

Research Drivers and Capabilities

Open Storage Network Concept Paper

December 7, 2020

Christine R. Kirkpatrick¹, Kevin Coakley¹, Melissa Cragin¹, James Glasgow², John Goodhue³

Affiliations: 1. San Diego Supercomputer Center, University of California San Diego, 2. National Center for Supercomputing Applications, University of Illinois Urbana Champaign, 3. Massachusetts Green High Performance Computing Center

The Open Storage Network is funded by the National Science Foundation under grants #1747552, 1747493, 1747507, 1747490, 1747483 and by Schmidt Futures

<http://www.openstoragenetwork.org/>

Executive Summary	2
Introduction	3
Case Studies	3
Terra Fusion	3
Sage	4
Integrated Hurricane Data	5
Conclusion	6

Executive Summary

The Open Storage Network (OSN) supports science and scholarly research that requires data storage and transfer at scale, by simplifying and accelerating access to data that is in active use by ongoing research projects. Deployment of the OSN is a response to the increasing importance of storage as the third component of national cyberinfrastructure, complementing investments in computing and networks. While other uses may emerge over time, the OSN is intended initially to serve two principal needs: (1) facilitate smooth flow of large data sets between data and computing resources such as instruments, synthetic data projects, campus data centers, national supercomputing centers, and cloud providers; and (2) make it easy to expose long tail data sets to the entire scientific community.

This is the first of a series of concept papers outlining the function and role of the OSN in the research infrastructure landscape. Here we address three cases: The Terra Fusion project, which has produced a massive dataset by fusing data from multiple instruments; the SAGE project is working to streamline data flowing from Internet of Things (IoT) devices; and a team of water and hazards researchers who have developed a new access point for hurricane data that will facilitate new research.

Each case study illustrates the use of the OSN to solve a specific problem. For Terra Fusion, the data is available from NASA as well as AWS Glacier. OSN solves the need for ready and speedy access for specific sample datasets. Whereas access to NASA datastores is restricted and somewhat slow, and Glacier access has a days-long delay and may incur extra costs, OSN delivers select Terra Fusion products quickly while leveraging existing infrastructure and investments. The use of OSN for Sage allows for computing on IoT devices, rather than copying all raw data to the cloud or manual retrieval. The availability of the Integrated Hurricane Data allows for on-demand reuse of models and integrated data as hurricanes strike. The OSN's object store provides ease of use and end-user access when every moment counts for keeping communities safe with data-driven advice. The highlighted projects demonstrate the research drivers and capabilities that are pushing the need for a distributed nationwide storage fabric. Single institutions do not have the resources needed to meet the storage needs for the research drivers of projects that are national scale and the commercial cloud cannot be relied upon to meet the needs of these research drivers due to costs and barriers to access.

Introduction

Scientific use cases are driving the need for distributed storage fabric. The OSN, established initially through funding from the Schmidt Foundation, just completed two years of research and development with support from the National Science Foundation. The third year of the project includes a four-part webinar series and companion concept papers to explore the impact of the project, place it in context of cyberinfrastructure related research challenges, and to disseminate findings for others to build upon. The webinar series kicked off in October 2020, the OSN with presentations from three projects and their scientific use cases¹.

- Donald Petravick (National Center for Supercomputing Applications (NCSA) at the University of Illinois) presented how the Terra Fusion project uses the OSN to share multi-petabyte data sets.
- Wolfgang Gerlach (University of Chicago and Argonne National Laboratory) presented how the Sage project uses the OSN with IoT and AI applications.
- Christina Bandaragoda (University of Washington) and Chris Lenhardt (Renaissance Computing Institute (RENCI) at UNC Chapel Hill) presented how they use the OSN to make Integrated Hurricane Data easily accessible to researchers and public health officials in order to make critical decisions before, during and after a hurricane.

A recording of the presentations can be viewed on the OSN website.² A closer examination of each case illustrates how the OSN is solving the storage needs for a wide array of research projects.

Case Studies

From Earth Science, to smart and connected cities, to hazards research, storage is a critical enabler for research teams, and communities.

Terra Fusion

Research Driver

The Terra project is the flagship of NASA's Earth Observing System (EOS).³ The 20 years of longevity for the Terra project has made it a critical resource for studying the Earth's climate and climate change. The Terra dataset has become one of NASA's most popular datasets due to interest from the scientific, government, commercial, and educational sectors. The selection of [the] five instruments for Terra were based, in part, on the potential for obtaining greater quality of information for Earth Science through the synergistic use (data fusion) of the five instruments compared to individual instruments alone.³ The goal of the Terra Fusion project has been to generate "fused" files from the 5 instruments and make the fused data available to the public. The Terra Fusion dataset is 2.4 PB and growing.

¹ <https://www.openstoragenetwork.org/abstracted-use-cases/>

² <https://www.openstoragenetwork.org/seminar-series/oct-22-2020/>

³ <https://digirolamo.web.illinois.edu/projects/terra-fusion/>

OSN Utilization

The OSN hosts ~150 TB of Level 1 Terra Basic Fusion files, about 6% of the total mission volume. This subset is known as a “Sampler” and provides the community with mission scale data for paths 11-14 (i.e., Bermuda), 108 & 233 (i.e., Greenland), 123-126 (i.e., all of China), and 143-147 (i.e., all of India). Users use a command line interface to download portions of the sampler. This is documented on the Terra Fusion project website².

How OSN Furthers the Terra Fusion Project’s Scientific Goals

The complete Terra dataset is not available to the public due to its size; instead users are provided access to the data “Sampler”, and there is the ability to serve out additional of these subsets. NASA has storage facilities in their Distributed Active Archive Centers⁴, however the separate instrument data is housed at different locations and sharing is limited. It took multiple years over high speed research networks to transfer the entire dataset from all five instruments for fusing. Once this fused dataset was ready, there were attempts to share the Terra Fusion data utilizing tape archives; while cost effective, this constrains the pace of the scientific process. The Terra Fusion project manages a copy of the whole dataset on NASA’s AWS S3 Glacier service, available to researchers upon request. In order to make access to the AWS S3 Glacier copy of the data cost effective and to work around S3’s restore limits, the Terra Fusion project had to implement a virtual Data Carousel,⁵ which mimics a physical tape carousel. This effectively constrains data calls due to restoration limits. AWS S3 Glacier imposes a limit of 35 restore requests per day; for Terra fusion data, this means that only 96 Level 1 files can be restored. Moving to AWS S3 Glacier didn’t remove the delays to the scientific process, users have to wait two weeks between the data access request and data availability. AWS S3 charges \$0.05/GB for data egress to the Internet. Hardware for three OSN pods could be purchased for the cost to download the whole Terra Fusion dataset 3.3 times. The OSN increases the time to discovery for researchers at a lower cost than traditional IT infrastructures, including commercial cloud.

Sage

Research Driver

The Sage project is a novel cyberinfrastructure created to exploit dramatic improvements in AI technology with the goal to build a continent-spanning network of smart sensors.⁶ Sage moves the advanced machine learning (ML) algorithms to “the edge,” to run on “Internet of Things” (IoT) devices, in order to streamline processing and analytics data flowing from the IoT devices themselves. Currently, data retrieval requires that technicians travel to those IoT devices to gather the data from onboard storage, or, that users access only the fraction of the data that can be uploaded to a cloud server due to bandwidth constraints. Running the ML algorithms on the IoT devices helps with the challenge of analyzing the large volume of data recorded by high-fidelity sensors that are located in the environment and part of a vast IoT network.

⁴ <https://earthdata.nasa.gov/eosdis/daacs>

⁵ <http://hdl.handle.net/2142/107186>

⁶ <http://sagecontinuum.org/about/>

OSN Utilization

The Sage project's object storage data repository is currently a work in progress. While development is mainly happening on developer workstations, the development data is being mirrored to the OSN to verify everything will work when in production. The Sage project's object storage data repository hosted by the OSN will store images, videos, sound files, LIDAR data and multispectral images from edge sensors across the United States. Sage data will be primarily housed on the SDSC OSN pod, cited to be near the SDSC (OpenStack) Cloud for just-in-time processing at the edge. Users will be given access to the data via a Sage-provided interface.

How OSN Furthers Scientific Goals

The Sage project has been on the forefront of using AI and Deep Learning (DL) algorithms to push the computing on to the IoT devices. Using edge computing and transmitting the derived results allows IoT devices to provide greater bandwidth efficiency, privacy/security, resiliency, energy efficiency and lower latency than saving to the device or continuously streaming the sensor output. Efficient and accurate AI and DL algorithms are required in order to accomplish computing on the IoT devices, which can be low powered and have slow Internet access. Bringing together many different types of IoT sensors on devices with different geographic scales, the Sage project works with projects like AoT⁷ (video, environmental sensors), HPWREN⁸ (video), and NEON⁹ (environmental, soil, surface water, and groundwater sensors) which operate at neighborhood, regional and continental scales. A goal of the Sage project's collaboration with the OSN is to provide domain experts and data scientists access to a large and diverse set of IoT sensor data in order to improve the efficiency and accuracy of their AI and DL models. An important role of the OSN is to be able to help projects in their early stages by providing no or low cost storage, direct access to the infrastructure engineers to solve bottlenecks, and to share the knowledge and experience of other projects that have been successful in delivering data at scale.

Integrated Hurricane Data

Research Driver

Due to climate change the negative impacts of hurricanes and floods is increasing. The National Oceanic and Atmospheric Administration (NOAA) and other government sources publish hurricane observations, flood maps, storm track forecasts, National Water Model (NWM) forecasts, de-identified drinking water quality sampling and more. Researchers need access to these datasets in order to project and educate themselves against water hazards. This project [creates] a synthesized data and software system to advance our understanding of how digital and physical infrastructure information reduce the impact of disasters using an integrated collaborative platform for ongoing research.¹⁰ The Integrated Hurricane Data project works with

⁷ <https://arrayofthings.github.io>

⁸ <http://hpwren.ucsd.edu>

⁹ <https://www.neonscience.org>

¹⁰ https://nsf.gov/awardsearch/showAward?AWD_ID=1902537

HydroShare to provide access to data, research findings and data models used by researchers and environmental engineers.

OSN Utilization

The OSN is hosting NOAA's NWM¹¹ data from Hurricane Matthew (8 TB) and Hurricane Harvey (20 TB). The data is stored in three types of NetCDF¹² files: 1km gridded (land surface variables and forcing), 250m gridded (ponded water depth and depth to soil saturation) and Point-type (stream routing and reservoir variables). The data is accessible on the web at https://matthew_nwm.renc.osn.xsede.org (Hurricane Matthew) and https://harvey_nwm_zip.renc.osn.xsede.org (Hurricane Harvey).

How OSN Furthers Scientific Goals

When there is a hurricane it is important to give pertinent data to scientists, engineers and public health officials as soon as possible. Questions of public safety, like identifying areas where flooding could occur and the potential impacts on quality of drinking water, need to be assessed during and immediately after a hurricane. Data collected from previous hurricanes can be used to model what will happen in a crisis. NOAA generates and publishes high volumes of NWM data that can be used to create these models; unfortunately, these data are not easily accessible as NOAA doesn't store this data on centralized servers and then they are archived offline after only a couple of weeks. The Integrated Hurricane Data project's aim is to gather a persistent dataset of integrated NWM data for significant events, like hurricanes, in a central location that is always available. Once the Integrated Hurricane Data project acquires the data, they need to make sure they are always available to the public for when the next hurricane occurs. The OSN excels at making datasets available to the public, through its S3 compatible API, via unauthenticated https, and through integrated third party services. The OSN removes data access barriers to allow researchers to access critical data and focus on the important issues around their research and public safety.

Conclusion

The Terra Fusion, Sage, and Integrated Hurricane Data projects demonstrate a few of the research drivers and capabilities that underscore the need for a distributed nationwide storage fabric. Single institutions do not have the resources necessary to meet the storage needs for the research drivers of projects that are national scale and the commercial cloud cannot be relied upon to meet the needs of these research drivers due to costs and barriers to access. From the beginning, the OSN has worked directly with these and other use cases in order to tailor the service around research drivers and capabilities instead of institutional or business drivers. The OSN has positioned itself as the large scale, distributed storage fabric to solve the storage needs of a wide array of research drivers.

¹¹ <https://water.noaa.gov/about/nwm>

¹² <https://www.unidata.ucar.edu/software/netcdf/>