

Reference ID: 11226726011_Honavar

Reference ID: 11226726011_Honavar

Submission Date and Time: 12/16/2019 2:41:14 PM

This contribution was submitted to the National Science Foundation in response to a Request for Information, <https://www.nsf.gov/pubs/2020/nsf20015/nsf20015.jsp>. 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>).

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 publicly available website, consistent with the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License (<https://creativecommons.org/licenses/by-nc-nd/4.0/legalcode>)."

Consent answer: I consent to NSF's use and display of the submitted information.

Author Names & Affiliations

Submitting author: Vasant Honavar - Pennsylvania State University

Additional authors: Jenni Evans, Pennsylvania State University; Guido Cervone, Pennsylvania State University; Ryan Gilmore, Pennsylvania State University; Wayne Figurelle, Pennsylvania State University

Contact Email Address (for NSF use only): (hidden)

Research domain(s), discipline(s)/sub-discipline(s)

Artificial Intelligence; Machine Learning; Meteorology; Geoscience; Bioinformatics; Materials Informatics; Data Science; Computational Sciences.

Title of Response

Harnessing the power and potential of data, computation, and artificial intelligence to advance science

Abstract

How should a large public institution such as Penn State take advantage of advanced cyberinfrastructure to assist in tackling some of the large challenges that impact our society and the world to enable

important scientific advances? In this response, we identify the emerging interdisciplinary research initiatives that rely on cyberinfrastructure, outline the significant challenges of building a responsive cyberinfrastructure ecosystem, and provide suggestions on the strategies and directions to address the sustainability of cyberinfrastructure for cyber-enabled scientific discoveries in the 21st century.

Question 1 (maximum 400 words) – Data-Intensive Research Question(s) and Challenge(s). Describe current or emerging data-intensive/data-driven S&E research challenge(s), providing context in terms of recent research activities and standing questions in the field. NSF is particularly interested in cross-disciplinary challenges that will drive requirements for cross-disciplinary and disciplinary-agnostic data-related CI.

Today's cutting-edge research tackles critical and complex problems facing our society, utilizing highly interdisciplinary data- and computation-intensive approaches. These approaches rely on big data, and on the cyber-ecosystem that can support model building and analyses of the data. Through cyber-enabled scientific advances, we might be able to: Alleviate the societal burden of chronic illnesses through advances in personalized monitoring and prediction of health risks and intervention outcomes; Improve learning and teaching through data driven advances in the learning sciences; Harness the power and potential of data to advance precision farming; Harness data and computation to enhance the effectiveness of responses to man-made or natural disasters; Accelerate materials synthesis and discovery by coupling modeling, experimentation, and prediction; Advance socio-technical systems in ways needed to harness data sciences and artificial intelligence to the benefit of all stakeholders. Many challenges stand in the way of these advances. One critical problem is that our ability to collect massive volumes and different varieties of complex data has outstripped our capacity to analyze them. Another is that we increasingly need simulations that can span vast ranges of scale. In explaining and contextualizing these challenges, we draw upon cases faced by researchers connected to Penn State's Institute for Computational and Data Sciences (ICDS), a university-wide research organization supporting interdisciplinary cyber-enabled science and engineering research. ICDS accomplishes its integrative mission both by growing an interdisciplinary community of cyber-enabled researchers and cyber-professionals, and by providing this community with a state-of-the-art high performance cloud (the ICDS Advanced CyberInfrastructure, ICDS-ACI).

Data Deluge: The Volume of Data Outstrips our Capacity to Analyze It

Advances in instrumentation and methodology in many fields have led to a data-rich, but knowledge-poor, research environment. Researchers now have unprecedented access to data, but challenges remain in integrating, analyzing, and understanding this data. In meteorology and Earth Science in general, the advances in our ability to observe and model the Earth system through the use of remote sensing, numerical models and social media has led to the generation of geographically distributed, asynchronous, dynamic, and diverse datasets. The rate at which these data are generated often greatly exceeds our ability to analyze them. Likewise, new research in the social and behavioral sciences increasingly depends on understanding the social networks that are characterizing a new era of human communication through social media. Social media data can provide critical insights into human behavior and organization, including political view and policy evolution, socio-economic development, population changes, and many more. This data is voluminous and accumulates by the second. The

emerging field of “precision health in context” relies on understanding the individual’s psychological, cognitive, and physiological states through predictive and causal modeling of health risks and health outcomes. This requires integrative analyses of clinical, genetic, environmental, socio-demographic, behavioral, and other types of data. Challenges span data access and use policy compliant infrastructure for analyses of sensitive data, machine learning, causal inference, and data integration methodologies for coping with complex data, among others. The life sciences are being transformed by various advances in method and instrumentation. High-throughput laboratory techniques have yielded vast amounts of genetic information, and have made possible system-wide measurements of the entire transcriptome, large portions of the proteome, interactome, metabolome, and diseasome. Advances in imaging, including the development of cryo-electron microscopes, have made possible the collection of diverse cell images at, which can be combined with measurement of the expression levels of thousands of genes to directly see the responses of individual cells and tissues under different perturbations. These advances, while promising, are also extremely data-intensive. Materials Discovery is an area where recent advances in machine learning, together with advances in computing and high throughput measurements of materials properties offer the potential for new data-driven approaches to predicting properties of materials based on their descriptors. This approach has a high potential as a complement to direct experimentation and to simulations that rely on explicit solutions to fundamental equations. When suitable data exists or can be generated, these data-driven methods can be useful to predict material properties that are too difficult to measure or compute using traditional methods, due to the cost, time, or effort involved. The descriptors may be of many types and scales, depending on the application domain and needs. Predictions may be interpolative or extrapolative, allowing the design of entirely new materials. These are only a few examples of the many fields that face the challenge of data deluge. Advanced analytical tools, algorithms, and computational models are imperative for ingesting and analyzing all of the available data. Expanding the Scalability of Simulations Computational sciences increasingly involve predictive modeling in the presence of uncertainty. From forecasting weather and natural disasters to disease outbreaks and financial risk, predictions are hampered by incomplete definition of the initial conditions – even in the face of massive advances in observational data streams – and our ability to represent the complete range of scales necessary to minimize uncertainty in predictions. For instance, an accurate simulation of a hurricane needs to span processes from molecular to planetary scales. Interactions at the atomic level can have significant impacts on a hurricane’s growth, as the interactions of water molecules and phase changes of water drive the convective clouds that govern hurricane intensity. But hurricane models must also be able to simulate the large-scale interactions of a hurricane with other weather systems, as these govern the storm’s motion and structure. Many other fields face similar multi-scale modelling challenges. Predicting the dynamics of cancer might require simulations that combine cellular-scale with organism-scale models. Predicting how new materials and chemical compounds interact might require combining models at different scales of time, from femtoseconds to seconds. Of particular interest are methods for data assimilation, for refining the predictions of simulation models by incorporating real-time data. Despite the promise of widely-scalable simulations, modeling frameworks that can handle vastly different scales and which can deliver accurate predictions in uncertain conditions are computationally intensive and demand personnel with a high level of expertise. This increasing level of demand but lack

of corresponding supply of human expertise may present significant challenges to many fields that rely on advanced cyberinfrastructure. Society faces challenges that are fundamental to the health and wellbeing of individuals, communities, nations and the planet itself. Increasingly rapid advances in observing and sensing technologies raise research questions that would have been beyond imagination in earlier eras, but developing the capability to process the data created by these technologies remains difficult. Integrating the data into multi-scale models that accurately predict uncertain events is also a major challenge. Solving these problems will likely lead to major transformation in a variety of fields and have crucial societal impacts.

Question 2 (maximum 600 words) – Data-Oriented CI Needed to Address the Research Question(s) and Challenge(s). Considering the end-to-end scientific data-to-discovery (workflow) challenges, describe any limitations or absence of existing data-related CI capabilities and services, and/or specific technical and capacity advancements needed in data-related and other CI (e.g., advanced computing, data services, software infrastructure, applications, networking, cybersecurity) that must be addressed to accomplish the research question(s) and challenge(s) identified in Question 1. If possible, please also consider the required end-to-end structural, functional and performance characteristics for such CI services and capabilities. For instance, how can they respond to high levels of data heterogeneity, data integration and interoperability? To what degree can/should they be cross-disciplinary and domain-agnostic? What is required to promote ease of data discovery, publishing and access and delivery?

The research challenges described above point to various cyberinfrastructure needs. While specific cyberinfrastructure requirements differ across fields, some needs are common, including: Standard tools, approaches, and frameworks for handling big data across different disciplines; Development of new tools that incorporate insights from emerging statistical, analysis and machine learning approaches, among others; Solutions for end-to-end scientific workflows for ingesting and analyzing massive amounts of data to developing data driven models, generating hypotheses, planning experiments, executing experiments, interpreting results, archiving the products of data; and Maintaining the security of data, especially as cyber-enabled researchers increasingly rely on clouds. **Standard Tools, Approaches, and Frameworks for Handling Big Data** Advancing the capability and capacity for realizing the promise and potential of data to advance science requires: computational abstractions of the relevant domains coupled with computational methods and tools for their analysis, synthesis, simulation, visualization, sharing, and integration; cognitive tools that leverage and extend human intellect, and partner with scientists on all aspects of science; agile and trustworthy data cyber-infrastructure that connect, manage a variety of instruments, multiple interrelated data types and associated metadata, data representations, processes, protocols and workflows; and enforce applicable security and data access and use policies; and organizational and social structures and processes for collaborative and coordinated activity across disciplinary and institutional boundaries. Many research fields rely on gathering, analyzing, interpreting, and integrating massive volumes of data. Across fields, there are widely different approaches to coordination, orchestration, and execution of these tasks. Many special-purpose tools exist for solving specific data problems, but their proliferation reflects a lack of best practices for developing, standardizing, and executing models. Tools providing a

framework for scalable and flexible execution of a large and collective set of tasks for high-level data mining and visualization of big data, applicable to a variety of data types, would be valuable to many fields. The lack of these standard tools and practices hampers interdisciplinarity. For example, in Earth Science piecemeal solutions have to be replaced by cyberinfrastructure solutions cognizant of the end-to-end application requirements and supporting different science problems that have similar computational characteristics. Several tools have been developed to assist Earth scientists to run their workloads on HPC infrastructures specifically with the requirements of solving a specific problem as the first-order design consideration. These solutions are often tied to a single problem domain and are not applicable to different domains that have similar computational requirements. Similarly, biologists and neuroscientists are often interested in similar issues related to neurodegenerative disease, each collecting their own data that could be mutually informative (e.g., a biology lab with large amount of neuroanatomical data from mice and rats and a psychology lab with huge human brain imaging data). Effective use of this data in science requires major breakthroughs in data and computing infrastructure as well as innovative machine learning methods and tools. These tools would have to handle the different types of data gathered by biologists and neuroscientists; only then would scientists in these fields arrive at cross-cutting insights that will lead to fundamental discoveries of biological/neural processes. While standard tools are important, future cyberinfrastructure advances must also be flexible enough to meet the needs of different kinds of research projects. Many research projects can benefit from exploiting HPC systems, but given the current design of HPC resources towards single monolithic jobs, not all tasks make optimal use of existing HPC systems. There is a need for robust and scalable approaches toward large-scale data analytics that can cover the wide-spectrum of application scenarios. We need cyberinfrastructure solutions that support different science problems with varying computational characteristics. These tools must be accessible, and support for them must be sustainable, especially in today's climate of tight funding. Software tools and data streams would need to be disseminated freely to the research community and frequently updated. The community would need to provide consultation on what hardware and software advances would need to be supported within a standard framework.

End-to-End Data Management and Scientific Workflow Solutions

Cyberinfrastructure could help to relieve the problems caused by the data deluge by streamlining many tasks in the data lifecycle—the stages through which data passes, from initial ingestion to dissemination and archiving. Each step in this lifecycle has its demands. Some raw data sources must be downloaded manually and passed through collection points before being transferred to the storage space where it can be analyzed and manipulated. Sharing datasets with collaborators can often cause delay and replication, and is not always done securely. Disseminating data requires tagging it with metadata—an often complex manual process. Penn State has begun to design a comprehensive end-to-end data management system that automates many tasks in the data lifecycle. The proposed system would allow for automatic downloading of data from designated sources, as well as automatic data tagging and classification to reflect security/privacy access and other use restrictions (HIPAA, etc.). Ingested data would be routed by the system to a secure shared research space, where it could be accessed and manipulated by teams of collaborators. Once results have been obtained, the system would automate the generation of metadata and push the data to outside portals, where it could be searched for and obtained by other researchers and the public. This system would help to alleviate strain on intra-

institution networks, especially at large institutions like Penn State. Currently, the limited capacity for transmission of terabytes (or more) of data from one site to another, even within the Penn State system, obstructs research progress and discourages collaboration. An integrated data flow framework such as that described here would reduce unnecessary repetitive data transfers and go some way to addressing this problem, but would not completely obviate it. Related to concerns about the data lifecycle, and to satisfying Federal and other data management plans, is the question of data archiving costs and searchability. Once a grant expires, there is no funding source to support storage of research products created from the work funded by that. Further, codified requirements on metadata and data formats are needed to ensure that data remains searchable in long term. This includes a need to refresh software for these data formats so they do not timeout. Compliance in the Cloud Researchers at many institutions are increasingly turning to cloud providers for high-performance computing services. How these cloud environments affect compliance with data security and privacy policies, however, is as yet unclear. Community guidance is needed on institutional responsibilities for data governance when individual researchers use the cloud and when cyber-institutes take research into the cloud. Addressing these cyberinfrastructure issues—finding standard tools and frameworks that aid collaboration while addressing the demands of different types of computational research; streamlining end-to-end data management; and ensuring cloud compliance—would allow researchers better to take advantage of massive quantities of available data and perform more accurate and complex simulations.

Question 3 (maximum 300 words) – Other considerations. Please discuss any other relevant aspects, such as organization, processes, learning and workforce development, access and sustainability, that need to be addressed; or any other issues more generally that NSF should consider.

Besides issues that can be addressed through cyberinfrastructure advances, there are other areas of action that NSF might consider, including: Support for training to develop skills needed to use cyberinfrastructure effectively; Data governance mechanisms for overseeing ethical use of data, including “restricted data”; A deliberate framework for coordinating collaborative community endeavors; Increasing NSF’s agility for planning new cyberinfrastructure-related initiatives; Leveraging the promise and potential of artificial intelligence to advance science as well as operation of computational and data infrastructure; Providing long-term support for cyber-research experts Researchers and students must be informed of the latest innovative cyberinfrastructure technologies and tools, and they must receive the training that can help them use these tools effectively. This education and training must be made available at all levels—from first degree earners to returning professionals. Stronger coordination between institutes such as ICDS (and academic institutions more generally) and industry would help keep researchers on the cutting edge of technology, and allow students to prepare for promising academic or industry careers. Currently, using restricted data at large academic institutions presents an administrative challenge. At Penn State, a researcher wishing to use data with regulatory and/or contractual restrictions on data use and/or dissemination must work with at least ten different offices, and often more, to remain compliant. Due to the dispersed services of these offices, we do not have a straightforward, routine way to recognize that we are receiving restricted data, categorize the data, implement the technical requirements of the restrictions, and

ensure that compliance is met throughout the research project. A coordinating body that could work to streamline the use of restricted data at large institutions would help address this problem. NSF funds a number of community endeavors for sharing disparate data across related communities (e.g. EarthCube). To make these activities more effective, a deliberate framework for coordinating these activities and sharing information is needed. What progress is being made across the community—and which areas should progress—must be made clear. Otherwise, we end up with unintended outcomes that may satisfy nobody, or a community comprised only of a very small core team (which defeats the original purpose of these endeavors). NSF’s community initiatives should include computational and data infrastructure specialists, not just domain experts. While domain experts may have strong programming and data management expertise, their skills are often limited to their research areas; cyber-ecosystem experts can bring valuable new approaches to research problems. NSF currently supports a number of activities, e.g., through regional Big Data Innovation Hubs, that offer excellent platforms for fostering community engagement and collaboration across organizational and disciplinary boundaries. Other initiatives include the midscale infrastructure projects that bridge the gap between institutional and national scale infrastructure. We strongly encourage continued investments in these and other similar initiatives. Given the increasing importance of Artificial Intelligence technologies in accelerating scientific discovery and enhancing cyberinfrastructure, we suggest additional activities and investments as they relate to (i) AI infrastructure at scale (with particular emphasis on cyberinfrastructure needed to scale up current AI methods and enable the next generation of AI techniques) (ii) AI for accelerating science (with particular emphasis on aspects of science that are currently human labor intensive and hence constitute bottlenecks, AI tools that augment and extend human intellect and abilities, and cyberinfrastructure to optimally support human-human, human-machine and machine-machine collaboration in science) and (iii) AI for self-managing adaptive cyberinfrastructure (including AI techniques that could enhance the ability of cyberinfrastructure to anticipate and adapt to changing workloads, hardware and software failures, etc.). Computational and data scientists are becoming an increasingly critical part of today’s complex cyber-ecosystem. We need a model for employing these individuals (likely PhD level, many with non-CS PhDs) that does not lead to “holes” in their funding. Institutions can’t afford to provide permanent salaries for all research areas, so we need a model that allows for a combination of agency and university funding that is flexible across individual grants. An institution-level grant to support the growth of research in new fields might be a way to provide consistent funding for these experts. This support would allow us to build a more diverse cyber-community, as we will be able to accommodate life transitions (parenthood, gaps in workforce, etc.). Finally, the cyber-ecosystem should also seriously consider the education of new generations of scientists who are CI-oriented graduate students today. This steady supply of technical expertise to the research fields will be critical for the sustainability of the future of cyberinfrastructure and the advances of frontier sciences and engineering.

-- End Submission --