

NESE: The North East Storage Exchange

James Cuff¹, Scott Yockel¹, John Goodhue², Saul Youssef³, Rajiv Shridhar⁴, Ralph Zottola⁵, Chris Hill⁶, Glenn Bresnahan²
Harvard University¹, MGHPC², Boston University³, Northeastern University⁴, University of Massachusetts⁵, MIT⁶

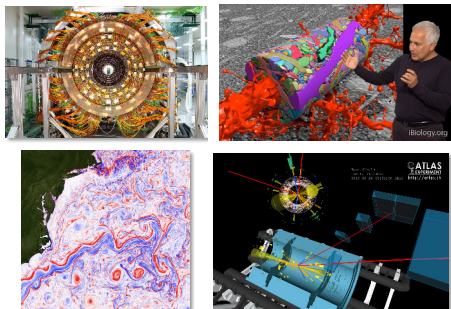
ABSTRACT

Research progress is increasingly dependent upon the available capacity of storage to flexibly exploit large volumes of digital information. The North East Storage Exchange (NESE) project creates a next-generation storage infrastructure specifically targeted at enabling new levels of collaborative research in projects regularly involving petabytes of information. This storage exchange will integrate with a computational and network infrastructure that links Harvard University, Boston University, the Massachusetts Institute of Technology (MIT), Northeastern University and the University of Massachusetts system. This project contributes to building a national data infrastructure to support advanced research in such priority topics as health care, epidemiology, physics, and earth science, among others.

NESE will provide a high capacity, highly networked, secure, cost effective, scalable, and accessible data store that lowers barriers to research, collaboration, and information sharing within and beyond the participating multi-university community. Some examples of NESE projects that will be early users of NESE include one of the four US Tier 2 centers that store and process ATLAS data from the Large Hadron Collider; the Center for Brain Science at Harvard University, which is generating 300 million micron-resolution images to map the billion neurons and synapses that make up a cubic millimeter of the human brain; and MIT collaborations with NASA and DARPA in next generation global ocean modeling and monitoring systems. NESE addresses several critical infrastructural challenges: the creation of a sustainable multi-institutional resource; advancement of methods for data retention, management, and access to sensitive research data; implementation of controls that simplify protection of sensitive data; and building a sustainable, collaborative operating infrastructure to support future research.

This award by the Advanced Cyberinfrastructure Division is jointly supported by the NSF Directorate for Biological Sciences (Division of Biological Infrastructure), and by NSF's Understanding the Brain and BRAIN initiative activities.

SCIENCE



NEWS

November 1st 2016: NESE AWARD:
https://nsf.gov/awardsearch/showAward?AWD_ID=1640831

November 28th 2016: HARVARD GAZETTE:
<http://news.harvard.edu/gazette/story/2016/11/for-bigger-data-more-storage/>

NESE GUIDING PRINCIPLES



CORE VALUES

- SECURE:** As consent based research data sets become standard practice (in particular within Health Science), security models are having to catch up with the Data Use Agreements required. From dbGap, to CMS/Medicare, our researchers manage significantly more human and health care subject data than ever before. For societal change to occur, and to produce better outcomes for patients through research and basic science, we need significantly more performant and secure data storage systems. We can't do this alone, or in isolation.
- ARCHIVE:** Scientists and researchers discuss data retention, archive and provenance on what seems to be a daily basis. We have multiple solutions to this challenge, but no unified overarching system that we can point to as a "standard". As funding agencies require more sophisticated "Data Management Plans", our research faculty are left with a bewildering array of options, each more confusing than the last. This has to stop.
- COST:** Storage is expensive. Many hundreds of millions of dollars are spent annually attempting to solve the challenge of reliable, available storage for science. The potential for economies of scale by collecting and coordinating resources here in what could well be argued as the most research data intensive part of the nation is vast. We are capable, and have proven by MGHPC that we can do more with less. Much more.
- CAPACITY:** We have heard this for many years now - there is quite simply an explosion of data in science, it is not being managed, and this proposal points to both technology and process to be able to manage unlimited capacity requirements.
- BANDWIDTH:** Science data requirements demand high performance storage. It is not sufficient to simply provide large capacity, as data access patterns vary dramatically across disciplines, and each NSF directorate. Fortunately, "object stores" (the technology we will deploy as part of NESE), are inherently designed to scale out for both speed and capacity.

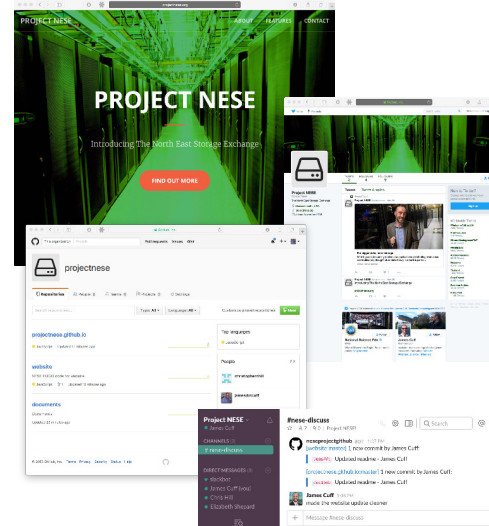
PRIMARY ROLES

Organization	Primary Role
Harvard	Build, install and operate CEPH object store hardware, software, and monitoring. DTN cluster configuration and deployment, Globus endpoint management, devops and systems engineering, cluster high availability with required network configuration and security of science DMZ. Assistance to Harvard research groups testing or adopting NESE storage.
Boston University	Planning, network configuration, testing and migration of NET2 storage to NESE. Integration of NESE storage into NET2 operations. File system interfaces to NESE. Assistance for BU research groups testing or adopting NESE storage. Federation with external CEPH clusters.
MIT	iRods overlay to NESE object store. Demonstration and evaluation using 4PiB heterogeneous ocean data. Support of iRods + NESE application to separately funded combined altimetry and ocean color research. Development of general cookbooks illustrating use of iRods and NESE for open data sharing and discovery science activities.
Northeastern	POSI file and block storage presentations of NESE object store. Evaluation of cost metering and allocation to researchers / projects. Help Northeastern researchers evaluate and use NESE object store.
University of Massachusetts	Policy, Standards, Cybersecurity and Security Operations Center. The resource will join a team that provides 24x7 monitoring, alerting and escalation; ensuring incidents are detected, investigated, communicated, and reported. Assistance to UMass research groups testing or adopting the NESE object store.
MGHPC	Federated authentication and access control; SDN access, operation of the data center that houses the hardware; planning for long term sustainability; physical security for sensitive data.

MANAGEMENT PLAN

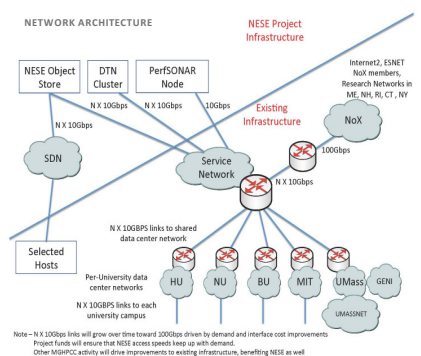


NESE PUBLIC COMMUNICATIONS

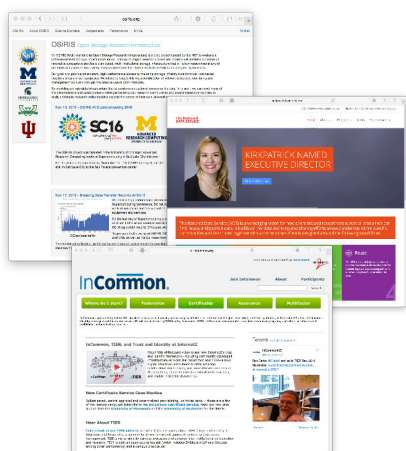


FIVE YEAR PLAN

Year	Primary Goal	Subgoals
1	• Purchase and install first tranche of equipment • Commence gateway work • Design network • Add BU ATLAS as first science case	• Verify initial build • Migrate test data • Appoint workgroup leadership
2	• Add HU and MIT science • Build a user community advisory group • Continue work on ATLAS integration • UMass Security Operations Center Support	• Design security model • Confirm year 1 network and system operation • Continue gateway work
3	• Add UMASS and NEU science • Freeze year one and year two builds, harden presentation layer • Support for data subject to HIPAA regulations	• Based on input from the user advisory group, build a first pass model for forecasting measurement, and allocation of costs and resources
4	• Commence "show backs" • Build appropriate fiscal model	• Build federation model
5	• Finalize the operating model for long term sustainability, including operations and capital refresh	• Add physical equipment from ATLAS • Complete external federation model



BUILDING ON PRIOR ART



WE ARE HIRING!

Business Title: Cyberinfrastructure Storage Engineer, Research Computing
School/Unit: Faculty of Arts and Sciences (FAS)
Location: USA - MA - Cambridge
Job Function: Information Technology
Time Status: Full-time
Schedule: Monday - Friday, 9am-5pm
Department: Science Division/Research Computing
Salary Grade: US\$

Duties & Responsibilities: The Cyberinfrastructure Storage Engineer will play a key part in the North East Storage Exchange (NESE) project, directly supporting a critical example of the significantly increased growth of Research Computing at Harvard. This position is entirely NSF funded (NSF-ACI-1640831) through October 2021. The NESE project is a joint venture between the five founding member universities and the Massachusetts Green High Performance Computing Center (MGHPC). The Cyberinfrastructure Storage Engineer will be responsible for creating and maintaining a multi-terabyte object storage platform and partnering with the other institutions regarding access, network and security aspects of this project. The NESE project is under the direction of the PI, James Cuff, FAS Assistant Dean and Distinguished Engineer for Research Computing, and will report to the Senior Team Lead for High Performance Computing.

Basic Requirements: A bachelor's degree in computational science or a related field with a background in engineering storage solutions. Minimum 7 years Unix/Linux system administration experience with at least 3 years experience as a storage solution architect.

Additional Qualifications: The ideal candidate will have established a background installing / monitoring / maintaining: (1) large scale (multiple PB) clustered storage (not single pointed NAS/SAN); (2) object storage, especially Ceph; (3) multi-tiered or hierarchical file systems; (4) high-throughput data transfer nodes, especially GridFTP; (5) PerSonAR server. Experience with secure data (HIPAA/FCRM) protocols and data encryption at rest for multi-tenant resources is a plus. Familiarity with software defined networks, federated access between authentication providers, Science DMZ, and other access related concepts are key to interacting with the partner institutions. Authentication and identity management across federated systems (certificate management, OAuth tokens, Shibboleth etc.) will form a key part of the final service delivery, the candidate will have shown previous implementations at scale.

Additional Information: When applying for this position please submit your resume and cover letter in our preferred format as one combined document (resume followed by cover letter). All formal offers will be made by FAS Human Resources

<http://news.harvard.edu/gazette/story/2016/11/for-bigger-data-more-storage/>