I/O. Both PSS and QSS should have a data protection scheme in place (e.g., a form of RAID or erasure coding). PSS and QSS will also be accessed via other systems (e.g., external data transfer nodes) and should not require the compute system for them to be accessible, and vice versa.

7.1 Platform Storage System (PSS)

- 7.1.1 The Offeror should propose a separately priced Platform Storage System solution, providing at least 120 PB of usable space with an aggregate sequential read bandwidth and an aggregate sequential write bandwidth of at least 20 terabytes per second each, as measured by the [IOR microbenchmark](#) described in Section 3.0. [TR-1]
- 7.1.2 The Offeror should describe the scalable unit for capacity, bandwidth, and IOPS, and provide details for increasing each of them as well as detail any associated limitations, including the number of concurrent clients. The Offeror should describe the projected characteristics of primary storage devices, such as media type, usable capacity, storage interfaces (e.g., NVMe, PCIe), and media durability. [TR-1]
- 7.1.3 The Offeror should describe all available interfaces to platform storage for the system, including but not limited to POSIX, Kubernetes CSI, object storage interfaces, and other APIs. The Offeror should describe any exceptions to POSIX compliance, time to consistency, and any potential delays for reliable data consumption. [TR-1]
- 7.1.4 The Offeror should propose a method(s) for PSS to be accessible even if the compute system is unavailable. The Offeror should describe any scenarios where a rebalance of resources is required after a reconfiguration and any scenarios where downtime for the entire PSS is required. [TR-1]
- 7.1.5 The Offeror should describe provided data protection schemes, mechanisms for recovery after service interruption or device failure, and related performance impact. The Offeror should describe any anticipated loss of performance over time as the file system ages or reaches capacity. [TR-1]
- 7.1.6 The Offeror should provide and describe features to enforce and report upon soft (accounting) and hard (enforcement) quotas based on uid, gid, project, or other constructs. [TR-1]
- 7.1.7 The Offeror should provide system features for metadata scanning, metadata processing and file/object purging across the entire PSS to allow for purging of older data, and provide data to feed user data management tools. The Offeror should describe the expected rate and elapsed time of a system scan of 10 billion files and any performance impact while a metadata scan or purge is happening. [TR-1]

- 7.1.8 The Offeror should provide the ability of the PSS to purge files, automatically removing them when not accessed within a certain time period. The Offeror should describe any capabilities based upon individual files' last access time, modification time, owner, group, location, and size, including the capability to explicitly include or exclude files and directories, and a dry-run reporting mode. [TR-2]

7.2 QoS Storage System (QSS)

- 7.2.1 The Offeror should propose a separately priced QoS Storage System solution, providing at least 80 PB of usable space, which shall provide storage and deterministic I/O to the system, to provide consistent performance for time-sensitive workloads.

Two modes of operation are anticipated, which could be simultaneous:

- Node-local disk emulation by assigning sufficient capacity and performance to a mounted directory (container, virtual disk)
- Parallel I/O and/or shared work directories by assigning capacity and performance to an allocation shared across nodes

Successful implementation may require coordinating the QSS with network QoS to guarantee end-to-end performance and minimize the effects of contention from other jobs. The Offeror should describe how the design goals will be achieved. [TR-1]

- 7.2.2 The Offeror should describe the scalable unit for capacity, bandwidth, and IOPS and provide details for increasing each and any associated limitations, including the number of concurrent clients. The Offeror should describe projected characteristics of primary storage devices such as media type, usable capacity, storage interfaces (e.g., NVMe, PCIe), and media durability. [TR-1]
- 7.2.3 The Offeror should describe any provided capability for the QSS to provide user and administrator control mechanisms for creating, serving, and querying the status of storage allocations, where an allocation is a directory (e.g., container, virtual disk) that has a capacity, capabilities (e.g., minimums and maximums for bandwidth and IOPS), and policies (e.g., persistence or duration, user quota, namespace mapping). The Offeror should describe any mechanisms that ensure that access to an allocation is limited to authorized users and to display current unallocated capacity and capability. [TR-1]
- 7.2.4 The Offeror should describe how allocation mechanisms could be automated (e.g., how the QSS could receive and serve requests for storage allocations via Slurm) or otherwise improved in the future, as well as detail the time required to prepare a new allocation, from request to ready-to-use. [TR-3]
- 7.2.5 The Offeror should describe all available interfaces (and any associated limitations) to the QSS, including but not limited to POSIX, Kubernetes CSI, NVMeoF, object storage interfaces, and other APIs. The Offeror should describe any exceptions to POSIX compliance, time to consistency, and any potential delays for reliable data consumption. [TR-1]

- 7.2.6 The Offeror should describe any provided data protection schemes, mechanisms for recovery after a service interruption or device failure, and related performance impact. The Offeror should describe any anticipated loss of performance over time as the file system ages or reaches capacity. [TR-1]
- 7.2.7 The Offeror should propose a method(s) for QSS to be accessible even if the compute system is unavailable. The Offeror should describe any scenarios where a rebalance of resources is required after a reconfiguration and any scenarios where downtime for the entire QSS is required. [TR-1]
- 7.2.8 The Offeror should provide system features for metadata scanning and processing across the entire QSS to produce data to feed tools for user data analysis and management. The Offeror should describe the expected rate and elapsed time of a full-system scan and any performance impact while the metadata scan is happening. [TR-2]

8.0 SYSTEM OPERATION

8.1 Resilience, Reliability and Availability Metrics

The ability to achieve the NERSC-10 mission goals hinges on the productivity of system users. System availability is, therefore, essential and requires system-wide focus to achieve a resilient, reliable, and available system. For each metric specified below, the Offeror should describe how it determined its estimate(s).

- 8.1.1 The system is available if it meets the system outage criteria in the glossary. The Offeror should propose a system availability. [TR-1]
- 8.1.2 The Offeror should propose a System Mean Time Between Interrupt (SMTBI). [TR-1]
- 8.1.3 The Offeror should propose a Job Mean Time To Interrupt (JMTTI) for a single job running on the entire system. [TR-1]
- 8.1.4 The system should complete power on and power off in a reasonable time. The Offeror should describe the sequence of steps and timings for full-system initialization and full-system shutdown. Include any dependencies and how timings may scale with the size of the system. [TR-1]

8.2 System Security

- 8.2.1 The Offeror should describe the security capabilities of the proposed compute node and service partition operating systems. [TR-1]
- 8.2.2 The Offeror should describe how the system may be configured to support Zero-Trust requirements as described in the CISA Zero Trust Maturity Model (<https://www.cisa.gov/zero-trust-maturity-model>) [TR-1]

- 8.2.3 The Offeror should describe how their implementation of Identity Management (including Federated Identity) for the system works, what protocols and standards are being utilized (e.g., SpiFFE, SPIRE), and what system services and/or service accounts will use that framework. [TR-1]
- 8.2.4 The Offeror should describe how they will implement a Root of Trust to assure a boot environment for the system that is secure. [TR-1]
- 8.2.5 Security vulnerabilities in the software supplied by the vendor should be addressed with patches and/or updates in a reasonable time, depending on the classification and severity of the vulnerability. The Offeror should describe the process for handling security vulnerabilities and the time to provide patches from confirmed availability for the following: [TR-1]
- Non-critical vulnerabilities
 - Vulnerabilities as defined by CISA
 - Common Vulnerabilities and Exposures (CVEs) in the National Vulnerability Database (NVD) as defined by NIST with a score of Critical or High as defined by the latest version of CVSS.
- 8.2.6 The Offeror should describe how the system software components are validated and data components are managed to be in compliance with the Presidential Executive Order on Improving the Nation’s Cybersecurity ([Cybersecurity Executive Order 14028](#)). For example, providing a Software Bill of Materials (SBOM), describing the management and control process for system data, etc. [TR-1]
- 8.2.7 The Offeror should describe how the baseline configuration is tested, validated, and documented. The Offeror should detail the process of revalidation, at any time, of a running system against the baseline, including how the system can be audited so its current state can be captured and documented. [TR-2]
- 8.2.8 The Offeror should describe the security model, tools, and functionality for observability for any system applications and services that are containerized. [TR-1]
- 8.2.9 The Offeror should describe the available security controls and how they would be deployed and used for access methods to the system and/or services (for example, interactive access using SSH, non-interactive methods using APIs, or Web interfaces). [TR-1]

8.3 Power and Energy

- 8.3.1 The maximum power consumed by the system and its peripheral systems, including the proposed storage systems, should not exceed 20 MW. The maximum power consumption includes all equipment provided by the proposal. The Offeror should describe how its proposal fits within the power budget. [TR-1]

- 8.3.2 The Offeror should describe the power management capabilities of the system available to users or administrators. The description should include a description of all system control capabilities to affect power or energy use (system, rack/cabinet, board, node, component, and sub-component level) and the reaction time of this mechanism. [TR-1]
- 8.3.3 The Offeror should describe any system data source or alarm that provides advanced warning of system power draw changes >1MW over a 15-minute time period, or a similar criteria. [TR-1]
- 8.3.4 The Offeror shall describe the overall system measurement capabilities that enable meeting Green 500 requirements at Level 2 or higher for power, current, and voltage. [TR-3]

8.4 Maintenance and Support (Hardware/Software)

- 8.4.1 The Offeror should propose and separately price maintenance and support for all systems for a period of five (5) years from the date of acceptance of the system. The maintenance and support will include all features outlined in the **Key Elements of the Maintenance and Support Plan** described in Appendix B.9 or detail how the proposed plan differs. [TR-1]
- 8.4.2 The system should include a means for tracking and analyzing all software updates, software and hardware failures, and hardware replacements over the lifetime of the system. [TR-2]
- 8.4.3 The Offeror should provide one (1) systems operations and advanced administration training for each system delivered at facilities specified by the University. The Offeror should describe their training for systems operations and advanced administration available for the lifetime of the system, including topics, duration, and proposed timing. [TR-1]

8.5 Documentation

- 8.5.1 The Offeror should describe the documentation provided to effectively administer and use the system, including types of documentation, format (e.g., user manuals, man pages, release notes, stable URLs, plain text versus PDF, vendor websites, any interactive elements, etc.), initial delivery, and frequency of updates, for the following: [TR-1]
- Documentation for each delivered system describing the configuration, interconnect topology, labeling schema, hardware layout, etc., of the system as deployed before the commencement of system acceptance testing
 - Documentation of the proposed solution to the operators and system administrators to effectively operate and configure the platform

- Documentation for users describing the programming environment and software tools, including compatibility across system software updates
- 8.5.2 The Offeror should demonstrate in its proposal how it will grant the University use and distribution rights for provided documentation, training session materials and recorded media to be shared with DOE Lab staff and all authorized users and NERSC support staff. The University may, at its option, make audio and video recordings of presentations from the Offeror at public events targeted at the NERSC user communities (e.g., user training events, collaborative application events, hackathons, best practices discussions) and make them available to all NERSC users. [TR-2]
- 8.5.3 The Offeror should demonstrate in its proposal how it will ensure all documentation is distributed and updated electronically in a reasonable time that maintains the productivity and performance of the system. For example, information about binary compatibility following system changes (e.g., major OS upgrades) should be provided in relevant documentation. [TR-1]
- 8.5.4 The Offeror should demonstrate in its proposal how it will ensure all documentation is maintained and up-to-date. Documentation of changes and fixes may be distributed electronically in the form of release notes. Reference manuals may be updated in a reasonable time. [TR-1]

9.0 FACILITIES AND SITE INTEGRATION

The following section addresses the facility and site-based requirements for the proposed system. It includes pertinent information and vendor requirements covering the physical, electrical, cooling, seismic, safety, and transportation aspects of designing, delivering, installing, and integrating the system.

- 9.0.1 The Offeror's proposed system and integration plan should include, but not be limited to, all key elements outlined in the **Facilities and Site Preparation Plan** in Appendix B.4 and should conform to the **Facility and Site Integration Specifications** provided in Appendix C. The Offeror should describe any limitations to meeting the specifications. [TR-1]
- 9.0.2 The Offeror should describe capabilities to improve the environmental sustainability and system efficiency, beyond energy efficiency, through the entire lifecycle of the system, including design, manufacturing, deployment, and operation, and the ability to reuse and recycle components into future systems, and in final disposal. [TR-1]

10.0 NON-RECURRING ENGINEERING (NRE)

The University expects to award a Non-Recurring Engineering (NRE) subcontract, separate from the system build subcontract. Cost sharing is a condition of the class advance waiver for a large business awardee. It is anticipated that the NRE subcontract could be approximately 10-20% of the NERSC-10 system price. The Offeror should provide proposals for areas of

NRE that will provide substantial value to the NERSC-10 system (e.g., TR-2 or TR-3 requirements). The goals of the NRE efforts may include but are not limited to the following:

- Optimizing the integration and usage of new hardware and software
- Enhancing the NERSC-10 system for workflows
 - Increasing performance and portability
 - Enabling composability, automation and seamlessness, debugging and profiling, and monitoring
- Enhancing the overall NERSC-10 system
 - Increasing resilience, reliability, and security
 - Enabling seamless integration into data center
 - Optimizing power and energy usage

Proposed NRE areas should focus on topics that may:

- Provide added value that enhance or go beyond planned roadmap activities
- Mitigate risk to ensure successful deployment of the system and/or future technologies to support the system schedule

Proposals should not focus on one-off point solutions or gaps created by the Offeror's proposed design for the HPC system that should be otherwise provided as part of a vertically integrated solution.

11.0 TECHNICAL OPTIONS

Technical Options will be evaluated as part of the proposal evaluation process; however, the University may or may not elect to include Technical Options in the resulting subcontract. Each proposed Technical Option should appear as a separately priced and identifiable item in the Offeror's proposal response.

11.1 Upgrades, Expansions, and Additions

The University expects to have future requirements for system upgrades and/or additional quantities of components based on the configurations proposed in response to this solicitation. The Offeror should propose separately priced options using whatever is the natural unit for the proposed architecture design as determined by the Offeror. For example, for system size, the unit may be the number of racks or some other unit appropriate for incrementally increasing the system. The Offeror should identify any thresholds requiring increased component infrastructure (e.g., extra spine switches), any technical challenges foreseen with respect to scaling, and any other production issues.

11.1.1 The Offeror should describe upgrades, expansions, or procurement of additional system configurations by the following fractions of the system, as measured by the Workflow-SSI metric. All additional system configurations should contain all the functionality of the production NERSC-10 system, including storage systems but scaled to the appropriate configuration. [TO-1]

- 25%
- 50%

- 100%
 - 200%
- 11.1.2 The Offeror should propose CPU-only compute nodes. Include a description of the natural unit (e.g., number of racks, number of nodes, peak PF, etc) for the proposed design as determined by the Offeror. [TO-1]
- 11.1.3 The Offeror should propose upgrades, expansions, or procurement of additional platform storage system (PSS) capacity in increments of 10% for the scalable units described in Section 7. [TO-1]
- 11.1.4 The Offeror should propose upgrades, expansions, or procurement of additional QoS storage system capacity (QSS) in increments of 10% for the scalable units described in Section 7. [TO-1]
- 11.1.5 The Offeror should propose non-volatile storage options (NVMe) and/or solid-state drives (SSD) internal to all node types, with a target size of at least triple the node memory capacity. Non-volatile storage performance characteristics and system software requirements to access and manage the storage should be described. [TO-3]
- 11.1.6 The Offeror should propose a novel AI acceleration partition that will both accelerate AI workloads and integrate within an HPC ecosystem to enable workflow capabilities. The partition may be delivered with the main NERSC-10 system or later and be made up of the scalable unit for the proposed design as determined by the Offeror. The partition should be made up of enough scalable units to achieve 1.0x Workflow-SSI for the DeepCAM workflow component benchmark described in Section 3.0. The Offeror’s description should include the following: [TO-3]
- An overall architectural diagram that shows all hardware, interconnect(s), compilation infrastructure, and Input/Output (I/O) subsystems, if applicable. The Offeror should work with the University to ensure the partition mounts the PSS and QSS file systems to achieve the necessary level of performance.
 - An overview of software architecture, including libraries and Software Development Kits (SDKs), support for frameworks (e.g., TensorFlow, PyTorch, etc.), usability, and programmability. The description should include terms of software licensing, the number of licenses included, and any support if applicable.
 - Available results or projections on [MLPerf benchmarks](#), in particular the “Training”, “Training: HPC”, and “Inference Datacenter” benchmark suites. Results and projections should specify the version of the benchmark used, whether the result was officially submitted to the MLCommons foundation, and any required modifications to the benchmark rules required to obtain the reported results or projections.
 - Performance results (actual, predicted, or extrapolated) for the proposed system for one or more of the workflow component benchmarks listed in Table 3.1.
- 11.1.7 The Offeror should propose additional maintenance and support extension for years 6-7. [TO-1]

11.1.8 The Offeror should propose to deinstall, remove, and/or recycle the system and supporting infrastructure at end of life. Storage media shall be wiped or destroyed to the satisfaction of the University and/or returned to the University upon request. [TO-1]

11.2 Early Access Systems

To allow for early and/or accelerated development of applications or development of functionality required as a part of the statement of work, the Offeror should propose options for Early Access Systems (EAS). The EAS should contain similar functionality to the final system, including storage systems, management, and workflow environment nodes, but scaled down to the appropriate configuration.

- 11.2.1 The Offeror should propose a **“Phase 1” EAS** that could be delivered in 2025 for production use. The primary purpose is to expose applications to the same programming environment as will be found on the final system. It is acceptable for the early access system not to use the final processor, node, or high-speed interconnect architectures. However, the programming and runtime environment must be sufficiently similar that a port to the final system is straightforward. The Offeror should propose and describe an option that will deliver at least a 1x performance improvement of the Workflow-Sustained System Improvement (Workflow-SSI) metric over *Perlmutter* as defined in Section 3.0. [TO-1]
- 11.2.2 The Offeror should propose and describe a small **“Pilot” EAS** that could be delivered in 2025 to aid in the integration and development of the NERSC-10 system, form the basis for NRE engineering efforts, or reduce risk and/or accelerate development. [TO-1]
- 11.2.3 The Offeror should propose other **EAS** (hardware and software) to aid in the integration and testing of the NERSC-10 system, form the basis for NRE engineering efforts, or reduce risk and/or accelerate development. Of interest are resources that are in support of any topics proposed for NRE and workflow readiness. The systems could be on-site at NERSC and/or accessed remotely and could include access to early software, simulators, and/or emulators. [TO-2]

11.3 Test Systems (TS)

A test system (TS) shall contain all the functionality of production NERSC-10 systems, including storage systems and all accelerator types, but scaled down to the appropriate configuration. Multiple TSs may be awarded to aid with the lifetime system management of any production system. It is desirable for NERSC-10 production systems and TSs to be able to dynamically attach and detach from the same resources to allow scale testing on the test system by temporarily moving these resources from the production system to the TS. The Offeror should propose TSs for any production system delivered in support of the NERSC-10 system. The TSs should be delivered before the production resource they are designed to support.

- 11.3.1 The Offeror should propose a Pre-Production TS, which should contain at least 32 compute nodes. [TO-1]
- 11.3.2 The Offeror should propose a System Development TS, which should contain at least 16 compute nodes. [TO-1]

12.0 DELIVERY AND ACCEPTANCE

The Acceptance Test Plan (ATP) will describe the steps needed to validate the system, including pre-delivery, post-delivery, and acceptance test plan for the NERSC-10 system and supporting systems, such as EASs and TSs.

Acceptance testing may comprise multiple components for which the overall goal is to ensure that the system as a whole is high performance, scalable, resilient, and reliable. Acceptance testing may exercise the system infrastructure with a combination of functionality tests, performance tests, forced failures, and stability tests. Any requirement described in the statement of work as awarded may generate a corresponding acceptance test element. The specifics of an acceptance plan will be determined before system delivery.

- 12.0.1 The Offeror should commit to work with the University to define and perform the pre-delivery, post-delivery, and acceptance testing as part of the NERSC-10 System Acceptance Test Plan. The Offer should describe any limitations to meeting the Sample Acceptance Test Plan in Appendix A, including scale and licensing limitations for any component of the ATP. [TR-1]

13.0 PROJECT AND RISK MANAGEMENT

The development, pre-shipment testing, installation, and acceptance testing of the NERSC-10 System and the management of the Non-Recurring Engineering (NRE) subcontract are complex endeavors and will require close cooperation between the successful Offeror and the University.

- 13.0.1 The Offeror should provide a proposed set of milestones and deliverables. Prior to award, the successful Offeror and the University will finalize the list of milestones. For each milestone, the Offeror should include a proposed associated payment that is applicable to the NERSC-10 system development and deployment timeline. The Offeror is encouraged to identify milestones for every year that the project merits revenue that the Offeror can legally recognize in that year. Milestones payments may be based on any of the following:

- Performance measured by objective quantifiable methods;
- Accomplishment of defined events; or
- Other quantifiable measures of results

Milestones should include dates identifying necessary Technical Decision Point (TDP) evaluations. The successful Offeror and the University will hold a TDP evaluation and joint planning meeting 9-12 months before system delivery. At the TDP meeting, the final configuration of the system will be determined based on technology status,

evaluations, and component pricing. Performance targets will be reevaluated and converted into requirements.

The following provide examples of milestones that would be of importance to the University. Other milestones may be needed for phased installations or deployments featuring major upgrades during the subcontract. [TR-1]

- Completion of Project Planning Deliverables
- Pilot and/or Phase I system delivery/acceptance
- Technical Decision Point
- Begin delivery and installation of system
- System installation and integration complete, including I/O subsystem
- System accepted, including I/O subsystem

13.0.2 The completed documents are not required in the RFP response; however, the Offeror should commit to develop, maintain, and submit to the University the Project Planning Deliverables described in this section within the timeframes required through a reliable and easily accessible mechanism that supports change control. The Offeror should describe any limitations to providing these plans and propose additional plans to understand and mitigate risk successfully. [TR-1]

- Project Plan B.1
- Communication Plan B.2
- Risk Management Plan B.3
- Facilities and Site Preparation Plan B.4
- Chemical Management Plan B.5
- Network Plan B.6
- Delivery and Installation Plan B.7
- Acceptance Test Plan B.8
- Maintenance and Support Plan B.9
- Center of Excellence for Workflow Readiness Plan B.10
- Training and Education Plan B.11

The Project Planning Deliverables listed above are described in more detail in Appendix B: Project and Risk Management - Planning Deliverables Descriptions.

13.0.3 The successful Offeror and the University will assess the project for areas in which additional tracking and collaboration are necessary, and form working groups (WGs). The WGs will serve as a key conduit to collaboratively identify, refine, and understand NERSC requirements in detail and to ensure that the delivered system meets those requirements to the greatest extent possible.

The Offeror should commit to form and actively participate in project working groups. WGs should establish a regular schedule for electronic meetings (e.g., teleconferences) including quarterly face-to-face meetings that may include working group breakout sessions. The WGs will be assessed on an ongoing basis including if some WGs are no

longer required or new WGs are required for additional topics. The Offeror should describe any limitation to working group participation. [TR-1]

- 13.0.4 The Offeror should commit to schedule and complete a Project Planning Kickoff Meeting. The Offeror should describe proposed topics to include in the agenda to review project management goals, techniques, and processes. The meeting should be held no later than 45 days after subcontract award. [TR-1]

DRAFT

End of Target Requirements and Technical Option Requirements

The following Appendices A, B, and C are intended to inform the successful Offeror of the performance activities and facility specifications associated with the NERSC-10 HPC System and NRE subcontracts. This information is also intended to assist Offerors in understanding what Offeror resources will be necessary, in the University's assessment, to successfully perform the project and within the Offeror's proposal provide sufficient detail of pricing information for evaluation by the University.

APPENDIX A: Sample Acceptance Test Plan

The Acceptance Test Plan (ATP) describes the steps needed to validate the NERSC-10 System. Some, or all, of these steps may be needed to validate supporting systems, such as EASs and TSs as agreed-upon by the University and Subcontractor.

A.1 Staffing and Safety

The Subcontractor must provide sufficient staffing to perform the installation, initial testing, and acceptance of the System.

The Subcontractor must conform to the University's safety protocols and policies and complete all necessary documents (e.g., approved safety plan) and required training. The system must be delivered and installed in accordance with the University's safety policies and the approved joint safety plan. The Subcontractor must designate a safety supervisor who will monitor Subcontractor staff for adherence to the University's safety policies and the approved safety plan.

A.2 Pre-Delivery Assembly, Quality Assurance, and Factory Test Plan

The Subcontractor must perform the pre-delivery assembly and quality assurance tests of the System and agreed-upon sub-configurations at the Subcontractor's location, demonstrating all hardware is fully functional prior to shipment. A factory test plan will be agreed on by the University and Subcontractor no less than 30 days prior to the factory test.

At its option, the University may send a representative(s) to observe testing at the Subcontractor's facility. Work to be performed by the Subcontractor includes the following:

- All hardware installation and assembly
- Burn-in of all components
- Installation of software
- Implementation of the University-specific production system configuration and programming environment necessary to complete required testing
- Perform tests and benchmarks to validate functionality, performance, reliability, and quality
- Run benchmarks and demonstrate that benchmarks meet performance commitments.

Pre-delivery testing includes the following:

- Demonstrate RAS capabilities and robustness using simple fault injection techniques, such as disconnecting cables, powering down subsystems, or installing known bad parts.
- Demonstrate functional capabilities on each segment of the system built, including the capacity to build applications, schedule jobs, and run them using a customer-provided testing framework. The root cause of application failure must be identified prior to system shipping.
- Provide a file system sufficiently provisioned to support the suite of tests.
- Provide onsite and remote access to the University to monitor testing and analyze results.
- Instill confidence in the ability to conform to the statement of work.

The pre-delivery test must consist of (but is not limited to) the following tests:

Name of Test	Pass Criteria
System Power Up	All nodes boot successfully
System Power Down	All nodes shut down
Monitoring	Monitoring software shows status for all nodes
Reset	“Reset” functions on all nodes
Power On/Off	Power cycle all components of the entire system from the console
Fail Over/Resilience	Demonstrate operation of all fail-over or resilience mechanisms
Benchmarks	The system must demonstrate the ability to achieve the required performance level on all benchmark requirements on the agreed-upon scale and configuration
72-Hour Test	High availability given by the system outage definition of the system for a 72-hour test period under constant throughput load

A.3 Site Integration and Post-Delivery Testing

The system must be delivered, installed, and fully integrated and must undergo Subcontractor stabilization processes. Post-delivery testing must include replication of all of the pre-delivery testing steps, along with appropriate tests at scale, on the fully integrated system. During post-delivery testing, the pre-delivery tests must be run on the full-system installation. Where applicable, tests must be run at full scale.

When the Subcontractor has declared the system to be stable, the Subcontractor must make the system available to University personnel for site-specific integration and customization. After the Subcontractor’s system has completed site-specific integration, security screening, and customization, the acceptance test must commence.

A.4 Acceptance Testing

The Acceptance Test Period shall commence when the system has been delivered, physically installed, undergone stabilization, site-specific integration, and customization, and conforms to all requirements in the statement of work designated “critical” priority. The Acceptance Test Period should target 60 days or until it has met agreed-upon exit criteria.

The Subcontractor will not be responsible for failures to meet the performance metrics or the availability metrics set forth in this plan, if such failure is the direct result of modifications made by the University to Subcontractor source code. Such suspension will be only for those requirements that fail due to the modification(s) and only for the length of time the modification(s) result(s) in the failure.

The Subcontractor must supply source code used, compile scripts, output, and verification files for all tests. All such provided materials become the property of the University.

All tests must be performed on the production configuration of the system, as it will be deployed to the University user community. The University may run all or any portion of these tests at any time on the System to ensure the Subcontractor's compliance with the requirements set forth in the Acceptance Test Plan.

The acceptance test must consist of a Functionality Demonstration, System Test, System Resilience Test, Performance Test, and an Availability Test.

A.4.1 Functionality Demonstration The Subcontractor and the University will perform the Functionality Demonstration on a dedicated system. The Functionality Demonstration must show that the system is configured and functions in accordance with the statement of work. Demonstrations must include, but are not limited to, the following:

- Remote monitoring, power control, and boot capability
- Network connectivity
- File system functionality
- Batch system
- System management software
- Program building and debugging (e.g., compilers, linkers, libraries, etc.)

A.4.2 System Test The Subcontractor and the University will perform the System Test on a dedicated system. The System Test must demonstrate that the system is configured and functions in accordance with the statement of work. Demonstrations must include but are not limited to the following:

- Two successful system cold boots to production state in accordance with required timings, with no intervention to bring the system up. Production state is defined as running all system services required for production use and being able to compile and run parallel jobs on the full system. In a cold boot, all elements of the system (compute, login, I/O, network) are completely powered off before the boot sequence is initiated. All components are then powered on.
- Single node power-fail/reset test: Failure or reset of a single compute node must not cause a system-wide failure. A node must reboot to production state after reset in accordance with required timings.

A.4.3 System Resilience Test The Subcontractor and the University will perform the System Resilience Test on a dedicated system. The System Resilience Test must show that the system is configured and functions in accordance with the statement of work.

All demonstrable system resilience features of the system must be demonstrated via fault-injection tests when running test applications at scale. Fault injection operations must include both graceful and hard shutdowns of components. The metrics for resilience operations include correct operation, any loss of access or data, and time to complete the initial recovery, plus any time required to restore (fail-back) a normal operating mode for the failed components.

A.4.4 Performance Test The Subcontractor and the University will run the NERSC-10 tests, workflow component benchmarks, and file system tests a minimum of five times each, as described in the Benchmark Run Rules. Benchmark answers must be correct, and each benchmark result must meet or exceed performance commitments.

Benchmarks must be run using the supplied resource management and scheduling software. Except as required by the run rules, benchmarks need not be run concurrently. If requested by the University, the Subcontractor must reconfigure the resource management software to utilize only a subset of compute nodes specified by the University. Performance must be consistent from run-to-run.

A.4.5 Availability Test The Availability Test will commence after successful completion of the Functionality Demonstration, System Test, and Performance Test. The Subcontractor must perform the Availability Test; at this time or before, the University will add user accounts to the system. The Availability Test shall be 30 contiguous days in a sliding window within the Acceptance Test Period. The NERSC-10 System must demonstrate the required availability of the system.

During the Availability Test, the University must have full access to the system and will monitor the System as if it was in production. The University and users designated by the University will submit jobs through the NERSC-10 resource management system. These jobs will be a mixture of benchmarks from the Performance Test and other applications.

The Subcontractor must adhere to the System Availability and Reliability requirements as defined in the following:

- All hardware and software must be made fully functional before the availability test can be declared complete. Any downtime required to repair failed hardware or software must be considered an outage unless it can be repaired without impacting system availability.
- Hardware and software upgrades must not be permitted during the last seven days of the Availability Test. The system must be considered down for the time required to perform any upgrades, including rolling patch upgrades.
- No critical bugs must be open during the last seven (7) days prior to the Availability Test.
- During the Availability Testing period, if any system software upgrade or significant hardware repairs are applied, the Subcontractor shall be required to run the Benchmark Tests and demonstrate that the changes cause no loss of performance. At its option, the University may also run any test or benchmark deemed necessary. Time taken to run the Benchmark and other tests shall not count as downtime, provided that all tests perform to specifications.
- Every test in the Functionality Test, Performance Test, and NERSC-defined workload shall obtain a correct result in both dedicated and non-dedicated modes.
- Each benchmark in the Performance Test shall meet or exceed the performance commitment and any variation requirement.

DRAFT

APPENDIX B: Project and Risk Management - Planning Deliverables Descriptions

The selected Subcontractor must develop, deliver, submit for approval, and maintain the following Planning Deliverables. Initial versions and updates of these plans must be provided in specific agreed-upon time frames. Each of the plans and any revisions will be submitted for comment and approval to the University's technical representative.

The specific details are designed to help the Subcontractor successfully meet its commitment, to help the University track the NERSC-10 project, and to help the University and the Subcontractor to understand and to mitigate risks successfully.

B.1 Key Elements of the Project Plan

The Subcontractor and University will have a joint project planning call no more than 30 days after subcontract award. The Subcontractor must provide the University with a detailed Project Plan no later than 90 days after subcontract award that addresses, at a minimum, the following:

- The Subcontractor must appoint a Program Manager (PM) for the purposes of executing the Project Plan for the system and NRE contracts on behalf of the Subcontractor. The PM must serve as the primary interface for the University, managing all aspects of the Subcontract in response to the subcontract requirements.
- Project Management Organization Breakdown Structure (OBS), with management team's roles and responsibilities, must be clearly defined.
- Points of contact to include the Subcontractor's Technical Contact(s) - for example, Service Manager, Contract Manager, and Account Manager - must be provided.
- Work Breakdown Structure (WBS) to include all major subsystems, each software product, and each major equipment deliverable to the University must be provided.
- The subcontractor must provide a full Project Schedule Gantt chart for the duration of the contract.

The Project Plan must be updated quarterly.

B.2 Key Elements of the Communication Plan

The project planning kick-off meeting must take place no later than 45 days after subcontract award. A Communication Plan must be developed and must describe the types of communications, meetings, and progress reviews as described below:

- Daily Communication (System Contract)
The Subcontractor's PM (or designate) is the owner of this meeting with a target duration of one-half hour. Both Subcontractor and the University may submit agenda items for this meeting. These daily communications must commence 30 days before expected system delivery and continue until both parties agree they are no longer needed
- Weekly Status Meeting (System and NRE subcontracts)

The Subcontractor's PM must schedule this meeting with a target duration of one hour. Attendees normally include the Subcontractor's PM, Service Manager, University's Procurement Representative, Technical Representative, and System Administrator(s) as well as other invitees.

- Quarterly Business Reviews (System and NRE subcontracts)

The Subcontractor's PM must schedule this meeting with a target duration of no less than six hours. Attendees normally include: Subcontractor's PM, Subcontractor's Senior Management, University's Procurement Representative, Technical Representative, selected NERSC Management, selected NERSC Technical Staff, and other invitees. Topics covered will cover both NRE and System subcontract issues that will include the following:

- Program status (Subcontractor to present)
- University satisfaction (University to present)
- Partnership issues and opportunities (joint discussion)
- Future hardware and software product plans and potential impacts for the University
- Participation by Subcontractor's suppliers as appropriate
- Other topics as appropriate
- Both Subcontractor and the University may submit agenda items for this meeting.

B.3 Key Elements of Risk Management (System and NRE contract)

The Subcontractor must provide the University with a Risk Management Plan (RMP) for the technology, schedule, and business risks of the NERSC-10 project no later than 30 days after award of the Subcontract. The RMP describes the Subcontractor's approach to managing NERSC-10 project risks by identifying, analyzing, mitigating, contingency planning, tracking, and ultimately retiring project risks. The Plan must address both the System and the NRE portions of the project. The purpose of this RMP, as detailed below, is to document, assess, and manage Subcontract's risks affecting the NERSC-10 project:

- Document procedures and methodology for identifying and analyzing known risks to the NERSC-10 project along with tactics and strategies to mitigate those risks.
- Serve as a basis for identifying alternatives to achieving cost, schedule, and performance goals.
- Assist in making informed decisions by providing risk-related information.

The RMP must include, but is not limited to, the following components: management, hardware, software; risk assessment, mitigation and contingency plan(s) (fallback strategies). Once the plan is approved by the University, the University will review the Subcontractor's RMP annually.

The Subcontractor must also maintain a formal Risk Register (RR) documenting all individual risk elements that may affect the successful completion of the NERSC-10 project (both System and NRE subcontracts). The RR is a database managed using an application and format approved by the University. The initial RR is due no later than 30 days after award of the Subcontract. The RR shall be updated at least monthly and before any Critical Decision

(CD) reviews with DOE. After acceptance, the RR must be updated quarterly. Items in the RR include, at a minimum, mitigation strategies, impact to the NERSC-10 project, severity rating, and probability of the risk occurring.

B.4 Key Elements of Facilities and Site Preparation Plan

Facilities and site preparation planning must be conducted by the Subcontractor's site engineering staff. Facilities and Site Preparation planning consists of the exchange of system specification documents and site floor plans and is followed by a physical inspection of facilities. The plan must include the drawings of the computer floor and the Machine Unit Specification (MUS). The plan will include the NERSC-10 system and supporting systems, such as EASs (e.g. Pilot and/or Phase 1) and TSs.

The Subcontractor must provide a Preliminary Facilities and Site Preparation Plan at least one year prior to the delivery of the first equipment or agreed-upon milestone date. The Final Facilities and Site Preparation Plan shall be delivered to the University for approval at least nine (9) months prior to the first equipment delivery or agreed-upon milestone date. Details of the Facility requirements are in Appendix C. Items in the plan, at a minimum, include the following:

- System layout in floor plan view
- Cabinet dimensions, diagrams, and cabinet weights
- Electrical requirements
- Cooling water requirements
- Environmental (room air condition) requirements
- Power cable and water hose requirements
- Raised floor requirements, tile cutout dimensions and locations
- Seismic restraint connection information
- Overhead cable tray information and requirements
- Electrical and mechanical safety features and requirements.

B.5 Key Elements of the Chemical Management Plan

The Subcontractor must develop and document a chemical management plan describing all chemical additives used for maintaining the NERSC-10 system and cooling water property requirements in coordination with the University. The plan must include all aspects of transport, storage, filling, draining, and disposal of the chemicals, safety data sheet (SDS) compliance, and required training and personal protective equipment needed for compliance with safety regulations. The plan must be documented and submitted for University review and concurrence nine (9) months prior to the delivery and deployment of the system.

B.6 Key Elements of the Network Plan

No less than six (6) months prior to installation of the system, the Subcontractor must provide an initial draft of the system network configurations for the University to review and approve, including all network types provided, and showing compute-to-compute, compute-to-storage, and system-to-external components connectivity.

No less than four (4) months prior to installation, the University and the Subcontractor must finalize the network design and the Subcontractor must provide an up-to-date copy of the network configurations reflecting the expected-at-installation design.

B.7 Key Elements of Delivery and Installation Plan

The Subcontractor must provide a Preliminary Delivery and Installation Plan to the University one year prior to the first equipment delivery and must update it as necessary. The plan will include the NERSC-10 system and supporting systems, such as EASs (e.g., Pilot and/or Phase 1) and TSs. The Subcontractor must provide a Final Delivery and Installation Plan no less than 90 days before each delivery. The plan must include the following:

- Installation team staffing plan
- Detailed delivery and installation schedule
- Installation sequence plan for each delivery
- Detailed integration and test plan addressing all equipment and software included in the delivery
- Safety documents covering the delivery and installation activities as required by the University
- Diagrams showing the internal layout of cabinets.

The University will review the plan and work with the Subcontractor to promptly resolve any issues or needed clarifications.

B.8 Key Elements of the Acceptance Test Plan

The University and the Subcontractor will create a detailed Acceptance Test Plan one year prior to the first equipment delivery that will be updated as necessary. The Acceptance Test Plan will describe the steps needed to validate the system, including pre-delivery, post-delivery, and acceptance test plan for the NERSC-10 system and supporting systems, such as EASs and TSs.

Items in the plan must, at a minimum, include the following:

- Pre-Delivery Assembly, Quality Assurance, and Preliminary Factory Testing. The plans will include a description of how the Subcontractor qualifies its vendors, factory burn in and validation test plans and a pre-ship test plan for the NERSC-10 system.
- Acceptance Testing. The plans must consist of Functionality Demonstrations, System Tests, System Resilience Tests, Performance Tests, and an Availability Test.
- EAS (e.g. Pilot and/or Phase 1) and TS Testing. Details for testing systems in support of the NERSC-10 production system.

A sample Acceptance Test Plan is included in Appendix A.

B.9 Key Elements of Maintenance and Support Plan

The Subcontractor must provide a maintenance, on-site support, and services plan for the life of the subcontract and must include the following features:

- **Maintenance and Support Period**
The Subcontractor must provide all maintenance and support for a period of five (5)

years from the date of acceptance of the system. Warranty as defined in the LBNL General Provisions, must be included in the five (5) years. The Subcontractor must also provide any component warranty period (e.g. processor, disk, etc.) within the maintenance period. For example, if the system is accepted on April 1, 2026 and the warranty is for one year, then the Warranty ends on March 31, 2027. The maintenance period begins April 1, 2026 and ends on March 31, 2031.

- **Maintenance and Support Solution**

The Subcontractor must propose a maintenance and support solution with full hardware and software support for all Subcontractor-provided hardware components and software. The principal period of maintenance (PPM) shall be for 24 hours by seven (7) days a week with a four-hour response to any request for service. The Subcontractor must provide/enable access to direct communication between NERSC staff and GPU vendor technical staff.

- **Concurrent Maintenance Techniques**

The Subcontractor must use continuous operations maintenance techniques (e.g., warm swap) that avoid service disruptions. Continuous operations comprise both hardware (including servicing node hardware and cabinet hardware) and software upgrades to systems management nodes, workflow nodes, storage, and compute nodes. These actions will not be deemed to cause a system outage if performed with the concurrence of the University and completed in a reasonable time. Six (6) hours are permitted for cabinet-level repairs and two (2) hours for all other repairs performed concurrently; node downtime due to concurrent maintenance is counted in calculating system availability.

- **General Service Provisions**

The Subcontractor must be responsible for repair or replacement of any failing hardware component that it supplies and correction of defects in software that it provides as part of the system. At its sole discretion, the University may require advance replacement of components that show a pattern of failures that reasonably indicates that future failures may occur in excess of reliability targets, or for which there is a systemic problem that prevents effective use of the system. Hardware failures due to environmental changes in facility power and cooling systems that can be reasonably anticipated (such as brown-outs, voltage spikes, or cooling system failures) are the responsibility of the Subcontractor.

When a component has failed in service, the Subcontractor must replace the component with a newly manufactured or remanufactured/fully-tested component. The Subcontractor must not place a component back into the main system in order to determine if a failure is transient. With University's concurrence, the Subcontractor may use the test system to test components.

- **Software and Firmware Update Service**

The Subcontractor must provide an update service for all software and firmware provided for the duration of the Maintenance and Support Period. This must include new releases of software/firmware and software/firmware patches as required for

normal use. The Subcontractor must integrate software fixes, revisions, or upgraded versions in supplied software, including community software (e.g., Linux or Lustre), and make them available to the University within six (6) months of their general availability. The Subcontractor must provide prompt availability of patches for cybersecurity defects.

- **Call Service**

The Subcontractor must provide contact information for technical personnel with knowledge of the proposed equipment and software. These personnel must be available for consultation by telephone and electronic mail with University personnel. In the case of degraded performance, the Subcontractor's services must be made readily available to develop strategies for improving performance, i.e., patches and workarounds.

- **Problem Escalation**

The Subcontractor must document severity classifications and response for hardware and software problems. The description must include the technical problem-escalation mechanism based either on time or the need for more technical support in the event issues are not being addressed to the University's satisfaction. Problem-escalation procedures must be the same for hardware and software problems. Problems must be searchable in a database and made accessible via a web interface or for download in a standard format (e.g., CSV). This capability must be made available to all individual University staff members designated by the University.

- **On-Site Parts Cache**

The Subcontractor must maintain a parts cache on-site at NERSC. The parts cache must be sized and provisioned sufficiently to support all normal repair actions for two weeks without the need for parts refresh. The initial sizing and provisioning of the cache must be based on the Subcontractor's Mean Time Between Failure (MTBF) estimates for each FRU and each rack and scaled based on the number of FRUs and racks delivered. The parts cache configuration must be periodically reviewed for quantities needed to satisfy this requirement, and adjusted, if necessary, based on observed FRU or node failure rates. The parts cache must be resized, at the Subcontractor's expense, should the on-site parts cache prove to be insufficient to sustain the actually observed FRU or node failure rates.

- **On-Site Node Cache**

The Subcontractor must also maintain an on-site spare node inventory, for each node type, of at least 1% of the total nodes in all of the system. These nodes must be maintained and tested for hardware integrity and functionality utilizing the Hardware Support Cluster defined below if provided.

- **Hardware Support Cluster**

A test system (TS) described in Section 11.3 will be used as a hardware support cluster (HSC). The HSC must support the hot spare nodes and provide functions such as hardware burn-in, problem diagnosis, etc. The Subcontractor must supply sufficient racks, interconnect, networking, storage equipment, and any associated

hardware/software necessary to make the HSC a stand-alone system capable of running diagnostics on individual or clusters of HSC nodes.

B.10 Key Elements of the Center of Excellence for Workflow Readiness Plan

The Subcontractor must provide a Center of Excellence for Workflow Readiness Plan to assist in transitioning select NERSC mission workflows to the system (e.g., [NESAP](#) focuses on workflows that contain simulations, data, and learning components) within 120 days of the subcontract award.

The plan must be updated quarterly throughout the life of the project as needed. The plan must include the following:

- Named support staff provided by the Subcontractor and the CPU and GPU vendor(s)
- Staff training and deep-dive interactions with a set of teams
- How application and workflow developers can begin porting and optimization activities using proposed early access systems described in Section 2
- Mutually agreed-upon duration and level of effort. For example: at least 1 full-time equivalent (FTE) support must be provided from the date of subcontract award through two (2) years after final acceptance of the system.

B.11 Key Elements of the Training and Education Plan

The Subcontractor must provide a Training and Education Plan within 120 days of the subcontract award. The plan must be updated throughout the life of the project to reflect the latest content, as needed. The plan must include the following:

- A description of available training activities that target effective use of the user environment, performance, and optimization. The description must include topics, frequency, format (such as classroom training or online training, hackathons, etc.), and pricing.
- A description of collaboration with CPU and GPU vendors, other key technology providers, and NERSC staff where appropriate for the proposed training activities.

APPENDIX C: Facility and Site Integration Specifications

The Subcontractor must provide documentation in the **Site Preparation Facilities Plan** unless otherwise noted.

C.1 NERSC Facilities Overview

The NERSC-10 system(s) will be sited at the NERSC data center in Building 59 (Shyh Wang Hall) on the Lawrence Berkeley National Laboratory campus of the University of California in Berkeley, California, hereinafter referred to as Building 59. The building has four stories, with the mechanical plant occupying the lowest level. The single computer room is located on the second floor at the south end of the building. There are two office levels located above the data center. The computer room is split into two areas: the south computing floor and the common area. The south computer floor accommodates the HPC and auxiliary systems. The common area accommodates conventional NERSC computing systems.

The data center floor is approximately 680 feet above sea level. Building 59 has a dedicated truck loading dock with a dock leveler and a freight elevator to facilitate equipment deliveries.

C.2 System Physical Requirements

The NERSC-10 system will be located at the north end of the south computing floor, north of the *Perlmutter* system at Building 59.

The approximate area of white space available for the system footprint and aisle space is 4,784 square feet (46 feet east-west, by 104 feet north-south). The Subcontractor must provide system and auxiliary racks that fit into this available space.

The Subcontractor must coordinate the optimal placement of the system with the University for alignment with available power and water cooling connections.

The Subcontractor must provide PDF and AutoCAD physical documentation that includes a complete description of the physical packaging of the system, including dimensioned drawings of individual cabinets, the definition of a scalable unit (if applicable), cooling distribution units, auxiliary racks, the row pitch, and the floor layout of the entire system. Drawings shall also include weights, center of gravity information, frame construction, and seismic anchorage points. Provide information with revision control as the design updates are made. Prior to acceptance, send final versions of these documents.

After award and six (6) months before delivery, the Subcontractor must validate the system physical size, configuration, construction, and weight are compatible with the dimension, weight, utility distribution, equipment pathway, and delivery requirements listed herein.

C.3 Floor System & Weight Requirements

The access floor system is a conventional 24” (609.6 mm) tile system with a top elevation of +48” above the suspended concrete floor slab below. The tile pedestals are approximately 6” high and are supported on a 2” wide x 6” deep hollow structural steel (HSS) frame that is part of the facility’s unique seismically base-isolated floor system. Below-floor clearance is different from conventional raised floor systems, especially at the perimeter of floor tiles where there are structural steel framing members below in all locations. This configuration limits tile cuts and utility routing near the outside ~1” of each tile.

Existing solid access floor tiles are concrete-filled steel FS400 units manufactured by ASM. Tiles are CISCSA load rated for a 2,000 concentrated design load, passes a 2x minimum safety factor, 6,000 pound concentrated load, 800 PSF (note structural floor live load design is for 500 PSF maximum), 200 pound impact load, 1,500 pound 10-pass rolling load, and a 1,250 pound 10,000-pass rolling load.

The Subcontractor must provide shipped and operational weights of all equipment as part of the physical documentation. Equipment weights must not exceed the capacity of the existing access floor tiles. The Subcontractor must provide uniform operational loads beneath cabinets, or individual foot loads, as applicable.

The Subcontractor must provide dimensioned drawings of required access floor tile cuts based on the proposed layout of the systems for coordination and approval by the University 12 months prior to delivery.

C.4 System Power & Cooling Configuration

The NERSC-10 system design must support power and cooling connections provided below the access floor.

Power distribution from perimeter wall subpanels to computer equipment on the floor must use underfloor cable trays supported by the HSS structural steel base isolation frames to cabinet connections at the base of the cabinets (HPC and conventional racks).

Cooling water supply and return piping is located below the access floor, distributed to equipment by 8” diameter piping manifolds. 10” diameter piping risers penetrate into the mechanical space below the computer access floor on a regular grid of approximately 12 feet. Water-cooled racks or cooling distribution units must accommodate water feeds from below the computer access floor.

Air cooling is supplied by air handlers in the mechanical space that supply the air into a common plenum below the entire data center access floor. Conventional air-cooled racks must be configured for consistent air flow direction to accommodate creation of hot aisle or cold aisle containment strategies.

C.5 Seismic Requirements

The NERSC facility is located in a seismically active area. The NERSC-10 system must be placed on a seismic isolation floor.

The Subcontractor must ensure physical integrity of the equipment (racks and all components) will be maintained for earthquake accelerations up to 0.49g in any direction.

The Subcontractor must ensure that racks are physically interconnected near the rack base and top to provide monolithic response to each individual row of racks. Racks must be equipped with reinforced seismic anchorage holes in the base of the frames that allow for positive, direct anchorage to the bottom of the access floor tiles using threaded fasteners and bolts. Seismic anchorage locations must utilize four corner anchor locations at a minimum to maximize horizontal distance between holes.

The Subcontractor must provide PDF and AutoCAD drawings showing all frame dimensions and centers of gravity of all racked equipment for use in determining seismic anchorage, as described above. Provide final seismic information nine (9) months prior to system delivery.

C.6 Power Requirements

The NERSC-10 electrical system must be compatible with 480 VAC 3-phase power fed from delta-wye transformers. Substations consist of 480V 1200A-rated distribution breakers feeding 1200A distribution boards. The Subcontractor shall provide HPC system electrical pin and sleeve connections coordinated and approved by the University to the facility power system as 480 VAC 3-phase (with 4 or 5 wires) or 277 VAC single phase. Provide electrical system information after selection to facilitate final construction changes to prepare for system delivery.

The auxiliary system components are typically fed by 208 VAC 3-phase or 208 or 120 VAC single-phase power by power distribution units provided by NERSC. The Subcontractor must coordinate layout and power needs of auxiliary systems with the University for final approval of whip lengths and other electrical system coordination items. The Subcontractor must provide auxiliary system pin and sleeve connections coordinated and approved by the University similar to the 480 VAC connectors.

The University will provide facility power feed conductors from room perimeter wall panels underfloor to termination locations. Feeders will use pin and sleeve receptacle termination connectors. The Subcontractor must coordinate the type and configuration of the receptacle system with the University.

The NERSC-10 system must be resilient to incoming power fluctuations at least to the level guaranteed by the ITIC power quality curve and SEMI F47-0706 (Reapproved 0812) - Specification for Semiconductor Processing Equipment Voltage Sag Immunity. The Subcontractor must provide documentation demonstrating adherence to these documents six (6) months prior to delivery.

The Subcontractor must ensure multi-phase power or equipment with multiple power feeds is balanced across multiple connections and phases as equally as practical.

All powered equipment must be Nationally Recognized Testing Laboratories (NRTL) certified and bear appropriate NRTL labels. NRTL certification must be verified prior to acceptance by the University:

- For server-level products, the NTRL marks will be verified at the time of unboxing and must be readily visible on the component case.
- For rack, custom, and hybrid systems, all NTRL marks must be easily visible when equipment is installed, under power, and is in normal operating position. For equipment where this is not practical, the Subcontractor must provide documentation that each item of equipment has NTRL certification prior to installation, acceptable to the University. For example, the Subcontractor may provide a list of such equipment, item serial number, physical location (e.g., rack and U location), NTRL marks (e.g., name of testing laboratory and standard citations), and either photos of the NTRL mark, NRTL compliance documents, or reference to vendor documentation for the specific model that provides a statement of certification (e.g., provide the URL or attach the document).

All racks and system cabinets that are hard-wired or that cannot be disconnected safely via plugs must provide a means for zero-voltage verification (ZVV) or provide a documented procedure for zero-voltage verification acceptable to the University and included in the specification submission. ZVV should include the following:

- Finger- and tool-safe systems: All systems must be electrically finger and tool safe, meaning that internal electrical distribution above 48V must be guarded to at least IP2X (12.5mm or larger intrusions) in accordance with IEC 60529, Degrees of Protection Provided By Enclosures (IP Code).
- Short-circuit current rating (SCCR): the Subcontractor must provide and coordinate the short-circuit current rating capabilities for all equipment with the University to ensure the equipment is adequately provisioned for equipment protection.

For conventional, air-cooled racks:

- Utilize power distribution units (PDUs) for intra-rack power distribution.
- Use in-rack PDU strips per University specifications to be provided to the Subcontractor:
 - ServerTech C3WG36RL-DQJE2MT2 Switched POPS PDU (Or current generation in same model series)
- PDUs must use 60A cables with 208V Delta 60A IEC 60309 3P+G 9h connectors.
- Utilize dual, redundant feeds with failover capability for connection to utility and UPS power.
- Peak equipment power in a single rack must not exceed 40kW.

C.7 Cooling Requirements

The NERSC-10 system must be designed to operate within the following cooling conditions.

Facility cooling water (Primary Loop) conditions:

- Cooling water supply temperature range: ASHRAE W32, 5 °C to 32 °C (41 °F to 86 °F), Table 5.3 of the ASHRAE Liquid Cooling Guidelines, Fifth Edition.
- Cooling water supply flow rate: Total loop flow capability of up to 16,000 GPM. This flow rate will be shared with other water-cooled systems (Perlmutter) and a new air-

cooled heat exchanger array. Coordinate flow rate system design and needs with the University.

- Cooling water supply pressure: Up to 20 PSI differential pressure at the system cabinets.

Subcontractor cooling water (Secondary Loop):

- The Subcontractor must provide operational ranges for water temperature, flow rate, and differential pressure at the CDU or cabinet level as it applies to the delivered systems after selection and as the system design progresses.

Cooling requirements for conventional air-cooled equipment:

- Conventional air-cooled equipment must have front-to-rear air flow when installed in the rack.
- Conventional air-cooled equipment must be compatible with a supply air meeting ASHRAE A2 ranges, Table 2.1 of the ASHRAE Thermal Guidelines for Data Processing Environments, Fifth Edition. [TR-1] Ability to operate at ASHRAE A3 levels is strongly preferred.
- Conventional air-cooled system airflow requirements must not exceed an average of 125 CFM per 1KW of load.
- The peak thermal heat load transferred to the computer room air from the system, including all rack types and radiant losses, shall not exceed 1.2 MW.
- Auxiliary racks with sustained operational power above 20 kW may require additional liquid cooling methods.

C.8 Network Cabling

Network cabling (e.g., system interconnect) in Building 59 must run above floor and be integrated into the system cabinetry. There are existing cable trays that may be able to accommodate limited network cabling. NERSC-10 system network cabling must be approved by the University and meet these requirements:

- Permanent power, network, or other cables must connect to the rear of the unit. Temporary connections for configuration or debugging are permitted on the front.
- Power, network, and other cables must be neatly organized. Necessary cable management accessories must be provided by the Subcontractor.
- All network cables, wherever installed, must be source/destination labeled at both ends (see other sections for specific requirements).
- Where necessary, under-floor cables must be plenum rated and comply with NEC 300.22 and NEC 645.5.
- All network cables and fibers over 10 meters in length and installed under the floor must also have a unique serial number and dB loss data document (or equivalent) delivered at time of installation for each cable, if a method of measurement exists for cable type.

C.9 System Labeling

The Subcontractor must provide labels for every rack, network switch, interconnect switch, node, and disk enclosure with a unique identifier visible from the front and rear of the rack when the doors are open. The labels must be high-quality plastic so that they do not fall off, fade, disintegrate, or otherwise become unusable or unreadable during the lifetime of the system. The Subcontractor must provide documentation on labeling conventions and update the documentation when changes are made.

C.10 Delivery and Transport Pathway

The Subcontractor must follow the requirements for deliveries and transport of equipment through Building 59.

- Deliveries must be scheduled and coordinated for specific dates and times with the University during business hours. Deliveries must occur between 7:00 a.m. and 10:00 a.m. For trucks in excess of 25 feet, delivery planning must be completed at least two weeks in advance of the shipment from the origin. Final delivery dates and times must be established at least three business days in advance for the University to ensure that security processing is given adequate review time.
- Tractor trailer delivery trucks must be escorted by University personnel through the Blackberry security entry gates and onto the site. University flaggers will manage traffic control. Trucks need to back into the loading dock driveway. Trucks need to perform a three-point turn to leave the site with the assistance of University flaggers and a pilot vehicle. Box trucks and shorter tractor trailers follow similar protocols but may be able to forego pilot vehicle and flagger requirements.
- System equipment and components must be unloaded at the Building 59 truck loading building.
- The loading dock is an exposed, raised concrete slab structure. There is no rain cover or canopy structure at the dock. The dock is accessible by a single drive aisle located off of Chu Road at the north end of the building.
- A roll-up door provides access to a concrete slab mechanical level.
- A freight elevator is located at the north end of the building immediately adjacent to the loading dock roll-up door. System equipment must use this freight elevator to reach the second (Computing) level of the facility.
- The delivery path to the data center from the freight elevator is completely on raised access floor tiles. The path runs east then south through a cold shell space into the north end of the data center. A path south leads to the HPC floor area at the south end of the building. The total travel distance is approximately 250 feet.
- The access floor panels and moat plates beneath the load travel pathways shall be protected using aluminum plates, or similar. The University has a stockpile of roughly 20 aluminum plates ¼" x 48" x 96".
- The elevator cab floor must be protected by placing an aluminum sheet inside the cab. A second aluminum sheet must be used at each floor and placed across the elevator door opening when loading and unloading the cab to protect the elevator thresholds and avoid the load from getting caught in the gap between cab and threshold angle.

- The Subcontractor must stabilize rolling equipment movers by using brakes or wheel chocks in the elevator to prevent the load from shifting during travel. This will reduce the risk of the elevator seismic sensor being triggered.
- The pathway must not cross louvered or grilled tiles. The Subcontractor must coordinate swapping out louvered and grilled tiles with solid tiles with the University, as needed.
- The Subcontractor must prepare a documented transport and delivery plan as described in the **Key Elements of Delivery and Installation Plan** for University review three months prior to the first scheduled system delivery.

C.11 Building Pathway Weight and Dimension Constraints

The Subcontractor must ensure the system design is compatible with following dimensional constraints for delivery and building pathway logistics:

- The 44” high loading dock has a leveler that is 6’-0” wide and can accommodate a vertical adjustment range of -0” to +4”. The dock leveler can support a 20,000-pound load.
- The main roll up door (1201A) into the facility is electronically powered (locally) and measures 10’-7” wide by 12’-0” high.
- The freight elevator door opening measures 6’-6” wide by 9’-0” (this opening is the smallest of all openings).
- The elevator cab plan dimensions are 8’-4” wide by 8’-0” deep.
- The elevator is a C-3 freight elevator. It has a rated capacity of 7,000 pounds. The maximum allowable load in the elevator, including weight of crating, moving equipment, and personnel in the elevator cab shall be 6,300 pounds based on recommendations from the service vendor to avoid tripping the seismic motion sensor.
- There is a single interior door (2102) from the shelled space into the main computer floor. It is a double leaf door that measures 8’-0” wide by 10’-0” tall.
- The clearance between the top of the access floor and the utility systems on the bottom flange of the structural steel truss bottom chord is approximately 9’-6”.

C.12 Packaging Materials and Handling

The Subcontractor must follow the requirements for packaging materials and handling in Building 59 as follows:

- **Packaging Recycling:** Packaging must be completely recyclable and actually recycled by the Subcontractor in a reasonable time. No packaging materials must be left behind or unaccounted for. Short-term storage of packaging materials at the Mechanical level of Building 59 can be coordinated with the University.
- Storage of materials in the computer room areas beyond a single work day must be coordinated and approved with the University. Per National Fire Protection Agency (NFPA) requirements, the LBNL fire marshal must approve all short-term storage plans of materials in packaging.
- Use of metal or fire-proof storage bins and containers is strongly preferred.

C.13 Safety and Training Requirements

The Subcontractor must document all emergency shutdown capabilities and internal capabilities designed to protect the system from physical damage due to electrical system faults and mechanical overheating

The Subcontractor must coordinate and document safety information for staff required to perform work on-site at LBNL with the University. The University requires the following:

- Subcontractor Job Hazard Analysis (SJHA) to be developed in coordination with the University. This document covers the safety steps and protections taken for all work activities that are required for the system delivery and installation.
- Based on activities to be performed, the University requires LBNL training classes, or an approved, documented equivalent, reviewed by the University.
- Electrical lock-out, tag-out (LOTO): The University maintains a rigorous electrical safety program that requires formal training for any person who will perform work on or near equipment under LOTO configuration, including coordination with University representatives for compliance with safety procedures and processes.

Hot or energized work is prohibited at the University.

Energization of equipment shall be made in coordination with University Qualified Electrical Workers (QEWs).

The Energization processes include University verification of appropriate voltage and phase rotation checks. Data will be shared with the Subcontractor.

DEFINITIONS AND GLOSSARY

ASCR	Advanced Scientific Computing Research
AI/ML	Artificial Intelligence/Machine Learning
API	Application Programming Interface
Baseline Configuration	A documented set of specifications for an information system, or a configuration item within a system, that has been formally reviewed and agreed on at a given point in time, and that can be changed only through change control procedures.
BLAS	Basic Linear Algebra Subprograms
CI/CD	Continuous Integration/Continuous Deployment
CLI	Command-Line Interface
CPU	Central Processing Unit
CSI	Container Storage Interface
CUDA	Compute Unified Device Architecture
CVSS	Common Vulnerability Scoring System
DDP	DistributedDataParallel
FFT	Fast Fourier Transform
Full Scale	All of the compute nodes in the system. This may or may not include all available compute resources on a node, depending on the use case.
gid	Group identifier
GPU	Graphics Processing Unit
HDF5	Hierarchical Data Format, Version 5
HPC	High Performance Computing
Idle Power	The projected power consumed on the system when the system is in an Idle State.
Idle State	A state when the system is prepared to but not currently executing jobs. There may be multiple idle states.
I/O	Input/Output
IOPS	Input/Output Operations per Second
Job Interrupt	Any system event that causes a job to unintentionally terminate.
Job Mean Time to Interrupt (JMTTI)	Average time between job interrupts over a given time interval on the full scale of the system. Automatic restarts do not mitigate a job interrupt for this metric.

LAPACK	Linear Algebra Package
LLVM	A collection of modular and reusable compiler and toolchain technologies. https://llvm.org/
MPI	Message passing interface
NIST	National Institute of Standards and Technology
Node Failure	Nodes must be considered to have failed if a hardware problem, including soft errors, or a defect in supplied software causes the node to be unavailable, unable to operate correctly or unable to perform at established levels. A node must be considered to have failed if a node is administratively taken offline (“admindown”) to avoid erroneous operation, or by automatic action of system management software, for example a node health checker.
NRE	Non-Recurring Engineering
NVMe	Non-volatile memory express
OpenMP	Open multi-processing
PCI-e	Peripheral component interconnect express
PMI	Process Management Interface
POSIX	Portable Operating System Interface
QoS	Quality of service
RAID	Redundant Array of Independent Disks
Rolling Upgrades/Rolling Rollbacks	A rolling upgrade or a rollback is defined as changing the operating software or firmware of a system component in such a way that the change does not require synchronization across the entire system. Rolling upgrades and rollbacks are designed to be performed with those parts of the system that are not being worked on remaining in full operational capacity.
SC	DOE Office of Science
ScaLAPACK	Scalable LAPACK
Slurm	Slurm is an open source, fault-tolerant, and highly scalable cluster management and job scheduling system for large and small Linux clusters.
Software Bill of Materials (SBOM)	A formal record containing the details and supply chain relationships of various components used in building software.
SpiFFE	Secure Production Identity Framework For Everyone
SPIRE	SPIFFE Runtime Environment
SSH	Secure socket shell

System Mean Time Between Interrupt (SMTBI)	Average time between system outages over a given time interval.
System Availability	$((\text{time in period} - \text{time unavailable due to outages in period}) / (\text{time in period} - \text{time unavailable due to scheduled outages in period})) * 100$
System Initialization	The time to bring 99% of the compute resource and 100% of any service resource to the point where a job can be successfully launched.
System Outage	<p>The system should be classified as down if any of the following requirements are NOT met by the system:</p> <ul style="list-style-type: none"> ● 99% of compute nodes are available. This includes full health of all the node- to-switch links on a given node ● At least 98.5% of all inter-switch network links are operational at full bandwidth (i.e., degraded links are not counted as operational.) ● At least 85% of the bandwidth out of the system from within the system is still available ● At least 85% of the workflow environment nodes must be available to users. ● Ability to run benchmark applications as defined in Section 3 ● All mounted file systems are fully operational. In other words, all users should be able to access all of their data from mounting resources to which they have access ● Administrators are able to access by provided APIs or ssh to the system control plane, and from there are able to manipulate the system as needed ● System monitoring or auditing is functional
TCP/IP	Transmission Control Protocol & Internet Protocol
UCX	Unified Communication X https://openucx.org/
uid	User identifier