

Submission in Response to NSF CI 2030 Request for Information

DATE AND TIME: 2017-04-05 14:27:06

PAGE 1

REFERENCE NO: 251

This contribution was submitted to the National Science Foundation as part of the NSF CI 2030 planning activity through an NSF Request for Information, https://www.nsf.gov/publications/pub_summ.jsp?ods_key=nsf17031. Consideration of this contribution in NSF's planning process and any NSF-provided public accessibility of this document does not constitute approval of the content by NSF or the US Government. The opinions and views expressed herein are those of the author(s) and do not necessarily reflect those of the NSF or the US Government. The content of this submission is protected by the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License (<https://creativecommons.org/licenses/by-nc-nd/4.0/legalcode>).

Author Names & Affiliations

- Harvey Newman - California Institute of Technology
- Maria Spiropulu - California Institute of Technology

Contact Email Address (for NSF use only)

(Hidden)

Research Domain, discipline, and sub-discipline

High Energy Physics; Global network and systems architectures

Title of Submission

Exascale Ecosystems for 21st Century Global Science

Abstract (maximum ~200 words).

The largest data- and network-intensive programs supported by the NSF, DOE and partner agencies, including the Upgraded High Luminosity LHC program, the Large Synoptic Space Telescope (LSST) [2] and the Square Kilometer Array (SKA) astrophysics surveys, photon-based sciences, the Joint Genome Institute applications, the Earth System Grid continue to face unprecedented challenges: in global data distribution, processing, access and analysis, in the coordinated use of massive but still limited computing, storage and network resources, and in the coordinated operation and collaboration within global scientific enterprises each encompassing hundreds to thousands of scientists. A new overarching concept responding to these challenges is "consistent operations", where the scientific workflow management systems are deeply network aware and proactive, responding to moment-to-moment feedback on actual versus estimated task progress, state changes of the networks and end systems, and a holistic view of workflows with diverse characteristics and requirements able to serve many fields of science. This will enable the science programs to develop a new more efficient operational paradigm based on software-driven bandwidth allocation, load balancing, flow moderation and on-the-fly topology reconfiguration where needed. The result will be optimal use of the available network, computing and storage infrastructures while avoiding saturation and blocking of other network traffic.

Question 1 Research Challenge(s) (maximum ~1200 words): Describe current or emerging science or engineering research challenge(s), providing context in terms of recent research activities and standing questions in the field.

Context

Recently the ASCAC Subcommittee on Synergistic Challenges in Data-Intensive Science and Exascale Computing produced an in-depth

Submission in Response to NSF CI 2030 Request for Information

DATE AND TIME: 2017-04-05 14:27:06

PAGE 2

REFERENCE NO: 251

report and recommendations for data analytics integration with exascale computation systems, towards new kinds of workflows that will impact both data-intensive science and exascale computing. The report points out the need to urgently simplify the workflow for data-intensive science. HEP is one of the most data-intensive sciences and has already been critically contributing in architecting Exascale ecosystems.

The science workflow framework of the US HEP program presents an excellent ground to apply Recommendations 1 and 2 of the report . The vision of this framework includes injection of intelligent methods and learning algorithms into a new class of autonomous global systems to steer and optimize workflows among hundreds of petascale computing and storage sites as well as establishing an integral role of the US Leadership Computing Facilities (LCFs) and other major HPC facilities in this ecosystem and securely processing petabyte data slices from Exabyte data stores at the US HEP laboratories.

Implementation of this vision will generate novel data-intensive workflows that will accelerate discovery of the major science programs by developing new modes of network operations that redefine the state of the art in high throughput while remaining compatible with the tide of smaller flows exchanged over the world's research and education networks. The development of the new high throughput workflow and global system control and optimization methodologies, coupled to novel proactive, reactive and predictive software defined network system designs, and the proto-validation and verification of these — can be used directly, or as models to guide and inform solutions, for other data-intensive sciences.

Challenges in Data Intensive Sciences

We are entering a new era of exploration and discovery in many fields, from high energy physics and astrophysics to climate science, genomics, seismology and biomedical research, each with its own complex workflow requiring massive computing, data handling and networks. The continued cycle of breakthroughs in each of these fields depends crucially on our ability to extract the wealth of knowledge, whether subtle patterns, small perturbations or rare events, buried in massive datasets whose scale and complexity continue to grow (super-)exponentially with time.

In spite of technology advances, the largest data- and network-intensive programs supported by the DOE and partner agencies, including the Upgraded High Luminosity LHC [1] program, the Large Synoptic Space Telescope (LSST) [2] and the Square Kilometer Array (SKA) [3] astrophysics surveys, photon-based sciences, the Joint Genome Institute applications, the Earth System Grid and any other data-intensive emerging areas of growth , will continue to face unprecedented challenges: in global data distribution, processing, access and analysis, in the coordinated use of massive but still limited computing, storage and network resources, and in the coordinated operation and collaboration within global scientific enterprises each encompassing hundreds to thousands of scientists.

The most data intensive program is currently the LHC high energy physics program, with more than 400 petabytes under management at 160 sites in the US and around the world, growing to an estimated volume of one Exabyte by the time the data taking run now underway completes in 2018. The data storage and computational requirements needs are both projected to grow by another two orders of magnitude by the HL LHC era in the middle of the next decade. The CPU requirements of one of the LHC experiments (CMS) are expected to grow by two orders of magnitude between now and the HL LHC. The affordable CPU power obtainable within a fixed budget, including Moore's law and possible code improvements, is estimated to be an order of magnitude less.

This leads to a critical need to develop and deploy new data- and network-intensive operational modes making effective use of the US Leadership Computing Facilities starting now, continuing to expand and refine the system driving data exchange with the LCF's as they progress through the pre-exascale and exascale stages and beyond, while taking full advantage of the state of the art in high performance storage and networks across several technology generations.

It is also to be expected that the SKA and the other programs cited, and the rise of societal developments that boost genomics in support of precision medicine and sequencing of other species' genomes, may match or eclipse HEP's needs within the next decade. This makes the developments highlighted here all the more pressing.

Addressing the challenges faced by the science programs requires the development of a new class of autonomous, intelligent site-resident services that dynamically interact with network-resident services, and with the science programs' principal data distribution and management tools, to request or command network resources in support of high throughput petascale to exascale workflows, using:

- (1) smart middleware to interface to SDN-orchestrated data flows over network paths with guaranteed bandwidth all the way to a set of high performance end-host data transfer nodes (DTNs),
- (2) protocol agnostic SDN-based QoS and traffic shaping services at the site egress that will provide stable, predictable data transfer rates, and auto-configuration of data transfer nodes, and

Submission in Response to NSF CI 2030 Request for Information

DATE AND TIME: 2017-04-05 14:27:06

PAGE 3

REFERENCE NO: 251

(3) host- and site agent systems coupled to machine learning methods

The development of such systems will exploit the synergy between deeply programmable software-defined agile and adaptive network infrastructures that are emerging as multi-service multi-domain network “operating systems” interconnecting next generation Science DMZs, and the systems developed by the data intensive science programs harnessing global workflow, scheduling and data management systems.

Specific R&D areas where development is required include:

(1) deep site orchestration among virtualized clusters, storage subsystems and subnets to successfully co-schedule CPU, storage and network resources; (2) science-program designed site architectures, operational modes, and priorities adjudicated across multiple network domains and among multiple virtual organizations; (3) seamlessly extending end-to-end operation across both extra-site and intra-site boundaries through the use of next generation Science DMZs; (4) funneling massive sets of streams to DTNs at the site edge hosting petascale buffer pools configured for flows of 100 Gbps and up, exploiting state of the art data transfers where possible; and (5) unsupervised and supervised machine learning and modeling methods to drive the optimization of end-to-end workflow involving terabyte to multi-petabyte datasets.

Question 2 Cyberinfrastructure Needed to Address the Research Challenge(s) (maximum ~1200 words): Describe any limitations or absence of existing cyberinfrastructure, and/or specific technical advancements in cyberinfrastructure (e.g. advanced computing, data infrastructure, software infrastructure, applications, networking, cybersecurity), that must be addressed to accomplish the identified research challenge(s).

The overall solution calls for network systems with embedded integration of i) novel modes of steering, use and reuse of data produced and consumed at multiple locations, ii) novel modes of propagating information on data availability, cost of delivery as a function of re-computation in real-time, and iii) interactions among user-groups and sites. To implement such system scientists and engineers will take advantage of the remarkable emerging synergy between deeply programmable, agile software-defined network (SDN) infrastructures -- which are evolving towards multi-service multi-domain network operating systems interconnecting science teams across regional, national and global distances-- and worldwide secure distributed systems riding on high capacity (but still-passive) networks developed by the data intensive research programs.

A new overarching concept is one of “consistent operations”, where the scientific workflow management systems are deeply network aware, reactive and proactive, responding to moment-to-moment feedback on actual versus estimated task progress, state changes of the networks and end systems, and a holistic view of workflows with diverse characteristics and requirements able to serve many fields of science. This will enable the major science programs to develop a new more efficient operational paradigm based on software-driven bandwidth allocation, load balancing, flow moderation and on-the-fly topology reconfiguration where needed. The result will be full and optimal use of the available network, computing and storage infrastructures while avoiding saturation and blocking of other network traffic. The systems to be developed should aim to serve multiple domain sciences. One very fertile area for development and progressive large scale field testing is the HEP LHC program, which is now on the cusp of its second three-year run, anticipated to yield a new round of groundbreaking discoveries, as well as a new level of global data and network intensity. In parallel different but equally challenging real-time workflows in diverse fields need to be anticipated and developed for science areas including bioinformatics, computational astrophysics, radio astronomy, and oceanic and atmospheric sciences.

Cross-cutting teams working at the intersection of their domain science, computational science and technology/engineering will have a crucial role in implementing this strategy.

Leadership Computing, Storage and Network (CSN) Ecosystems for Next Generation Data Intensive Science

The use of Leadership Computing Facilities (LCFs), and folding them into the data- and network intensive ecosystems of the HL LHC thus presents us with a great opportunity in terms of providing the needed CPU resources for the HEP science program, with longer term benefits for a wide range of other science programs with similar needs.

Taking HEP as the guiding use-case, the Argonne, Oak Ridge LCFs and other major HPC facilities such as NERSC, can develop the next generation envisioned systems as described below from the science Virtual Organization point of view. The key steps are to:

- Recast HEP’s generation, reconstruction and simulation codes, case by case, to adapt to the emerging HPC architectures, addressing

issues of memory, dataflow versus CPU etc.

- Identify and match the units of work in HEP's workflow to the specific HPC resources or sub-facilities well-adapted to the task (after the code recasting step)
- Build dynamic and adaptive "just in time" systems that respond rapidly (on the required timescale) to offered resources as they occur.
- Develop algorithms that effectively co-schedule CPU, memory, storage, IO port, local and wide area network resources
- Develop an appropriate security infrastructure, and corresponding system architectures in hardware and software, that meet the security needs of the LCF
- Apply machine learning to optimize the workflow of the HEP experiments, using self-organizing system methods which are well-adapted to such problems while also taking the special parameters, conditions, and restrictions of LCFs into account as part of the workflow
- Exploit the intense ongoing development of virtualized computing systems, networks and services in the research community and in industry: in the data center, campus and wide area network space aimed at coherent distributed system operations (including software defined networking, network function virtualization, and service chaining, along with emerging higher level concepts)

From the LCF and HPC facility vantage point the following mirrored response is envisioned:

- Identify and match the units of work in HEP's workflow to the specific HPC resources or sub-facilities well-adapted to the task (after the recasting step)
- Build dynamic and adaptive "just in time" systems that respond rapidly (on the required timescale) to offered demands as they occur; including client-side/server-side coordination for a consistent outcome
- Develop algorithms that effectively co-schedule CPU, memory, storage, IO port, local and wide area network resources; with the necessary coordination as above
- Develop an appropriate security infrastructure, and corresponding system architectures in hardware and software, that meet the security needs of the LCF. For the LCFs this means adopting a new mode of ongoing service to a major client in quasi-real time, in a way that can be adapted to meet the LCF's requirements.
- Apply machine learning coupled to game theory and system modeling, to optimizing the workflow of the HEP experiments, using self-organizing system methods which are well-adapted to such problems; while also taking the special parameters, conditions, and restrictions of LCFs into account as part of the workflow.
- Exploit the intense ongoing developments of virtualization of computing systems and services in the research community and in industry: in the case of the LCFs, the recent developments of "site orchestration" of virtualized resources, and even newer concepts of secure ways to bridge the site edge, such next generation Science DMZs [4] or similar edge-bridging methods are relevant.

LCF-Edge Data Intensive System Operational Model

A new class of LCF-Edge Data Intensive Systems is emerging as a promising direction for the envisioned future global systems architectures. The use of secure systems at the site perimeter means that security (both human and AI) and countermeasures where needed can be focused on a limited number of subsystems and entities (proxies), so that the manpower burden may be acceptable. The operational concept for HEP calls for the data to be brought into the edge systems in chunks (a petabyte per chunk for example), far enough in advance so that the data is always on standby and ready when the corresponding jobs are scheduled to start. Multiple chunks for different stages of the overall workflow are foreseen, with each having a definite provenance and attributes (such as the ratio of CPU to I/O requirements) so that clusters of chunks can be matched to an appropriately configured HPC subsystem. At a later stage, one can also foresee dynamic restructuring of the HPC resources, especially if they are virtualized in logical sectors. Considering the parameters in this problem yields interesting consequences. As of today, a 1 petabyte chunk would occupy a 100 Gbps link if used to 100% capacity for a full 24 hour day. Given the 400 petabytes currently stored by the LHC experiments and the fact that approximately 850 petabytes flowed over the networks in and out of the US in the past year, the 1 petabyte chunks foreseen each will represent a relatively small "data transaction" compared to the whole task at hand, and so one would like to transport many chunks to and from the LCF. A typical near-future configuration would thus preferably include several 100 Gbps links, migrating to several 400 Gbps links within approximately 5 years and to several 1 Tbps links by the startup of the High Luminosity LHC a decade from now, depending on the demand evolution and the cost evolution during this period.

Question 3 Other considerations (maximum ~1200 words, optional): Any other relevant aspects, such as organization, process, learning and workforce development, access, and sustainability, that need to be addressed; or any other issues that NSF should consider.

Outlook

Submission in Response to NSF CI 2030 Request for Information

DATE AND TIME: 2017-04-05 14:27:06

PAGE 5

REFERENCE NO: 251

We described significant and novel parts of the vision for network integrated systems for global data intensive science using HEP's LHC as a working initial use-case. The architecture described can be used or be adapted to many domain sciences from astronomy to genomics and beyond. The implementation of the vision in the form of R&D and scalable demo systems is ongoing in the context of collaborative initiatives with National Laboratories and University partnerships (ANL, FNAL, LBNL, Caltech, Stanford, Yale etc). With further R&D and future production and deployment of such systems we shall be able to respond to the big-data challenges we are faced with in almost all areas of scientific exploration today and for the next decade and beyond.

Consent Statement

- "I hereby agree to give the National Science Foundation (NSF) the right to use this information for the purposes stated above and to display it on a publically available website, consistent with the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License (<https://creativecommons.org/licenses/by-nc-nd/4.0/legalcode>)."
-