

# The Open Storage Network as a Storage Solution for NSF CC\* and Other Program Solicitations

Dr. Christopher S. Simmons  
Deputy Director, MGHPC  
[csim@mghpcc.org](mailto:csim@mghpcc.org)

# US Research Cyberinfrastructure Today

## Computation

*Local, Regional and  
National Resources*

*Standardized*

## Networking

*Over 200 universities  
with 40/100Gb  
Connectivity*

*Standardized*

## Storage

*Largely Balkanized*

*Many standards to  
choose from; typically  
tied to a single system*

# CC\* Program Requirements

- The program supports open-source platforms and solutions
- Software license fees are not allowed
- Budget request for professional services are allowed
- At least 20% of the disk/storage space on the proposed storage system should be made available as part of the chosen federated data sharing fabric
- OSN provides all of this as a service via ACCESS
- Pod owners use our coldfront portal to request projects and buckets
- OSN has been the storage solution for 6 successful CC\* Data Storage Awards

# The Open Storage Network

National resource for sharing open scientific data

With Distributed Infrastructure and Governance

With Simple and Flexible Access and Management methods

# Distributed Infrastructure



# FAIR Data Guiding Principles

- Findable
  - Globally unique persistent identifier, metadata registered or indexed in a searchable resource
- Accessible
  - Standard protocols; authentication and authorization when needed; persistent metadata
- Interoperable
  - Metadata vocabularies and formats that support exploration and use of available data sets
- Reusable
  - Well described attributes, provenance, terms of use, and domain-relevant community standards

OSN Focus



Wilkinson, M. D. et al.

The FAIR Guiding Principles for scientific data management and stewardship.

Sci. Data 3:160018 doi: 10.1038/sdata.2016.18 (2016).

# Simple and Flexible Access and Management

- Access from anywhere
  - RESTful API (S3 riding on https)
  - Think of OSN as a publicly-routed data lake
- Authentication and Authorization services
  - Federated identity protocols and services
- Security and data integrity
  - Open-source Ceph object store
- High performance
  - Research networks and commodity servers

# Key Ideas

- Every data set is a collection of objects
- Every object is accessible from anywhere
  - site-name.osn.xsede.org/bucket name/data set name/object name
  - e.g. [https://mghp.osn.xsede.org/osndemo1/AHM18\\_OSN\\_Poster.pdf](https://mghp.osn.xsede.org/osndemo1/AHM18_OSN_Poster.pdf)
- Physical storage is distributed across multiple sites
  - Located in science DMZs for fast access preferably on Internet2



# Access and Curation

- Open Access Data Sets
  - Readable by anyone
  - Writeable by anyone with access to the RW API key pair
- Protected Access Data Sets
  - Readable by anyone with the RO API key pair
  - Writeable by anyone with the RW API key pair

# Beyond “just” sharing data

- OSN was originally designed as “just” a low-friction platform for sharing data across organization boundaries
- However, projects and organizations are starting to use OSN as part of their larger cyberinfrastructure plans
- We now allow Universities, Non-Profit Organizations, and Government Labs to purchase their own OSN “appliance”

# OSN summary

- You own the infrastructure
- Distributed OSN Implementer's Team does the sysadmin/operations work
- Operations transparent; Open weekly meetings and GitHub-hosted playbooks
- Cost-recovery model run by non-profit organizations using commodity hardware and open-source software

# Workflows enabled by OSN

- Check in / Check out locally or on campus and national systems
  - Rclone, cyberduck (S3 native)
  - Globus (via the AWS S3 Connector); Also adds support for ACLs
- Compute locally while accessing storage remotely
  - Python - Zarr, boto3, fsspec + kerchunk
  - Julia - AWS.jl and AWSS3.jl
  - R - aws.s3 package
- Mountable file system
  - S3FS, rclone, iRods, juicefs
- Use as backend storage w/ a Research Data Management Platform
  - Dataverse, InvenioRDM, Clowder, iRods

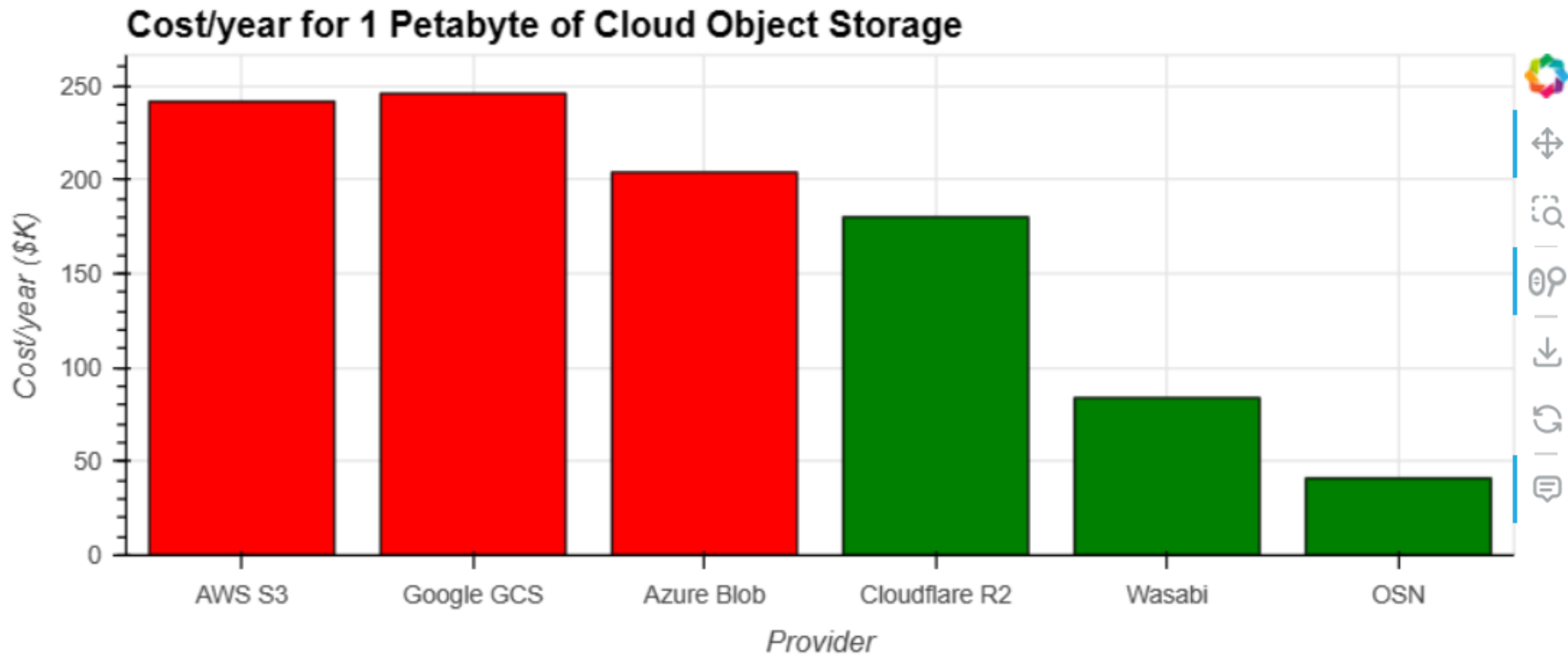
# Workflows enabled by OSN for AI/ML Pipelines

- AI/ML and LLMs require traceability and reproducibility
- Data “just” sitting on a POSIX file system and consumed doesn’t capture history and metadata
- Multiple Data Version Control solutions help tackle this problem
  - Git LFS – versioning data directly with Git but stored on OSN
  - DVC – versioning of model and data based on GitOps principles
  - LakeFS – GitOps for data pipelines with support for multiple APIs including Apache Airflow, Hive, Spark and native R and Python with a total of 22 integrations

# OSN Cost Model

- You purchase an OSN Gen2 Pod
  - ~\$100K today but expected to go down over the next ~6 months
  - 3 x 1 U Dell Services
  - 1 x 4U Seagate Exos AP with 106 SAS drives
  - 5 years of support
- You send us \$10K per year for 5 years
  - Current cost model is \$10K per organization NOT per Pod
  - This is subject to change for next CC\* cycle or future purchase
- You rack and stack in your data center and boot each node once with a provided USB drive or host with us for an additional fee
- We do the rest!

# Independent OSN Cost Analysis by USGS



Questions?