

National and International Trends in Research Storage at Scale

Open Storage Network Concept Paper

March 12, 2021

Christine R. Kirkpatrick¹, Kevin Coakley¹, Melissa Cragin¹, James Glasgow², Kenton McHenry², Uwe Jandt³, Lene Krøl Andersen⁴, Ilya Baldin⁵ Anita Nikolich⁶, and Alex Szalay⁷

Affiliations: 1. San Diego Supercomputer Center, University of California San Diego, 2. National Center for Supercomputing Applications, University of Illinois Urbana Champaign, 3. Deutsches Elektronen-Synchrotron, 4. The Danish eInfrastructure Cooperation, 5. Renaissance Computing Institute, University of North Carolina at Chapel Hill, 6. University of Illinois Urbana Champaign, Johns Hopkins University

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
State of the Art	3
EOSC-Nordic	3
Helmholtz Federated IT Services (HIFIS)	5
FABRIC	6
Key Takeaways	8
Conclusion	8

Executive Summary

The Open Storage Network (OSN) is a pilot designed to provide a national storage substrate for access and use of active scientific and scholarly data. This distributed service leverages existing investments in network infrastructure to support high speed transfer among institutions and computational facilities. Over the last five years the U.S. National Science Foundation (NSF) has funded more than 200 high-speed connections to the Internet-2 backbone operating at 10-100Gbps speeds. The goal of this project is to develop a prototype module for a high performance distributed storage system that extends the usability of the existing high-speed interconnects. OSN pods each provide 1 PB of active storage at low cost and with little maintenance and leading edge automation; these are linked together to form a larger storage pool that supports access, data movement, replication, and local caching near computational resources to better support high performance. At full scale the OSN could serve Petabyte-scale research, which is generally beyond the capacity of individual academic institutions and independent research organizations. The scientific and scholarly communities are facing a major challenge dealing with the increasing amount of research data emerging from projects produced at all scales-- from large facilities to small research labs.

This is the second in a series of concept papers outlining the function and role of the OSN in the research infrastructure landscape. Three national and international projects that serve research storage at scale are examined: EOSC-Nordic, which brings together research institutions, research infrastructure providers and policy makers, in the Nordic and Baltic region; HIFIS, which manages private cloud services for the Helmholtz Association in Germany; and FABRIC which is deploying its own network of compute and storage resources across the US, Asia and Europe.

Each project illustrates the state of the art in trends for research storage at scale. EOSC-Nordic established policies and use cases for enabling sharing across national borders. HIFIS integrated authentication and authorization infrastructure is integrated to present distributed storage as an easy to use service. FABRIC created a tiered based solution for achieving optimal performance of finite resources, so active data storage is near compute services. These projects provide several key takeaways related to the importance of the sociotechnical including policy, international participation, and the need to design and tune services to specific use cases and researcher needs. These exemplars demonstrate that the current trend is to present storage to users as a service instead of a file system. There is still great need in the US; providing research storage at scale has not been solved. Several organizations are working on national and international efforts to remove issues that researchers encounter when working with active data at scale.

Introduction

Many national and international projects have been established to address the challenge of storing active research data at scale. 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. This is the second concept paper resulting from the OSN seminar series. Three presentations from projects employing state of the art technologies and sociotechnical approaches were featured in the second webinar, “National and International Trends in Research Storage at Scale.”

- Lene Krøl Andersen, EOSC-Nordic project manager, from the Danish eInfrastructure Cooperation (DeiC), presented how EOSC-Nordic brought together dozens of institutions across the Nordic and Baltic regions to utilize capabilities including shared storage resources in the European Open Science Cloud (EOSC).
- Uwe Jandt (Deutsches Elektronen-Synchrotron) discussed how the HIFIS initiative federated the IT services among the 18 Helmholtz Association research centers to simplify access to storage services for their researchers.
- Ilya Baldin (Renaissance Computing Institute (RENCI) at UNC Chapel Hill) & Anita Nikolich (University of Illinois Urbana-Champaign) presented the challenges of managing multiple tiers of data storage at scale across geographic and network borders.

A recording of the presentations can be viewed on the OSN website.¹

State of the Art

EOSC-Nordic

Research Profile

The EOSC-Nordic project², funded by the European Commission, aims to facilitate the coordination of EOSC relevant initiatives within the Nordic and Baltic countries and exploit synergies to achieve greater harmonisation at policy and service provisioning across these countries, in compliance with EOSC agreed standards and practices. EOSC-Nordic brings together a strong consortium of 24 complementary partners including e-Infrastructure providers, research performing organisations and expert networks, with national mandates and experience with regards to the provision of research data services, and a unique capacity to realise the outcomes of the EOSC design as outlined by the EOSC Implementation Roadmap.

¹<https://www.openstoragenetwork.org/seminar-series/nov-12-2020-national-and-international-trends-in-research-storage-at-scale/>

² <https://eosc-nordic.eu/>

The EOSC-Nordic project is coordinated by The Nordic e-Infrastructure Collaboration (NeIC)³, which facilitates development and operation of high-quality e-infrastructure solutions in areas of joint Nordic interest. By working together to realise EOSC at the regional level, the countries can provide “Nordic Added Value” (*Nordisk Nytte*) a key concept for the Nordic Council of Ministers and Nordic cooperation at large and a way to achieve more together than separately. EOSC-Nordic is currently working with over 200 participants in 10 countries in the Nordic and Baltic regions. The complete EOSC-Nordic workplan, consortium, deliverables etc can be accessed via the public EOSC-Nordic Wiki page, hosted by NeIC⁴.

Data Sharing Capabilities

The EOSC-Nordic project takes part in designing an Analysis and Data Cloud Infrastructure, leveraging EOSC resources, which ensures that researchers have easy and secure access to working with large amounts of data seamlessly under uniform conditions across the countries.

Users can search for the list of data services provided by EOSC through the European central gateway; i.e the EOSC Portal⁵ or through the regional gateway; via. ESOC-Nordic Knowledge Hub⁶, co-hosted by NeIC, increasing regional sustainability.

The EOSC-Nordic Knowledge Hub is the Nordic and Baltic regional EOSC enabling gateway! The Knowledge Hub contributes to a wider adoption of processes and practices and optimal use of the resources available. In its present and initial phase it also serves as a regional and cross border support mechanism helping to identify and remove barriers in the whole uptake of EOSC serving several stakeholder groups.

Storage Insights

EOSC-Nordic contributes to environmental and socially important areas, mainly through the activities of its demonstrators by widening the user base of the thematic services operating in the areas of biodiversity, Human Language Technology, climate science and precision cardiology.

EOSC-Nordic has developed four thematic demonstrations, via ten use cases⁷, helping researchers enhance the usage of the data being stored via a European federated data infrastructure; i.e. EOSC.

1. **Discovery and Re-Use of Research Data** through tools to harvest metadata and update the EOSC metadata catalogue⁸.
2. **Analysis and Post-Processing** service provided by community specific portals hosted at large scale computing facilitates across EOSC-Nordic.

³ <https://neic.no/>

⁴ <https://wiki.neic.no/wiki/EOSC-Nordic>

⁵ <https://eosc-portal.eu/>

⁶ <https://www.eosc-nordic.eu/knowledge-hub/>

⁷ <https://www.eosc-nordic.eu/demonstrating-eosc-nordic/>

⁸ <https://www.eosc-nordic.eu/open-science-will-help-us-better-understand-the-vikings/>

3. **Data Management Sharing⁹ and Archiving** service that facilitates sharing of data across EOSC-Nordic's distributed environment.
4. **Sensitive Data and Orchestration** service to allow data analysis of sensitive data without moving the data away from the data custodian.

These four thematic demonstrators show EOSC-Nordic's commitment to meet the challenges of delivering data services to meet the needs of the researchers.

Key Insights for the US Research Community

EOSC-Nordic demonstrates how a coordinating organization (i.e. NeIC) can facilitate access for existing services to researchers across national borders. While the US research community does not cross national borders, the boundaries between states and institutions can have many of the same legal and policy challenges. Legal and policy requirements often prevent data and services from being shared, not technical limitations. Identifying and sharing political challenges in doing Open Science in practice, is facilitated via workshops interlinking technical with non-technical (alias policy) experts. Thus creating understanding in where and how to bring technical solutions to life across country borders and policies. Open Science landscape analyses and practical barriers are communicated via a regional research infrastructure portal like the EOSC-Nordic Knowledge Hub. This understanding between technical and policy level is critical to breaking down non-technical barriers to national research efforts. EOSC-Nordic provides other important sociotechnical best practices including regular and frequent check-ins with institutional members. This gives all partners a chance to elevate issues and to ensure all partners have a voice in the initiative. It is a key component of modeling continuous improvement and feedback from end users as well. The successful cross-border and cross-institution collaboration model adopted in EOSC-Nordic, builds upon the many years of expertise in the Nordic infrastructure collaboration, facilitated and continuously being developed via the organisation NeIC.

Helmholtz Federated IT Services (HIFIS)

Research Profile

The German Helmholtz Association is an alliance of 18 research centers in Germany with an annual budget of approx. €5 billion. The platform “Helmholtz Federated IT Services” (HIFIS)¹⁰ has been established to provide common access to IT resources within Helmholtz¹¹, as well as training and support for professional and sustainable scientific software development. The Helmholtz Cloud Services offered by HIFIS include services for large data transfer, high performance computing, as well as documentation and collaboration tools of all kinds – the need for the latter is in very high demand.

⁹ <https://www.eosc-nordic.eu/wp-5-wants-to-improve-the-entire-data-management-process/>

¹⁰ <https://hifis.net>

¹¹ <https://www.helmholtz.de/en/research/information-data-science/helmholtz-incubator/>

Data Sharing Capabilities

Although the main focus of HIFIS are the common, secure and easy access to cloud services as well as support for professional scientific software development, there is obvious demand for reliable and versatile data storage and data transfer solutions. In the context of HIFIS, the dCache project has been integrated into the Helmholtz Cloud, allowing access to storage by users from all Helmholtz centres – and their collaboration partners – either directly or via connected services, e.g., HPC. Other HIFIS data sharing services include for example Sync&Share services and permanently identifiable data storage for scientific data.

Storage Insights

The dCache identity management allows to seamlessly integrate into the authentication and authorization infrastructure (AAI) of the Helmholtz Cloud using OpenID-Connect. On top of this, dCache features advanced authorization delegation, thus facilitating to share access rights between research groups at definable levels, without any need to share secrets, creating additional accounts, etc. dCache provides a user-transparent workflow to enable buffering and pre-fetching of data transfers, thus minimizing performance losses due to latency. It also supports integration with large data transfer services such as CERN's File Transfer Service and Globus.

Key Insights for the US Research Community

From the start of HIFIS and all HIFIS Services, the requirements of all scientific communities assigned to Helmholtz - including scientific platforms like the Helmholtz Imaging Platform and the Helmholtz Artificial Intelligence Cooperation Unit - have been surveyed extensively, allowing to select and shape services according to their needs.

Strict compatibility with existing supranational IT frameworks like EGI and the European Open Science Cloud (EOSC) is maintained, e.g. by following the AARC framework for authentication and authorization. All HIFIS activities are being accompanied by international experts in our scientific advisory board, thus incorporating the experience and knowledge of the best experts in the field.

FABRIC

Research Profile

FABRIC¹² is a unique, national research infrastructure to enable cutting-edge and exploratory research at-scale in networking, cybersecurity, distributed computing and storage systems, machine learning, and science applications. When completed FABRIC will become a widely

¹² I. Baldin and A. Nikolich and J.Griffioen and I.Monga and KC Wang and T. Lehman and P. Ruth, "FABRIC: A National-Scale Programmable Experimental Network Infrastructure," in IEEE Internet Computing, vol. 23, no. 6, pp. 38-47, 1 Nov.-Dec. 2019, doi: 10.1109/MIC.2019.2958545.

distributed instrument both for computer science research and for the many science domains that want to explore faster and more capable distributed computational and data infrastructures. FABRIC Across Borders (FAB) extends the FABRIC network to add nodes in Asia and Europe for expanded scientific impact. Both FABRIC and FAB are funded by the NSF via Mid-scale Research Infrastructure-1 (Mid-Scale RI-1) and International Research and education Network Connections (IRNC) programs and will be available for use by academic and industry researchers spanning multiple domains including internet architecture and protocols, network measurements, artificial Intelligence/machine learning, cyber-security as well cyberinfrastructure (CI) research for many scientific disciplines, like astronomy, high-energy physics, weather/climate and many others.

Data Sharing Capabilities

Experimental Data is stored by FABRIC for a very short time and FABRIC encourages users to share their data more broadly. FABRIC will facilitate data sharing and is currently exploring mechanisms such as a portal (i.e., the FABRIC web page) for users to access data sets. The goal is to serve as a conduit for data sharing while not storing data long-term. FABRIC will work with other data sharing projects such as SEARCCH¹³ (Sharing Expertise and Artifacts for Re-use through Cybersecurity Community Hub), the Research SOC¹⁴ and others to enable wide sharing of data.

Storage Insights

FABRIC faces the same storage challenge as many other federally-funded projects, namely that the project does not fund a large amount of storage and certainly not enough storage to keep every type of data needed or wanted by the community for even a modest amount of time. Since information collected about both the instrument and by user experiments will quickly exceed the amount of storage the project provides, FABRIC developed a tiered storage and backup model to determine what elements to keep and for what period. The tiers are based on three elements: required retention times, data characteristics (shared vs. internal) and data function (operations vs. experimenter data). FABRIC's approach is to develop a network of campus, regional, and potentially commercial cloud partners who will work with them on long-term data storage challenges.

Key Insights for the US Research Community

FABRIC demonstrates that constructing scientific CI that spans the globe will always have storage challenges which the national research community must creatively solve until there is a national solution. Storage for large projects remains largely unfunded, leaving projects to partner with each other and campuses to leverage available storage. Not all data can be treated the same. Projects should take a community-based yet aggressive approach to determining what data is needed versus desirable, and develop a tiered model based on these factors to

¹³ <https://www.flux.utah.edu/project/search>

¹⁴ <https://researchsoc.iu.edu/>

help alleviate the storage challenges. FABRIC is a model for how the US Research Community can deal with multiple data dimensions, types of data - data collected from experiments, data collected from CI, data imported into the testbed to support experiments, and data-as-a-service from a testbed.

Key Takeaways

A brief survey of the state of the art in distributed storage for researchers yielded several key points that the Open Storage Network should adopt or maintain:

- **Policy** plays a key role in architecting service flow (EOSC-Nordic).
- It is not enough to provide storage, a service must have robust **authentication and authorization** (HIFIS).
- For services that align well with stakeholder needs, one must begin with those use cases and keep tuning the service so it remains **researcher focused** (HIFIS).
- Services should be **researcher, capabilities facing and not technical facing**, offering data sharing and interfaces that assist with this rather than outreach about storage technical specifications (EOSC-Nordic, HIFIS).
- Planning for **future integrations begins with adhering to standards** and 'supranational IT frameworks' where available, e.g. EOSC compatibility (HIFIS, EOSC-Nordic).
- **International advisory boards** are a key way to stay abreast of developing services and initiatives that can inform an initiative (HIFIS).
- As long as there is a lack of US national research infrastructure, projects will be forced to create their own solutions and use institutional resources. This has the side effect of pushing effort to create **tiered data structures** to mitigate shortfalls in storage (FABRIC).
- None of the state of the art examined provide **archival storage**.

Conclusion

EOSC-Nordic, HIFIS and FABRIC show that active data storage has not been solved and researchers cannot rely solely on the commercial cloud for their storage needs. Funding agencies have made large investments in computing and networking resources, but data storage infrastructure has not seen the same consistent, long term investment, leaving individual projects to create their own solutions and leverage institutional resources where they exist. As well, it would be more efficient to keep active data near these prior NSF investments. A variety of approaches is needed to solve the active data needs of researchers. These solutions can include locally managed institutional storage, nationally managed distributed storage and commercial cloud storage. EOSC-Nordic, HIFIS and FABRIC demonstrate the trend to hide the storage layer behind services to simplify the storage interactions for researchers. Services, like iRods, Clowder and NextCloud, provide simple web interfaces for users to access and organize their data. Providing this storage service layer makes it easier for researchers to access their data and allows storage providers to choose from a wide variety of open source and commercial

storage options. The state of the art efforts underway at EOSC-Nordic, HIFIS and FABRIC, show that OSN is well positioned within the national and international trends in research storage at scale by providing storage for active data that is distributed nationally and by providing a platform for researcher-friendly storage services to run on.