

**ADVANCED SCIENTIFIC COMPUTING ADVISORY COMMITTEE  
to the  
U.S. DEPARTMENT OF ENERGY**

**MEETING MINUTES**

**September 29-30, 2021**

**Videoconference**  
**ADVANCED SCIENTIFIC COMPUTING ADVISORY COMMITTEE**

The U.S. Department of Energy (DOE) Advanced Scientific Computing Advisory Committee (ASCAC) convened a Videoconference on Thursday, July 29, 2021 via Zoom. The meeting was open to the public and conducted in accordance with the requirements of the Federal Advisory Committee Act (FACA). Information about ASCAC and this meeting can be found at <http://science.osti.gov/ascr/ascac>.

**ASCAC Members Present**

Daniel Reed (Chairperson)  
Richard Arthur  
Keren Bergman  
Martin Berzins  
Vinton Cerf  
Barbara Chapman  
Jacqueline Chen  
Silvia Crivelli  
John Dolbow  
Jack Dongarra

Timothy Germann  
Roscoe Giles  
Susan Gregurick  
Anthony Hey  
Bruce Hendrickson  
Sandy Landsberg  
Richard Lethin  
John Negele  
Vivek Sarkar  
Krysta Svore

**ASCAC Members Absent**

Thom Dunning  
Gwendolyn Huntoon

Satoshi Matsouka  
Jill Mesirov

**Also Participating**

Alex Aiken, SLAC National Accelerator  
Laboratory  
Steve Binkley, Office of Science (SC)  
Betsy Bizot, Computing Research  
Association (CRA)  
Ben Brown, Advanced Scientific Computing  
Research (ASCR)  
Christine Chalk, ASCAC Designated  
Federal Officer, Oak Ridge Leadership  
Computing Facility (OLCF), Advanced  
Scientific Computing Research (ASCR)  
Jack Deslippe, Lawrence Berkeley National  
Laboratory (LBNL)  
Sandhya Dwarkadas, CRA

Barbara Helland, Office of Science for  
ASCR  
Tammy Kolda, MathSci.ai  
Pat McCormick, Los Alamos National  
Laboratory (LANL)  
Bronson Messer, Oak Ridge National  
Laboratory (ORNL)  
Rob Ross, Argonne National Laboratory  
(ANL)  
Erik Russell, CRA  
Amanda Stent, CRA  
Andrew Siegel, ANL  
Ceren Susut, ASCR  
Burçin Tamer, CRA

**Attending**

There were approximately 269 attendees in total for all or part of the virtual meeting.

**OPENING REMARKS FROM THE COMMITTEE CHAIR**, Daniel Reed convened the meeting at 11:01 a.m. Eastern Time and welcomed attendees.

**VIEW FROM GERMANTOWN**, Barbara Helland, Associate Director of the Office of Science for Advanced Scientific Computing Research; Ceren Susut, Computational Science Research and Partnerships Division Director; Ben Brown, Facilities Division Director

Speakers reviewed changes in ASCR personnel. Several staff members are serving in acting roles.

Reconciliation allocates ~\$10.3B for the DOE SC from FY22 to FY26. This figure includes ~\$7.8B for construction projects, with ~\$1.3B for the Exascale Computing Project (ECP), National Energy Research Scientific Computing (NERSC) Center, Argonne Leadership Computing Facility (ALCF), Oak Ridge Leadership Computing Facility (OLCF) and Energy Sciences Network (ESnet); ~\$1.5B for Major Items of Equipment (MIEs), with \$302M for the High Performance Data Facility (HPDF); and \$2B for Research, Development and Demonstration activities, with \$116M for the Computational Science Graduate Fellowship (CSGF) and \$340M to facilitate access of researchers to U.S. quantum computing facilities.

The House passed a stopgap funding measure on September 22, 2021 to maintain government operations through December 3. On September 27, 2021, the Senate blocked the proposed Continuing Resolution (CR) bill to suspend the debt limit through December 16, 2022.

The \$1T Infrastructure, Investment and Jobs Act passed in the Senate in August 2021 and includes billions of dollars in funding for applied Research and Development (R&D) and technology demonstration programs at the DOE and nearly \$3B for the National Ocean and Atmospheric Administration (NOAA).

The 2021 United States Innovation and Competition Act (USICA) passed the Senate on June 8, 2021. The ~\$200B bill includes funding for the Creating Helpful Incentives to Produce Semiconductors (CHIPS) for America Act and National Science Foundation (NSF) research. The NSF for the Future Act and the DOE Science for the Future Act both passed the House on June 28. The latter invests \$50B over five years in the DOE SC and national laboratories for renewable energy and emerging technology research. If passed, it would boost the DOE SC budget to \$8.8B in FY22, a \$1.8B increase from the enacted FY21 levels. The DOE SC's budget would reach \$11B in 2026.

FY22 Funding Opportunity Announcements (FOAs) will debut a new collaborator template to be submitted by the lead applicant. The SC is also working with SciENCv to deliver a biosketch template that can be linked to applicant's Open Researcher and Contributor ID (ORCID) accounts by January 2022.

Guidance updates to the Suggested Element of a Data Management Plan (DMP) and the new Guidance for DMP Reviewers will be effective for all solicitations issued after January 1, 2022. Reviewers must annually certify reading the guidance document. Solicitations have no changes to formal DMP requirements. DMPs are reviewed as part of the SC research proposal merit review process, and proposals may request funding to implement a DMP.

In FY21, 129 new awards were issued for 12 computational partnerships solicitations. Six DOE Nuclear Computational Low-Energy Initiative Scientific Discovery through Advanced Computing-4 (SciDAC-4) scientists, including two from ASCR, were selected for Early Career Research Program (ECRP) awards. Additionally, ASCR community members were conferred the Ernest Orlando Lawrence Award or American Association for the Advancement of Science

(AAAS), Society for Industrial and Applied Mathematics (SIAM) or Association for Women in Mathematics (AWM) fellowships in 2020 or 2021.

FY22 funding areas may include but are not limited to Data Visualization Beyond 4D; Storage Systems and Inputs/ Outputs (I/ O); Parallel Discrete Event Simulation (PDES); Explainable Artificial Intelligence (AI); High-Productivity Environments for Scientific Computing; Randomized Algorithms for Scientific Computing (RASC); Federated Scientific Machine Learning (ML); Mathematical Multifaceted Integrated Capabilities Center (MMICCs); and SciDAC Partnerships.

FY21 community events included the Terahertz-6G Wireless Communications Roundtable; DOE SC Community of Interest Workshop; RASC Workshop; Roundtable Discussion on Operating-Systems Research; Data Reduction for Science Workshop; SciDAC-4 Virtual Get Together; ASCR Workshop on Reimagining Co-Design; and ASCR Roundtable on PDES.

The ASCR Software-Stewardship Taskforce has met with key groups, including the former ASCAC subcommittee on *Transitioning ASCR after ECP*; ECP leadership; ASCR facilities leadership; and the Computational Research Leadership Council (CRLC). Meetings with other science funding agencies are upcoming and the task force is preparing a Request for Information (RFI) to seek feedback from the wider community.

Super Tech is improving gate speed and fidelity by optimizing the decompositions of quantum circuits for Berkeley Advanced Quantum Testbed's native gates. The Quantum Science Open User Testbed at Sandia National Labs (SNL) and the Advanced Quantum Testbed at LBNL (AQT@LBNL) are open to external collaborators, subject to merit review.

The ASCR's Leadership Computing Challenge (ALCC) Working Group highlighted significant ALCC program improvements. The ASCR Integrated Research Infrastructure Task Force published a white paper positioning ASCR for deeper engagement with non-ASCR user facilities and program offices. Requirements reviews found that ASCR facilities continue to innovate in leadership computing for clean energy, advanced manufacturing and biodefense efforts and ESnet High Energy Physics (HEP) and Fusion Energy Sciences (FES) programs.

In core operations, ALCF exceeded target metrics for systems availability, Innovative and Novel Computational Impact on Theory and Experiment (INCITE) hours delivered and capability hours delivered. Theta GPU user availability and production started in January 2021. Storage system growth investments continue. The modification to the Intel Build Contract for the ALCF-3 Project was completed in February 2021, and authorization was provided to Intel for phase 1, with completion anticipated in 2022. Installation of cabinets and the cooling system for the Polaris Testbed is on schedule, and Polaris will be available for testing and porting applications in 2022.

OLCF core operations likewise exceeded target metrics for Agency Priority Goals and system availability, delivering INCITE and capability hours. OLCF effectively managed Coronavirus Aid, Relief and Economic Security (CARES) Act funding to support COVID-19 research and maintained uninterrupted Summit operations and Frontier site preparation. A secure container, Citadel, was developed for personal health data. OLCF maintained user access to the Atmospheric Radiation Measurement testbed and quantum computing resources as well as other data visualization and processing tools. New high memory racks were added to Summit for COVID-19 applications, and Summit will continue to support INCITE, ALCC, and COVID-19 research in FY22. Frontier is currently being installed at ORNL, meeting the DOE Agency

Priority Goal of beginning the deployment of at least one Exascale Computing System by September 30, 2021. Hardware installation will be complete in October 2021, and Frontier will be available for testing and porting applications in 2022. Competitive allocations are scheduled for July 2022 with ALCC receiving priority access. General access through INCITE will be available in January 2023.

ESnet6 substrate has been installed across 300 locations with optical service fully transitioned from ESnet5 to ESnet6. ESnet5 decommissioning was completed in the first quarter (Q1) of FY21. Software automation of network workflows are fundamentally transforming deployment of ESnet6, with automation threshold Key Performance Parameter (KPP) to be completed by Q2, FY22. There has been significant progress in network deployment capacity and new routers. The network threshold KPP is expected to be complete by Q2, FY22.

The NERSC facility upgrade provided an additional 12.5 MW of power, along with the associated cooling and management infrastructure. The physical placement of the Phase I system was completed in April 2021 and a Perlmutter dedication ceremony held in May. The NERSC data center received the DOE Sustainability Award for Innovative Approach to Sustainability.

Selected Computational Partnership science highlights include contributions of FlexFlow and Legion to accelerating deep learning at scale and the Mochi project, which provides multiple data services in a flexible framework for exascale computing.

## DISCUSSION

**Cerf** asked about the HPDF budget and inquired if DOE will check whether DMPs are implemented. **Helland** obtained the \$302M figure from the legislative line. Regarding DMPs, **Reed** observed the Association of American Universities (AAU) and other consortia are discussing unfunded mandates and who must bear the costs of maintaining valuable data resources after funding ends.

**Reed** urged standard templates for applicant biosketches, collaborators, conflicts of interest and funding across federal agencies. **Helland** replied other agencies are using SciENCv. **Cerf** suggested a cross-agency effort to build a data management system. **Helland** replied that Brown's talk will offer such a vision.

**Cerf** sought an explanation for the 15x speedup observed with FlexFlow and Legion. **Hal Finkel** (ASCR) said the speedup primarily resulted from increased scaling. Legion is an asynchronous task-based environment that scales out to many nodes and covers execution on central processing units (CPUs) and graphics processing units (GPUs). FlexFlow was layered on top of Legion, enabling the training of more nodes to decrease the time to solution.

**Cerf** asked about the ESnet6 optical fiber backbone capacity. **Brown** explained that capacity is generally a terabit/s, with 15.5 terabits/s expected in aggregate. The quarterly capacity build-out per link entails hard decisions regarding resource allocation.

**Svore** inquired about preparatory activities at facilities for quantum machines. **Brown** said an OLCF quantum testbed has been paving the way for building future machines in-house through community collaboration. **Helland** clarified that the OLCF does not have a testbed but does offer access to quantum computers through the Quantum Computing User Program (QCUP), which is preparing scientific staff to support quantum users. While there is no coordinated program at present, future technological needs are being examined across facilities. New testbeds will be needed, and the future connection of a quantum accelerator to a high performance computer is likely along with other efforts to flexibly broaden the number of accelerators on machines following NERSC-10 and other upgrades.

**Arthur** asked about communication to clarify the roles of high performance computing (HPC) and quantum computing solutions, stating that quantum computing has limitations. As the ambassador of HPC to the government, ASCR is in an excellent position to communicate the areas where quantum computing has the most potential. Referencing a report by the National Academies of Science (NAS), **Brown** opined that there is a natural progression of quantum science to quantum technologies to engineered systems, but these technologies will not put HPC out of business. Aggregation of technology offers powerful opportunities. **Susut** reiterated that quantum will not replace but augment HPC. The SC views quantum computing as an exciting basic research challenge. The current five national Quantum Information Science (QIS) Centers are working with different industry partners to examine each step in the technology innovation chain that will link HPC and quantum efforts. **Svore** appreciated Arthur's comments. Microsoft has studied architectures and applications extensively and believes that scaling to industrial commercial value will come only after moving beyond 1M physical qubits. R&D is critical to this path's roadmap and it is wonderful to see investments from DOE and other agencies. With 72 qubits per chip, **Cerf** observed that the community is a long way from the 1M qubit goal. The 1M target is important because it might take 1K physical qubits to make a solid logical qubit that lasts long enough to provide results. Additionally, two wires are needed per qubit. Quantum networking is important if 1M qubits cannot be built on a single machine. However, entangling machines will require a quantum relay, which has not been developed. Thus, conventional HPC will be needed while quantum machines are developed. **Helland** acknowledged Arthur's remarks, adding that the Hill and Office of Management and Budget (OMB) do still ask why exascale machines are needed if quantum computing is coming. The SC consistently emphasizes that both are needed.

**VIEW FROM WASHINGTON**, Steve Binkley, Principal Deputy Director of the Office of Science.

Binkley reviewed the status of political appointees and SC staff. Secretary Jennifer Granholm has been active in articulating the vision of the Biden Administration. She engaged in decisions regarding the DOE SC FY23 budget. Deputy Secretary David Turk is on board and receiving updates on SC activities. Geraldine Richmond and Asmeret Berhe are the nominees for Under Secretary for Science and SC Director, respectively. Both nominees had their confirmation hearings in August 2021 before the Senate Committee on Energy and Natural Resources and are awaiting being voted out of committee to the full Senate. Chief of Staff Tanya Das recently took a position outside of government and Natalie Tham, Special Assistant, will soon rotate out of the SC.

Under the Biden Administration, the DOE Applied Energy Programs were returned to the purview of the Under Secretary for Science and Energy. The SC worked closely with the Applied Energy Programs on preparation of the FY23 budget which has been sent to the OMB.

The FY22 President's Budget Request (PBR) seeks \$7.44B for the SC, which is a 5.89% (\$414M) increase over FY21's PBR. The House Energy and Water Development Subcommittee issued a lower \$7.32B markup. The Senate markup is higher at \$7.49B.

SC activities are not covered in the Infrastructure bill, but the Reconciliation bill provides uplift for the SC. Congressional negotiations are ongoing.

FY21 ends at midnight on September 30, 2021. If Congress does not pass a CR, the SC and most national laboratories expect to operate on carryover funds. The DOE is also considering actions in case of a shutdown related to the U.S. debt ceiling.

## DISCUSSION

**Giles** requested a mechanism for receiving updates before the next meeting. **Binkley** will communicate to all FACA committees through the chairs. **Reed** agreed.

**Lethin** requested priorities from the FY23 budget. **Binkley** remarked on the overarching clean energy and climate priorities of the Biden-Harris Administration. Budget details are embargoed until released by the White House, typically in January. **Reed** added that members of the President's Council of Advisors on Science and Technology (PCAST) were announced last week. Additional high-level materials addressing administration priorities are available on the Office of Science and Technology Policy (OSTP) website. **Binkley** commented that the first of the monthly PCAST meetings is this week.

### **A VISION FOR ASCR FACILITIES**, Ben Brown, Facilities Division Director, Advanced Scientific Computing Research

ASCR's world-leading capabilities span supercomputing, data analysis, testbeds and data transport and are driven by scientists advancing the DOE mission while supporting users from academia, industry and other stakeholder groups. Today, ASCR is not only entering the exascale and ESnet6 era, but also a new era of complexity in scientific computing as highlighted in several recent ASCAC reports. Couplings between modeling/ simulation, experimental/ observational data, advanced algorithms, and AI/ ML tools now create complex workflows. Meanwhile, managing risk and opportunity in hardware choices has become increasingly challenging as technology evolves. These challenges are embedded in a broader DOE SC community context. Nearly every non-ASCR facility has users who leveraged an ASCR HPC facility, and ESnet connects to every SC user facility. ASCR computing resources and workforce expertise are in high demand. Yet, numerous talented individuals do not enter or are lost from the workforce pipeline. ASCR and DOE SC must implement a strategy to maximize DOE investments, yielding results greater than the sum of the parts.

A future ASCR facilities vision pictures each facility thriving and possessing agency, but also collectively working to advance scientific computing across the DOE and beyond by driving the state-of-the-art with the ASCR research and vendor communities; catalyzing discovery and innovation; responding to national needs; delivering on stakeholder priorities with balance and equity; fostering scientific ecosystems; broadening the diversity of individual, institutional and domain participation; and demonstrating excellence in project management and program operations. To realize this vision, a shared salience must emerge from user programs; collaborations for priority DOE and national needs; and scientific ecosystems. Continuous system evaluation and reciprocal feedback will guide investments and determine the impacts of decisions. To advance this vision, user programs must continue efforts to broaden their user base; collaboration portfolios must be visible and consider new models; ecosystems must be defined and made visible while fostering community participation and sustaining careers; and shared salience activities must capture and aggregate insights via annual synthesis. Resulting shared salience will enable stakeholders to systematically manage risk and opportunity.

Incipient scientific ecosystems are already emerging through the interconnected DOE SC user facilities, even as ASCR facilities contemplate the operational implications of integrating computing across experimental and observational facilities. The ECP generated a software ecosystem that greatly reduces barriers for ASCR fundamental research maturation and impactful delivery. Early shared salience from user facilities has emerged through review

activities that produced the Exascale Crosscut Report and the Lessons Learned from the COVID Era Report. Indeed, the arc of the DOE's history bends towards open ecosystems; the DOE was conceived in the secrecy of the Manhattan Project but has moved towards ever more open practices, recently culminating in the National Virtual Biotechnology Laboratory (NVBL).

## DISCUSSION

**Reed** observed that complexity offers diverse challenges and opportunities.

**Cerf** commented that the challenges created by heterogeneity in underlying hardware reach upwards through the stack, resulting in new challenges for generating compilers and operating systems. **Brown** agreed, emphasizing the broad nature of these challenges and communitywide opportunities for shaping future operational paradigms.

**Sarkar** appreciated the presented vision. Not only high complexity, but also high uncertainty present challenges. DOE SC must prepare for unknown outcomes, which means taking lots of risks rather than betting on one approach. Per earlier remarks, compilers can no longer function as workhorses, but must innovatively manage heterogeneity.

**Cerf** said the ecosystem metaphor accurately captures the unity and complexity of environment. **Brown** recognized that scientists are currently contending with uncertainty and complexity; the risks taken by individual researchers in surrendering control of portions of computer code must be aggregated across scales to inform shared salience.

**Sarkar** continued discussion of complexity and uncertainty, recalling an atmosphere of high anxiety when industry shifted from vector to NVIDIA Performance Primitives (NPP) followed by a period of company expiration. DOE must be prepared to innovate. When exploring, research groups often assemble their own testbeds, only later realizing they are not equipped to manage them. DOE SC user facilities have the opportunity to stand up testbeds and supply staff knowledge with the understanding that the level of service may be less than that at user facilities. Through this interface, facility staff can advise scientists on the most appropriate types of testbeds. **Brown** agreed. To be most effective, such messaging may emerge from many community locales. Currently, many view the user programs as the interface for interacting with the user facilities. However, community pockets are evolving towards a more organized approach that will enable researchers to advance workflows beyond those realized by individuals' limited resources.

**Giles** reflected on ecosystem interfaces with universities and emphasized the importance of supplying technical expertise at all levels of the management chain as discussed in the Exascale Transition Report. What will be the role of ASCAC and other FACA committees when this vision is realized? **Brown** called attention to the competition of ideas stewarded by the DOE. If the presented vision is implemented, researchers must surrender some autonomy to participate. Such decisions must be undertaken at a community level, and FACA committees can help ensure discussions are visible while offering guidance. **Giles** commented on building direct collaborations across the DOE SC and gathering perspectives across offices, as done when compiling information for the Exascale Transition Report. Changes will also need to occur at the university level for labs to engage researchers who do not know how to get started. **Brown** resonated with the message that moving forward will require all voices to be heard. Synthesizing directions cannot be burdensome.

**Bergman** stressed the need for testbeds that take advantage of emerging heterogeneous technologies. Scientists lack an intermediate-scale way of advancing research at the technological level. This is a difficult challenge, and considering how to implement interfaces



with larger systems is worthwhile. **Brown** concurred. The ASCR portfolio must offer a coherent strategy and provide multiple paths for entry. The Advanced Computing Technology (ACT) Division, with other partners, dreams of interweaving mainline user facility programs with testbeds of the future. Ongoing conversations with National Nuclear Security Administration (NNSA) and Small Business Innovation Research/ Small Business Technology Transfer (SBIR/ STTR) partners are exploring the best models for renewing vendor engagements.

Also calling attention to testbeds, **Reed** asked how to expand the ecological space to allow for a risk continuum. It is hard to take crazy risks when on the hook to deploy a productive infrastructure. The coming environment will likely be more heterogeneous than that fostered by the ECP technology stack. **Brown** said that the federal U.S. R&D system is manifold, spanning many agencies and operational models, conferring an international competitive advantage. At the same time, agencies must maintain awareness of what others are doing to balance risk and opportunity within their own portfolios. **Helland** concurred, noting that ASCR must foster partnerships across the SC, especially with BES, to track the new materials and science that will give rise to the next generation of computers.

**John Shalf** (NERSC) inquired about the complementary roles that facilities, research and ACT can contribute to ASCR's continuous improvement. **Brown** does not yet have a blueprint. Roles must be determined to fuse where the world is to where it is going.

**Sarkar** emphasized frequent failure as a basic tenet of science. Industry has the motto of failing fast in order to succeed sooner. To encourage this attitude, ASCAC could invite presentations from those that stood up testbeds in order to learn what works. Success should be tied to exploration.

**Reed** dismissed the meeting for a break at 1:32 p.m. and reconvened at 1:45 p.m.

**ECP UPDATE**, Andrew Siegel, Applications Development Director, Argonne National Laboratory

The \$1.8B, 7-year ECP project will conclude in 2023. It has engaged DOE labs, universities and industry in Application Development (AD) spanning mission critical science and engineering across nine DOE program offices and the NIH. AD comprises 24 applications and six co-design projects. Its 78 separate codes were mostly started with MPI or MPI+OpenMP on CPUs. Each AD project is defined by a challenge problem and characterized by algorithmic innovations that are not focused on benchmarks.

The current AD porting focus is on diverse CPU/ Multi-GPU and GPU-resident hardwares. Early hardware is available from Intel, AMD and NVIDIA. Each application has its own figure of merit (FOM) with performance measured relative to baseline values. Thus far, six key ECP applications have surpassed a FOM performance measure of 50 on Summit.

Porting applications to new hardware to achieve GPU speedups has been non-trivial. Challenges have included evolution of new fundamental methodologies; navigating computational motifs preferred by GPU hardware; time needed for programming models/ analysis tools to mature and obtain community buy-in; and development of application-level libraries. GPUs perform best for codes that are characterized by fine-grain parallelism; can be made GPU-resident; can operate in the weak scaling regime; have high arithmetic intensity; can be formulated as single instruction, multiple data (SIMD) processing with minimal branching logic; require extreme performance with a relatively high floating point operations per second (flops) to byte (of storage) ratio; and can make use of specialized tensor core instructions.

Progress in addressing GPU challenges presented by application case studies (ExaBiome: tools for metagenomics; ExaFEL: real-time particle imaging from light sources; ExaSGD: tools for the power grid; and ExaSMR: tools for nuclear reactor design) illustrate how hardware and programming models can drive algorithms and methods as much as the reverse. Finally, since the beginning of the project, there has been significant movement in languages and GPU programming models, mostly towards C++ and abstraction layers/ libraries, respectively.

Emergent themes from the 2020 ECP AD report include use of mixed precision; strong scaling; optimized libraries on early access machines; performance of OpenMP offload; GPU-resident and unified virtual memory; and relative increased cost of internode communication.

## DISCUSSION

**Berzins** requested additional information on how AD will meet the challenge of moving between CPUs and GPUs with different programming models. Are there attempts to standardize? **Siegel** said the ecosystem is currently self-organizing with no clear winner. Many GPU programming models are in use, including Kokkos, RAJA, Hip and CUDA as well as OpenMP. The latter presents a risk because not much is known about OpenMP offload, especially on new machines. However, many PIs prize portability as a top priority although FOMs are used for project metrics. **Berzins** remarked that OpenMP appears to have potential for running well on Aurora. The ecosystem for Aurora and Frontier is less well-defined, and there must be some attempt at achieving standard portable performance across architectures. This situation is in some ways analogous to work conducted with MPI 30 years ago.

**Cerf** inquired about open source library validation. **Siegel** said more effort is needed. Good open source libraries, like the math libraries, have lots of verification problems built in. Other libraries are more boutique with less engineering. Open source libraries come with many risks ranging from verification to sustainability. Consequently, people often use libraries from trusted, collaborative sources. **Doug Kothe** (ORNL, via chat) added the ECP is making a concerted continuous integration/ continuous delivery (CI/ CD) effort to harden and improve the quality of its libraries.

## SATISFYING THE DATA DEMANDS OF DOE SCIENCE, Rob Ross, Argonne National Laboratory

DOE science exhibits many types of data models, interfaces for data access and data organizations. The Mochi project, an R&D 100 finalist, is a state-of-the-art open source tool for the rapid development of customized data services involving HPC, big data and large-scale learning. Mochi's plug and play components for filtering, sorting and processing data fall into three categories (Core, Utilities and Microservices) to address needs spanning HPC, distributed computing and cloud computing. For example, DeltaFS is a data service for kinetic-plasma simulations that leverages Mochi components and LevelDB technology from Google to improve query speed by three to four orders of magnitude. The High-Energy Physics Event Store (HEPnOS) temporarily stores event data, accelerates the put/get of event data and integrates easily with analysis tools to support Neutrinos at the Main Injector Off-Axis Electron Neutrino Appearance workflows. HEPnOS is also being integrated into the Imaging Cosmic and Rare Underground Signals (ICARUS) workflow. National laboratories, universities and two industry partners (Intel and Kitware), have used Mochi technology to build several data sources. Users have adopted tools because they perform well, are portable and customizable.

The 2018 workshop in Storage Systems and I/O underestimated the potential for nonvolatile memory and the complexity of exploiting it. The same workshop was too skeptical that smart devices would appear in the near term. Research is needed to improve I/O operations per second (IOPS) to pave rapid pathways to services, with high-speed networks and nonvolatile memory providing a basis for solutions. Hardware trends have led to smart devices with capabilities much closer to those of CPUs; there is potential to offload service capabilities to such devices rather than to application tasks.

Data management can both benefit and benefit from AI/ ML. HPC can accelerate AI by providing faster access to data. New data formats offer opportunities to match AI workflow needs while capturing information enabling Findability, Accessibility, Interoperability and Reusability (FAIR) data principles. AI may also improve initial configurations for data services, assist in adapting these to changing workloads and environments and detect anomalies. Encouraging communication between the “Storage I/O” and “Data Management” communities presents rich opportunities to advance concepts and technologies. Examples include annotations capturing workflow details for reproducibility; “smart stores” connecting relevant datasets at the exabyte scale; and “collaborative stores” sharing information across institutional boundaries.

## DISCUSSION

Citing earlier discussions of testbeds and prototypes, **Reed** inquired about smart devices. **Ross** recommended building networking and storage components into the earliest testbed installments. Having early technology access, as ECP participants did, can make a big difference.

**Pouchard** (BNL, chat) asked about bridging I/O and data management worlds. **Ross** cited language as a challenge. Terms like provenance can mean different things to those focused on storage versus data management. A venue for discussion, such as a workshop, could help communities to identify priorities and move forward. The storage community can contribute algorithms, scaling and architectures to solve problems while the data management community can inform solutions to reproducibility and other FAIR challenges.

**Reed** dismissed the meeting for the day at 3:02 p.m.

## THURSDAY, SEPTEMBER 30, 2021

**Reed** convened the meeting at 11:00 a.m.

**COVID-19 HPC CONSORTIUM**, Bronson Messer, Director of Science for the Oak Ridge Leadership Computing Facility, Oak Ridge National Laboratory

The COVID-19 HPC Consortium is a public-private consortium that supplies researchers with access to the world’s most powerful HPC resources to advance scientific discovery in the fight to stop the virus. OLCF has supported this endeavor from the Consortium’s inception in March of 2020, joining partners from DOE national laboratories, academia, industry, affiliates, other federal agencies and international entities. To date, the Consortium has funded >100 proposals in basic science (e.g., viral structure, viral-human interactions, viral evolution, environmental effects and tools); therapeutics (e.g., target discovery, small molecule design, protein design, drug repurposing and technologies for development); and patients (e.g., trajectory and outcomes, medical technology, supply chain, epidemiology and detection).

The OLCF has provided a unique set of capabilities for a variety of projects, allocating 1.42M Summit node-hours to date to 11 projects. Several of these projects have also leveraged resources provided by other Consortium members (e.g., Google, IBM) and the projects enjoy queue priorities equal to those of INCITE projects. Summit's unique capabilities enable it to excel at problems related to viral structure, small molecule design and drug repurposing. CARES Act funding received in July 2021 allowed the installation of 54 new nodes that can be run as a separate partition, enabling jobs that require larger on-node memory. This installation has greatly benefitted several research areas, including computational chemistry, molecular dynamics and AI for medical imaging analysis.

Highlighted studies include the role of glycosylation in SARS-CoV-2 S-protein conformational dynamics; prediction of synergistic drug combinations for treatment of COVID-19; physical models of COVID-19 related proteins; and structural modeling of COVID-19. An AI workflow to simulate the virus's spike protein in numerous environments, including in the viral envelope comprising 305M atoms, received the Gordon Bell Special Prize for HPC-Based COVID-19 Research. More projects are anticipated in the coming months. The OLCF will continue to deploy computational and data resources to tackle global challenges.

## DISCUSSION

**Cerf** asked whether messenger (mRNA) vaccines can trigger the immune system to recognize the human ACE2 receptor, the viral spike protein's target. **Messer** did not know.

**Hendrickson** requested information about institutionalizing HPC resources for unanticipated national needs. **Messer** stated that the Consortium demonstrates proof of concept and could serve as a model. Moving forward will require an understanding of all members' needs, policies and missions. HPC resources are valuable, and creating a strategic reserve should not mean holding resources back but deploying them quickly. **Helland** said the Consortium's response to a recent RFI is under discussion. The Consortium is holding a workshop series to understand how HPC can support the response to a variety of national disasters.

**Cerf** asked about animating the assembly of SARS-COV-2 viral particles inside a host cell. **Messer** said scientists have an understanding of the most meaningful assembly stages. The importance of animating the entire process is an open question; equating structure with function and understanding the serialized dynamics of key components may be more critical.

**Cerf** inquired whether a commonly supported abstraction layer would be necessary to allow users to move computations between clouds. **Messer** affirmed. This is under discussion. The process should be made as simple as possible for end users.

**Reed** posed a question about broadening the user base to non-HPC experts. **Messer** replied the Consortium opened the OLCF's aperture to new project topics and engaged users that would not normally think of using HPC resources.

## REPORT FROM COV, Alexandra Landsberg, COV Chair, Office of Naval Research

The ASCR Research Committee of Visitors (COV) was charged with reviewing the Applied Mathematics, Computer Science, Computational Partnerships and Research and Evaluation Prototypes (REP) programs for FY16-FY19. The international and institutionally diverse COV met virtually August 18-19, 2021. The report is currently in draft form.

Key findings pointed to 1) The impact of the ECP as well as new AI/ ML and QIS efforts to ASCR's research portfolio, related funding cuts and the need for a holistic budget plan; 2) The need for clarity in how programmatic shifts are made and communicated; 3) The ECRP's

positive impact on awardees' careers and the negative impact on award numbers of the congressional mandate to fully fund awards <\$1M; 4) The need for increased use of preproposals; and 5) The COV's need for additional time to digest dense presentations and expert help in navigating Portfolio Analysis and Management Software (PAMS).

Key recommendations suggest that ASCR Research 1) Identify and document their North Star, including a clear vision and mission statement with an accompanying five-year plan; 2) Develop procedures to better communicate the impact of programmatic shifts; 3) Investigate strategies to identify early (and early-mid-career) researchers with significant promise and ways to enable them to develop into PIs of large DOE projects; 4) Implement a preproposal process to reduce the community burden of writing and reviewing proposals with little chance of being funded; and 5) Provide solicitation summary statistics in COV presentations to facilitate evaluation.

Applied Math findings addressed solicitation decisions, proposal procedures and quantum testbeds. Comments touched on ECRP awardees, navigation of funding challenges and emerging research areas. Recommendations discussed use of preproposals, mechanisms to diversify PIs and effectiveness measures for math centers and long-term lab projects.

Computer Science findings referred to solicitations, award numbers, program management, the impact of limited funding on the ASCR portfolio, university principal investigator (PI) opportunities and workshops. Comments dealt with unsolicited awards, R&D research risk and workshop fatigue. Recommendations considered communication, PI diversification, identification of emerging technologies and program assessments.

Computational Partnerships findings reviewed solicitation workload, the reviewer pool, responsive program management, inter-office and -program demand for partnerships and SciDAC's impact and software migration to exascale platforms. Comments referred to the SciDAC-4 Coordination Committee and documentation of feedback loops from partnerships to ASCR base research programs. Recommendations considered preproposals, communication of SciDAC strategic goals and external SciDAC review.

REP findings discussed solicitations as well as program activity and progress. Comments remarked on award recipients, program flexibility and CSGF growth. Recommendations dealt with encouraging scientists to utilize quantum testbeds and expanding CSGF diversity.

## DISCUSSION

**Chapman** appreciated the recommendation to support early career researchers. They face many challenges, and much talent may be lost if they are not able to fully contribute.

**Lethin** requested more information about the recommendation to increase communication. **Landsberg** explained that the ECP caused many programmatic shifts. Communication of these changes was variable, potentially leading to different levels of awareness for individuals at labs, in academia or in industry. If opportunities are ramping down, ASCR can direct the community to resources within the DOE or other agencies.

**Arthur** (chat) observed that GE has adopted nomenclature similar to the recommended North Star.

**Bergman** sought further explanation of the North Star recommendation. **Landsberg** said ASCR priorities shifted to accommodate ECP, AI/ ML and QIS initiatives. There was also program officer turnover during this time, and it was unclear why different officers made different decisions. Clear, centralized documentation outlining priorities was lacking. Moving forward, ASCR can incorporate new emerging areas as focal points in its North Star.

**Sarkar** raised the reduction in research funding coinciding with the start of ECP and inquired if the COV identified budget trends in FY20 and FY21 that were more reassuring. **Landsberg** replied that while the COV is not allowed to comment on recent trends, there was consideration of actions ASCR might take to support rising research stars if funds were restored. The COV plans to elevate an Applied Math recommendation to re-establish the small group and single PI program. Those small teams in academia and industry feed the entire ASCR pipeline; if dropped, this talent may be lost forever. **Sarkar** appreciated recommendations to support PIs from early to mature career stages. In academia, PIs can work up from small NSF grants to large projects. It would be great to have comparable options for ASCR researchers.

Citing North Star discussions, **Cerf** cautioned that while it is important to know organizational directions, it is also important to leave room for innovation. **Landsberg** agreed. The COV encourages ASCR to consider emerging research areas.

**Giles** commented on the overlap between COV recommendations and the ECP report. The ECP shifted some scientists' research focus from longer-term horizons to projects with more immediate relevance. **Landsberg** highlighted recommendations to develop measures of success for long-term programs, providing a way to showcase the value of long-term research.

Noting funding limitations during the review period, **Cerf** revisited discussion of HPC heterogeneity. ASCR must prepare for the even more diverse supercomputing engines of the future. **Chapman** reinforced Cerf's remarks, highlighting the need for longer-term research investment. **Sarkar** agreed. Some DOE researchers left during this time, and lack of long-term research investment may have contributed to attrition. **Landsberg** concurred. The COV cannot offer advice on ASCR investments but can comment on how funding impacted programs. The COV may elevate remarks about re-establishing the small group and single PI program to a key recommendation.

**Reed** observed that blue sky explorations are needed. If constrained to ideas that are only compatible with existing systems, then a large fraction of the idea space has been precluded. Incrementalism and low-risk studies predominate when funding is tight. **Landsberg** said the EXPRESS program receives blue sky ideas, but a better down-select process is needed.

**Arthur** inquired about the projected inclusion of industry awards relative to university contracts. In the context of AI, data has tremendous value as a raw material in industry; though work may be conducted in a proprietary context, practices can still be advanced. **Landsberg** stated FYs 16-19 were difficult. There was not a large number of independent industry awards, but there were partnerships in quantum testbeds and prototypes. **Helland** explained that ASCR transitions research to industry through the SBIR program. Prior to 2016, most vendors were focused on developing HPC components and software. Vendor programs continue to support industry relationships with facilities. A current project with strong industry partnerships is examining how to maintain data privacy on big machines.

**Cerf** asked if entangled photons can be sent in the ESnet6 optical infrastructure. **Helland** said ANL, Fermilab and BNL have been trying to use existing fibers. **Brown** concurred, stating that marginal fibers are being used. The issue of quantum repeaters is a basic research challenge. **Kleese van Dam** (BNL, chat) said BNL demonstrated, with ESnet support, that different strands of an existing fiber connection can transport quantum and classical information at the same time. However, interference levels suggest combined communication is not desirable. In tests, the classical communication was used to remote operate and synchronize the quantum hardware. **Brown** (chat) clarified that the recent quantum/ ESnet collaborations at BNL and Fermilab/ ANL

did not use the backbone network, but depended critically on ESnet's infrastructure presence at those labs and the physical presence of an ESnet employee at the BNL site.

**Giles** inquired whether the COV had any systems recommendations for PAMS.

**Landsberg** relayed that PAMS was great for university proposals. Lab proposals are still manually uploaded. PAMS is a valuable system, but it would be helpful for ASCR to provide user support for outsiders, especially during COV activities.

**Sarkar** discussed Fast Forward and Design Forward investments to industry for accelerating technology transition and suggested ASCR conduct a retrospective review of their impact. **Arthur** distinguished computing industry investments for leadership facilities from those that support industry in advancing the DOE mission. **Sarkar** agreed, noting many potential AI/ML partners, including those in the life sciences. **James Ang** (Pacific Northwestern National Lab) said conversations via the Semiconductor Research Corporations Decadal Plan Committee show the computer and microelectronics community invests in R&D. DOE programs like Path Forward cover gaps between fundamental research and advanced product development for mission needs. Industry would like to support future Path Forward efforts, but it does not appear to be in the budget.

**Chen** suggested leveraging ECP application codes as testbeds for *in situ* as opposed to data-driven AI/ML methods. Applying FAIR principles to data generated from these codes would support development of future algorithms enabling more complicated workflows.

**Hendrickson** voiced concerns about the oversubscription of ASCR calls and low proposal success rates. He invited solutions beyond implementing pre-proposals, noting that having more money to distribute would help. **Landsberg** agreed; in Applied Math and Computer Science especially, the large number of pre- and in some cases full proposals submitted is burdening the community. Researchers do not propose when there is little chance of funding. The COV recommends providing a thoroughly documented review process with down-select mechanisms. Budgets are going back up, and hopefully there will be more future opportunities. **Reed** commented on the similarly low NSF proposal success rates. The message this sends to the research community is worrisome. The country must address support for basic research.

**Reed and Helland** thanked Landsberg and the COV for their work.

## **RANDOMIZED ALGORITHMS FOR SCIENTIFIC COMPUTING (RASC) WORKSHOP REPORT**, Tammy Kolda, MathSci.ai

Some of DOE's instruments and tool suites are now producing data at rates that outpace current abilities to process and move data, necessitating randomized algorithms for efficiency. These methods employ randomness in making internal decisions, generating sketched data from the original, and are de rigueur in AI/ML. Randomized algorithms accelerate time to solution; increase scalability; improve reliability; and expand impact across the DOE.

The virtual two-part w engaged 453 participants in a RASC Bootcamp in December 2020 and 204 participants in a Brainstorming and Writing session in January 2021. Both the two-part and online format worked well, and effort was made to make the meeting engaging and inclusive. Indeed, the high Bootcamp attendance illustrates community appetite for RASC knowledge.

Attendees of the second portion identified six challenge/ opportunity areas: 1) Randomized algorithms to enable future computational capacity; 2) Novel approaches by reframing long-standing challenges; 3) Randomness intrinsic to next-generation hardware; 4)

Technical hurdles requiring theoretical and practical advances; 5. Reconciliation of randomness with user expectations; and 6) Need for expanded expertise in statistics and other areas. Recommended future research directions address theoretical foundations; algorithmic foundations; application integration; next-generation architectures; interdisciplinary research/outreach; and workflow standardization. The workshop report, *Randomized Algorithms for Scientific Computing*, is available on arXiv.org.

Following the meeting, the EXPRESS Randomized Algorithms for Extreme-Scale Science FOA was posted in March of 2021. Six awards totaling \$2.8M, one addressing each recommended research direction, were made in August of 2021.

## DISCUSSION

**Cerf** (chat) asked about sketched data. **Dongarra** (chat) said the process is similar to down-sampling to obtain simpler or lower-rank approximations. **Kolda** added the term typically means a random linear combination of data observations. There are also nonlinear sketches.

**Svore** (chat) appreciated Kolda's comments about engagement. Online or hybrid formats present opportunities for being more inclusive and encouraging more diverse participation.

**Lethin** requested more information about the relationship between randomized algorithm compressed sensing to information and communication theory and how to draw upon U.S. leadership in 5G or 6G. **Kolda** explained the signal processing community generated much of the technology that has been important to randomized algorithms, including compressed sensing. This topic presents a tremendous opportunity for the Applied Math program, and DOE has the opportunity to show leadership in this area, perhaps establishing this field as a North Star. **Lethin** urged ASCR to tap into the Department of Defense's (DOD's) and telecom industry's expertise, which show federal and commercial leadership, respectively, in signal processing.

**Dongarra** remarked that randomization presents a sea change in the use of linear algebra and agreed DOE has an opportunity for a leadership role. However, ~\$3M is a small amount of funding for six projects. **Steven Lee** (chat, ASCR) clarified that projects are two years in duration. **Kolda** said \$2.8M is hopefully the start of more funding to come. It would have been ideal to fund a much larger set of small projects with more representation from the statistics community. Two of the funded projects begin to address statistics topics, but the community will not apply if their funding chances are low. Future funding needs to engage larger projects.

**Landsberg** concurred that DOE could take a leadership role in this area and advised engaging other agencies. The DOD Office of Naval Research (ONR), for example, has a small math budget and could address topics in theoretical foundations. There is opportunity for community groups to work together on different pieces. **Kolda** agreed. ONR, NSF and DOD have been involved, and it has been exciting to engage with the theory community to begin addressing some of DOE's unique challenges.

**TAULBEE, TALENT AND TRENDS** Betsy Bizot, Sandhya Dwarkadas, Erik Russell, Amanda Stent and Burçin Tamer, Computing Research Association

The nonprofit Computing Research Association (CRA) has >200 North American member organizations spanning academia, professional societies, industry and government. CRA conducts research to effect beneficial change for the computing research community and society at large.

Since 1970, the CRA Taulbee Survey has been issued to all U.S. and Canadian institutions offering doctorates in computer science (CS) and has tracked key national trends.



From 2004-2019, the annual number of CS undergraduate degrees conferred by U.S. doctoral institutions fell from ~17K to its lowest point of ~8K in 2009 due to the dot-com bust and then surged to the current high of ~32K in 2020. Surveyed institutions confer ~1/3 of nationwide CS degrees, and >80K undergraduate degrees total were awarded nationwide in 2020. This boom in undergraduate degrees is not translating to increased PhD program enrollment which has remained relatively level over the last 15 years. During this period, the majority of incoming students have been international with temporary visas. Recently, the percentage of women has hovered at ~20%. When considering 2020 PhD enrollment by residency, ethnicity and gender, international students outnumber domestic students. Within the domestic population, White and Asian students outnumber Black, Hispanic and Native American students. Men outnumber women in all categories. In 2020, the number of graduating PhD students specializing in AI/ ML far outstripped numbers in other specialty categories, and enrollment trends from 2010-2019 show the largest percent increases in Human-Computer Interaction, Robotics/ Vision and AI/ ML specialties while Programming Languages/ Compiles, Software Engineering and Networks show the largest percent declines.

The CRA-Widening Participation (WP) Committee aims to broaden the participation and improve the access, opportunities and positive experiences from underrepresented groups (URGs) in computing research and education. The CRA-WP Board consists of ~30 volunteers that are practicing CS researchers. They work with additional volunteers and CRA-WP staff to support research, career mentoring, community building and recognition programs targeting community members from undergraduate to mid-career stages. Programs are supported by funding from DOE, other government agencies, industry, academia and foundations.

CRA's Center for Evaluating the Research Pipeline (CERP) is a research and evaluation resource for the CS community. Established in 2012, CERP's work assesses the effectiveness of >20 programs by CRA-WP, select Broadening Participation in Computing (BPC) Alliances and other organizations. Evidence-based insights are used to improve existing programs and create new ones. Since 2011, CERP's Data Buddies Project has conducted an annual national survey of undergraduate and graduate students to inform research evaluating education. More recently, the survey has included alumni and professionals in higher education, and data on departmental policies will soon be collected. Survey data illustrates that URGs in CS feel significantly lower levels of a sense of belonging. Women and non-binary gendered computing professionals report lower confidence in their ability to negotiate for resources. Short-term and long-term assessments show that CRA-WP programs positively impact participants.

## DISCUSSION

**Reed** (chat) commented that the NSF is working to address the missing millions from URGs in U.S. science and engineering.

**Cerf** (chat) requested information about CERP statistical charts. **Stent** (chat) relayed that all charts are available on CRA's webpage.

**Crivelli** (chat) asked if first-generation students are classified as an URG. **Stent** (chat) said the CRA-WP mission statement aims to address all forms of underrepresentation in computing. **Dwarkadas** commented that CRA should collect this demographic information.

Building on presentation remarks, **Russell** (chat) stated that CRA is proud to partner with many organizations, including AccessComputing at the University of Washington, the Institute for African American Mentoring in Computer Sciences, Center for Minorities and People with Disabilities in Information Technology, National Center for Women in Information Technology,

Computing Alliance of Hispanics Serving Institutions, the Students & Technology in Academia, Research & Services (STARS) Alliance, Expanding Computing Education Pathways and others.

**Crivelli** asked about outreach to high schools and community colleges, pointing out that community Colleges enroll many URGs and can be a pipeline entry point. Though CRA-WP is focused more on the pipeline beginning at the undergraduate level, **Dwarkadas** said the current strategic plan considers ways to engage high schools and community colleges as well as Minority Serving Institutions (MSI), including representation on the board. Scaling programs to high schools will require volunteers. Enlisting undergraduate help offers a solution: University of Rochester undergraduates that participated in CRA-WP programs recently began mentoring high school students of their own volition. CRA-WP can collate best practices for the community.

**Sarkar** asked how to use booming undergraduate enrollments to boost the percentage of URGs in CS. For example, the percentage of female undergraduates rose to 36% during the 1980s CS boom, but subsequently declined to 20%. Also, what guidance and resources can CRA provide to the national laboratories to increase engagement with MSIs and Historically Black Colleges and Universities (HBCUs) to augment URG participation in undergraduate research programs? **Russell** said DOE has access to CRA's graduating class directory which is populated with the applicants and participants of CRA's graduate cohort workshop programs. These individuals are interested in learning about employment options. National laboratories can engage directly with workshops by sending speakers and mentors. Forging personal connections with students is very valuable. **Stent** (chat) would like to see CRA's Distributed Research Experiences for Undergraduates (DREU) program, currently restricted to mentors at PhD granting institutions, opened to the national labs. **Dwarkadas** agreed with broadening CRA program scope to national labs and industry. DREU emphasizes research mentoring rather than providing just an internship. **Svore** (chat) advised that Harvey Mudd College, under the leadership of Maria Klawe, successfully increased women in the CS program after changing the curriculum and requirements. **Dwarkadas** (chat) added that Dr. Klawe initiated the Building, Recruiting And Inclusion for Diversity (BRAID) program at the University of Rochester. **Sarkar** (chat) said finding ways to scale successes more broadly at larger institutions is important.

Contrasting the urgency and resources used to address COVID or the ECP with those allocated to increasing participation of URGs, **Giles** asked when women will comprise 50% and minorities 12% of graduating CS researchers. The starting position is always one without much funding or the ability to make big changes. Maybe this could get done if someone were willing to spend \$100M. **Stent** said CERP and the Taulbee Survey show some progress, funding limits faster change. For example, the last CRA-WP graduate cohort had 400 participants, but there were 4,000 graduate women in the U.S. More individuals could be reached with more funding.

**THE LEGION PROGRAMMING SYSTEM**, Alex Aiken, SLAC National Accelerator Laboratory and Pat McCormick, Los Alamos National Laboratory

Legion is a task-based, data-centric programming model for parallel, accelerated, distributed (PAD) machines. Initiated in ~2011 by ASCR, Legion has received support from several DOE offices, programs and NNSA as well as the Defense Advanced Research Projects Agency (DARPA). Legion is now part of the ECP's Software Technology portfolio, has attracted industry investment from Facebook and NVIDIA and was a 2020 R&D 100 Award recipient.

Due to hardware changes, nearly every DOE machine is now a PAD machine. PAD machines have a complex memory hierarchy and significant memory capacity constraints, making task-based programming a better fit than previous models. Legion task graphs define

nodes as tasks. Edges are dependencies; they are subject to ordering constraints and inferred automatically as tasks are launched. This allows for parallel asynchronous execution, enabling graceful management of variable latencies. Importantly, task graphs are machine independent, with no commitment to the size of the machine, where tasks will execute or where data is placed.

During program execution, task graphs form behind a construction wavefront, followed subsequently by an analysis and mapping wavefront, data movement window, execution wavefront and destruction wavefront. Executed Legion tasks work on collections of data, and these collections can be partitioned into subcollections. Partitions are first class objects in Legion and can be hierarchical; multiple partitions of the same data can exist simultaneously, allowing composition. The application selects where tasks run and where collections are placed and then infers the needed communication. Legion mapping is computed dynamically and decouples correctness from performance.

Legion has successfully supported several applications at scale, including S3D, a combustion simulation; Soleil X, a multiphysics simulation of a solar collector heating nickel particles in a channel; FlexFlow, a deep neural network; and CANcer Distributed Learning Environment (CANDLE), a deep learning platform for accelerating cancer research.

Legate, developed by NVIDIA, is an open source PAD support for Python built in Legion that offers people with little to no HPC background access to supercomputers. Legate serves as a drop-in replacement for NumPy or Pandas libraries. Legate can automatically partition NumPy arrays and create tasks for NumPy operators, relying on Legion's partitioning support and task graph creation. Additional current and future work includes developing libraries that support Legion capabilities; supporting Legion interoperability with Fortran, C++, Python and MPI languages; continuing Legion deployment on exascale machines; writing CPU/ GPU portable kernels in Kokkos, OpenMP and Regent; and working on automating the mapping process.

## DISCUSSION

**Cerf** asked about partitions, If data is replicated, how does Legion track the effects of distinct processes executed with the same data? **Aiken** explained that Legion uses sequential semantics. The runtime system identifies task intersections and adds edges, forcing the correct order of execution and inserting communication. If two tasks reference overlapping data, the system must compute the intersection and move the updated data to the second task when the first task finishes. The runtime system has considerable caching and is able to learn overlaps quickly in repetitive computations, alleviating the need to recompute on every iteration. The runtime system also builds the task graph several iterations ahead of the application execution.

**Dongarra** complimented Legion for attracting users; it has been difficult to attract users to a similar system called Parallel Runtime Scheduling and Execution Controller (PaRSEC). What is the best way to educate the community about this new programming model? **Aiken** uses demonstrations to show that results can be obtained with less and easier code. Hopefully, transition will snowball with a critical mass of users. In the case of Legion, more documentation is needed, and education will be a crucial focus over the next couple of years.

**Hey** recommended using Legion to process curated experimental data to obtain measures of improvement on select applications for scientific ML benchmarking.

**Chapman** requested more information about the status of Regent. **Aiken** said Regent is a language for Legion and is an ongoing project. Its goal is to provide all Legion features to support the full programming model. Regent has its own compiler and some aspects of the programming model can be automated to generate a higher level of abstraction from the user.

Regent is currently being used at the Predictive Science Academic Alliance Program (PSAAP) III Center where the code has been stood up quickly and run on supercomputers.

**PERLMUTTER EARLY SCIENCE**, Jack Deslippe, Application Performance Lead, Lawrence Berkeley National Laboratory

The Perlmutter supercomputer (NERSC-9) follows the Cori system (NERSC-8) in transitioning users to energy efficient architectures. Perlmutter ranks sixth on the Green500 and is the most energy efficient of the Top500 Top 10 at 25.55 GF/W.

Perlmutter features an HPE Cray System with four times the capability of Cori's, GPU-accelerated (GPU/CPU) and CPU-only nodes, HPE Cray Slingshot high-performance network, and an All-Flash file system. Completed in the spring of 2021, the first of two installation phases stationed 1,536 GPU-accelerated nodes,. The second installation phase is scheduled for the winter of 2021 and will add 3,072 CPU-only nodes with an upgraded high speed network. This CPU partition will match or exceed performance of the entire Cori system.

The NERSC Exascale Applications Program (NESAP) is partnering with a total of 58 application development teams and vendors to port and optimize key applications. Provided resources include NERSC staff technical liaisons, performance postdocs, access to vendor application engineers and early access to hardware. Teams are assigned to one of two tiers, with the first tier having dedicated staff. First tier teams were chosen to represent different science and algorithm areas as well as different kinds of codes. Applications span all six DOE SC programs and cover NERSC's broad workload, including electronic structure; molecular dynamics; data; learning; particles and grids; and lattice quantum chromo dynamics (LQCD). Across these application dimensions, Perlmutter achieved a projected 20x system-wide throughput increase relative to performance on Edison (NERSC-7).

Perlmutter supports every GPU programming model, including Fortran/C/C++; CUDA; OpenMP 5.x; Kokkos/Raja; MPI; HIP; and DPC++/ SYCL. NERSC and NVIDIA partnered to develop a production OpenMP offload compiler released in April 2021.

Virtual hackathons have been effective in preparing applicaitons for new architectures and systems. NESAP organizes private quarterly DOE Hackathons for application teams engaging Cray and NVIDIA support as well as public GPU Hackathons, reaching NERSC teams around the world. NESAP also ensures that widely used community codes such as VASP, LAMMPS, WEST, CP2K and Quantum Espresso are available to load as system modules while collaborating to integrate vendor tools that lower the barrier to entry for non-experts. Lessons learned are documented on the NERSC hub and NESAP is holding trainings.

GPU optimizations, such as increasing parallelism and minimizing data movement, have led to enhanced application performance on Perlmutter for diverse science applications. The Dark Energy Spectroscopic Instrument (DESI) Spectral Extraction Python image processing code was ported in 2020. Optimization has resulted in a 25x improvement in per-node throughput using Perlmutter compared to the Edison baseline. ExaFEL will require real-time data analysis for decision making during experiments; in two years, NESAP developed a highly scalable application that contributed to significant runtime improvements, decreasing from 12.3 hours on Edison to two minutes on Perlmutter. Two materials science applications have been recognized as Gordon Bell Finalists: the Record Scale Molecular Dynamics with LAMMPS simulates 20B atoms, and a run achieved 11.24 PF on Perlmutter on 1024 nodes. The BerkeleyGW NESAP team achieved unprecedented simulation sizes with accuracy beyond density functional theory. The team is seeing speedups on the Perlmutter system. The Exabiome

NESAP team wrote the world's fastest GPU aligner for genetic sequences; there is still significant work to address sensitivity to communication and latency at large scale. The NESAP Generative Adversarial Networks (GANs) team achieved a 2.9x performance improvement over Cori on a ML scaling workflow.

## DISCUSSION

**Jack Wells** (chat) clarified that [gpuhackathons.org](https://gpuhackathons.org) is owned and under the supervision of OpenACC.org.

**Richard Gerber** (chat, NERSC) said NESAP experiences were used to create a user guide available at <https://docs.nersc.gov/performance/readiness/>

**Gregory White** requested more information about the user interface when implementing Perlmutter for real-time experimental data analysis. **Deslippe** explained that NERSC's super facility project has released an initial REpresentational State Transfer (REST) API that allows scientists to track system status as well as reserve and execute jobs. In future practice, when a researcher applies for beam time, they will concurrently apply for a certain amount of compute at a partner HPC facility and utilize the API to generate initial results during data collection.

## PUBLIC COMMENT

None.

**Reed** adjourned the meeting at 3:23 p.m.

Respectfully submitted October 29, 2021  
Holly Holt, PhD  
Science Writer, ORISE/ORAU