# Globus Research Data Management: Introduction and Service Overview

Steve Tuecke
tuecke@uchicago.edu

Vas Vasiliadis
vas@uchicago.edu

Presentation material available at

globus.org/events/sc15
# bit.ly/globus-sc15

# Thank you to our sponsors!

U.S. DEPARTMENT OF **ENERGY**

NSF

ALFRED P. SLOAN FOUNDATION 1934

NATIONAL INSTITUTES OF HEALTH

THE UNIVERSITY OF CHICAGO

Argonne NATIONAL LABORATORY

powered by amazon web services

3

# Agenda

- **Research data management challenges**
- **Globus: a high-level flyover**
- **File Transfer and Sharing: Accelerating and streamlining collaboration**
- **Data Publication: Enhancing reproducibility and discoverability**
- **Our sustainability challenge**
- **Globus campus deployment & intergation**
- **Deployment best practices: the Science DMZ**
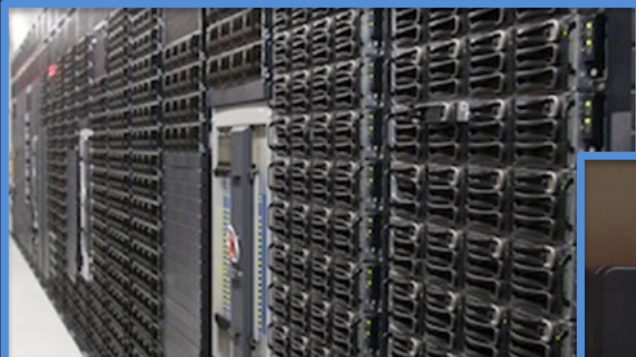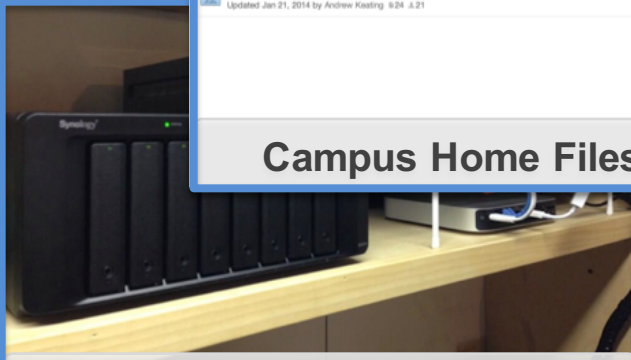- **Leveraging the Globus platform**

# Who are you?

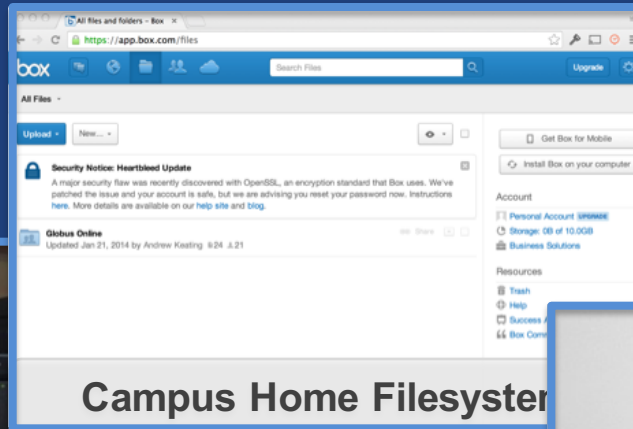# Research data management scenarios and challenges

# "I need to easily, quickly, & reliably move portions of my data to other locations."

**Research Computing HPC Cluster**

**Campus Home Filesystem**

**Lab Server**

**Personal Laptop**

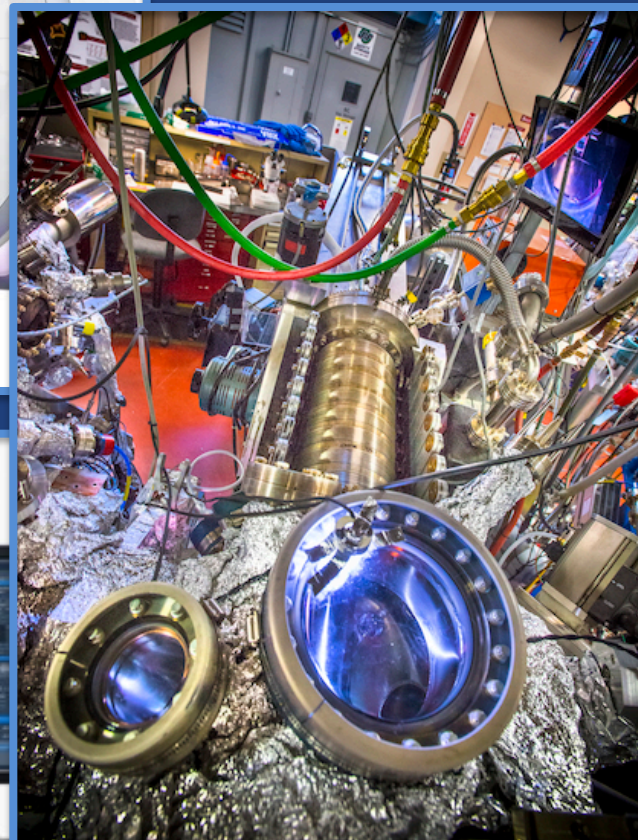**Desktop Workstation**

**XSEDE Resource**

**Public Cloud**

# "I need to get data from a scientific instrument to my analysis system."
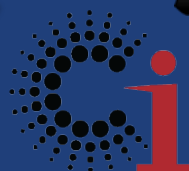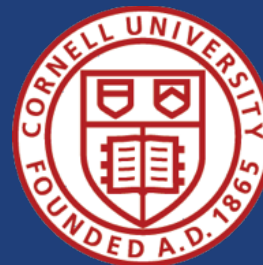
MRI

Advanced Light Source

Next Gen Sequencer

Light Sheet Microscope

"I need to easily and securely share my data with my colleagues at other institutions."

# "I need to publish my data so others can find/use/validate/reproduce it."

**Reference Dataset**

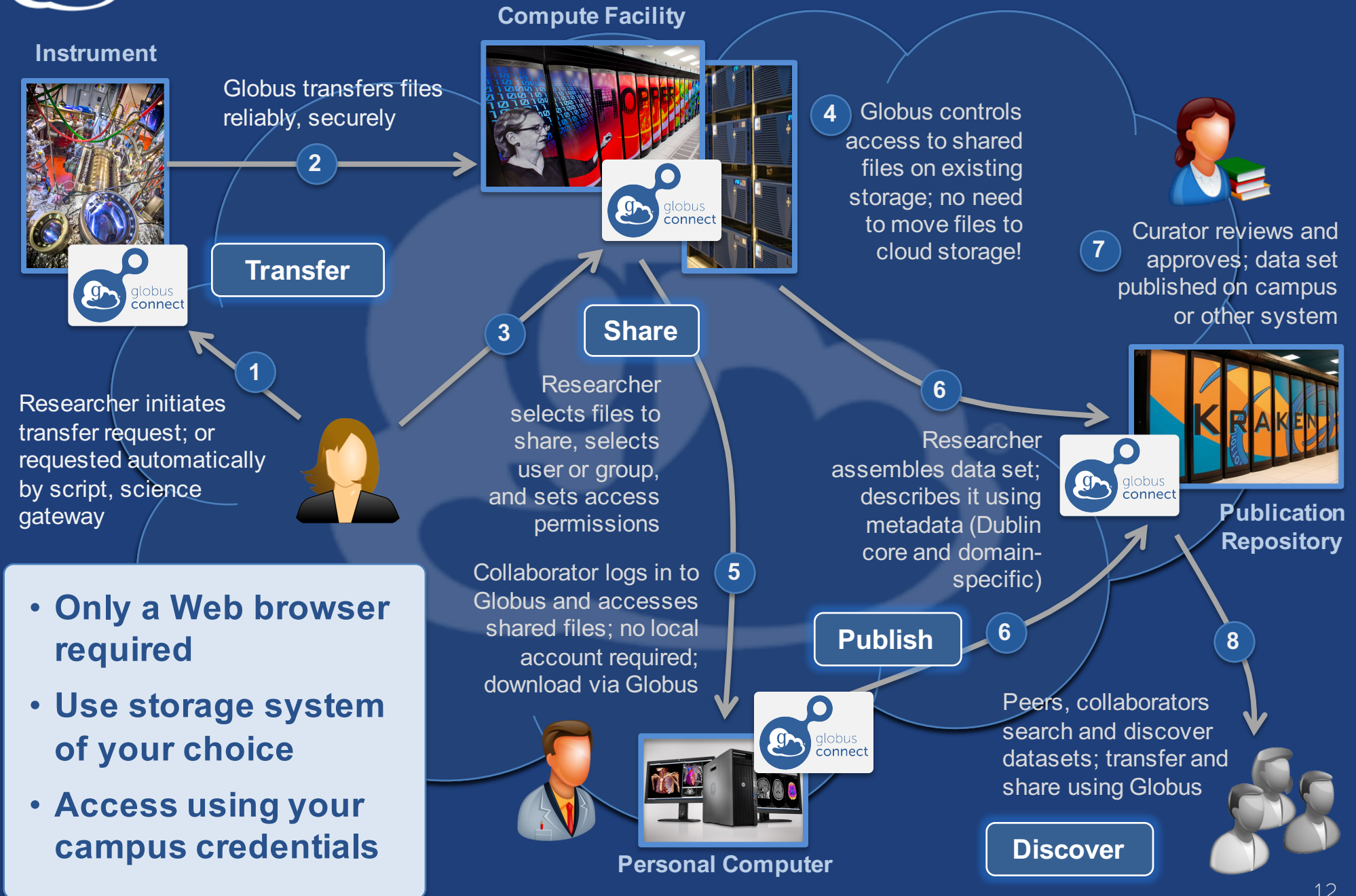**Scholarly Publication**

**Research Community Collaboration**

# Research data management today



Index?

# Globus and the research data lifecycle

**Instrument**

**Compute Facility**

**Transfer**

Globus transfers files reliably, securely

**2**

**4** Globus controls access to shared files on existing storage; no need to move files to cloud storage!

Curator reviews and approves; data set published on campus or other system

**7**

**1**

**3**

**Share**

**6**

Researcher initiates transfer request; or requested automatically by script, science gateway

Researcher selects files to share, selects user or group, and sets access permissions

Researcher assembles data set; describes it using metadata (Dublin core and domain-specific)

**Publication Repository**

Collaborator logs in to Globus and accesses shared files; no local account required; download via Globus

**5**

**6**

**Publish**

**8**

- **Only a Web browser required**
- **Use storage system of your choice**
- **Access using your campus credentials**

Peers, collaborators search and discover datasets; transfer and share using Globus

**Discover**

**Personal Computer**

Globus delivers…

Big data transfer, sharing, publication, and discovery…

…directly from your own storage systems…

…via software-as-a-service

# Globus is SaaS

- **Easy to access via Web browser**
  - Command line, REST interfaces for flexible automation and integration

- **New features automatically available**

- **Reduced IT operational costs**
  - Small local footprint (Globus Connect)
  - Consolidated support and troubleshooting

# Our focus: User Experience

**flickr** …for your photos

Google …for your office docs

NETFLIX …for your entertainment

globus …for your research data

# Accessing Globus and Moving Data

# Sign up & transfer files

1. **Go to: www.globus.org/signup**

2. **Create your Globus account**

3. **Validate e-mail address**

4. **Optional: Login with your campus/InCommon identity**

5. **Install Globus Connect Personal**

6. **Move files from vas#sc15 endpoint to your laptop**

# Sharing Data

# Lowering collaboration overhead

- **Grant collaborators access to data on systems without requiring local accounts**

- **No need to replicate or move data to separate system/cloud just for sharing**

- **Researchers manage "virtual" ACLs…**

- **Respect local system access controls**

# Share files

1. **Join the "Tutorial Users" groups**
   - Go to "Groups", search for "tutorial"
   - Select group from list, click "Join Group"

2. **Create a shared endpoint on your laptop**

3. **Grant your neighbor permissions on your shared endpoint**

4. **Access your neighbor's shared endpoint**

# Group Management

# Exercise 3: Create/configure group

1. **Create a group**
   – Go to globus.org/groups
   – Click "Create New Group"
   – Enter the group name and a short description
   – Set visibility to "all Globus members"

2. **Configure your group policies**
   – Select your group and click the "Settings" tab
   – Set requests to "a logged in Globus user"
   – Set approvals to "automatically if all policies are met"

3. **Ask your neighbor to join your group**

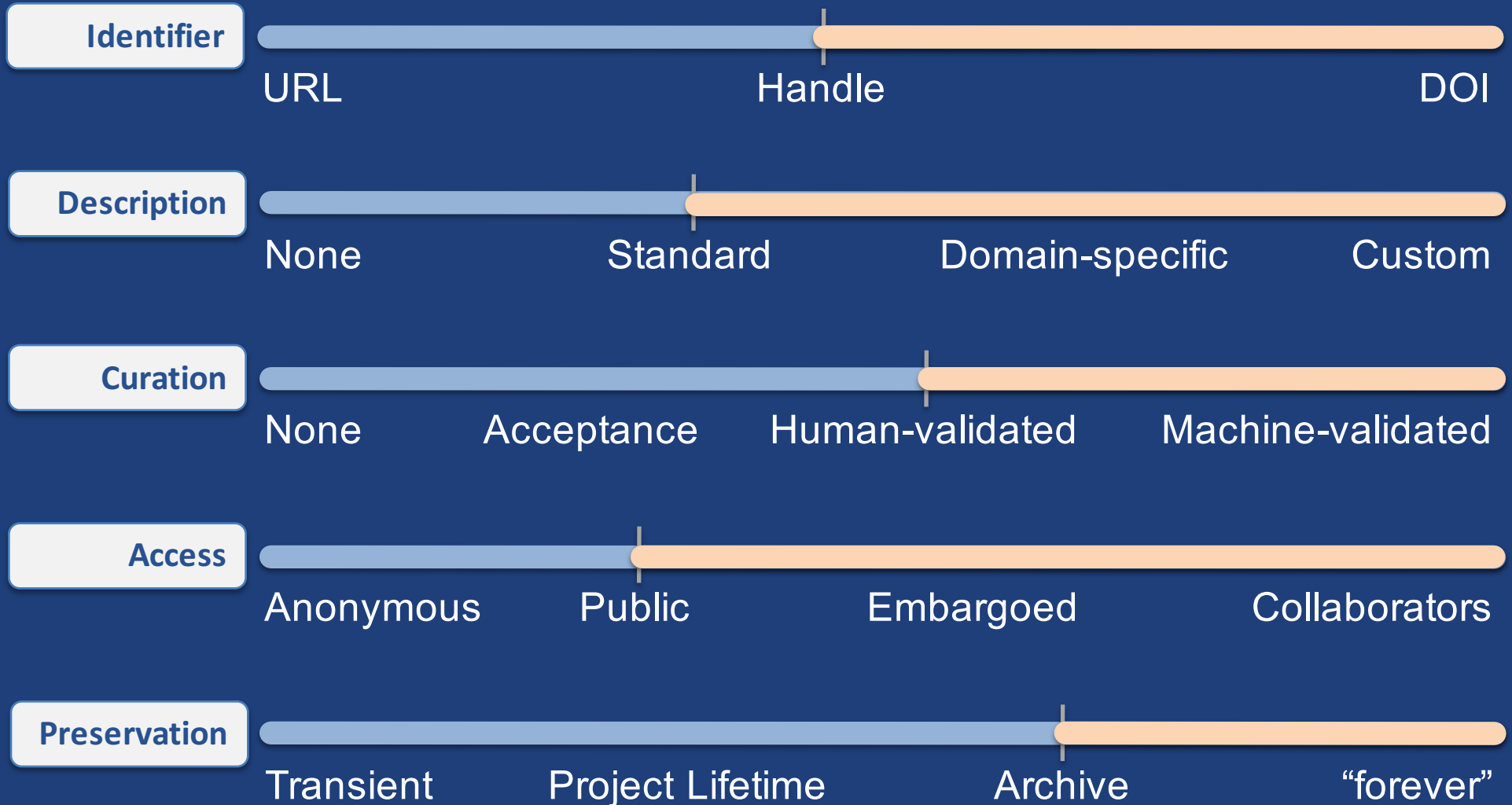4. **Grant permissions to the group on your shared endpoint**

5. **Confirm your neighbor can access your shared endpoint**

# Enhancing reproducibility and discoverability

# Globus data publication framework

**Identifier**

URL                Handle                DOI

**Description**

None         Standard       Domain-specific      Custom

**Curation**

None       Acceptance      Human-validated     Machine-validated

**Access**

Anonymous     Public       Embargoed      Collaborators

**Preservation**

Transient     Project Lifetime      Archive      "forever"

# Raw NGS output

Minimal metadata…

- Source environment
  - Instrument, timestamp,…
- Unique ID

High durability, low cost store

No curation

- Automated dataset acceptance

Identify…

- Handle

**Glacier**

# Upstream analysis

globus **genomics**

**Campus HPC**

```
##fileformat=VCFv4.0
##fileDate=20110705
##reference=1000GenomesPilot-NCBI37
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Desc
##INFO=<ID=DP,Number=1,Type=Integer,Desc
##INFO=<ID=AF,Number=.,Type=Float,Descri
```

Processing metadata…

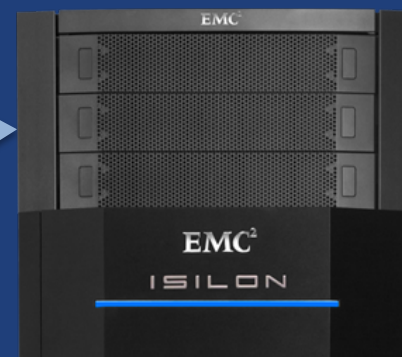- Pipeline description
- Tool parameters
- Exec environment

Moderate durability/cost

Automated curation

- Machine validated
- Exception review

Identify…

- URL

# Downstream analysis

Optional metadata…

- "Implicit" metadata
- Description through organization

Widely accessible stores

Team review

- Any collaborator may approve

Identify…

- Globus share

# Peer reviewed paper

(Re)format…

- PDF/A
- HDF
- …

Fully described…

- Dublin core metadata
- Domain metadata
- Provenance info

Replicated, public repositories

Formal, multi-step review

- Review → Update → Resubmit cycle

Persistent identifier

- DOI

# Globus publication - Initial release

| | Supported in GA release | | Consulting support | | Planned |
|---|---|---|---|---|---|

**Identifier**

URL        Handle        DOI

**Description**

None    Standard    Domain-specific    Custom

**Curation**

None    Acceptance    Human-validated    Machine-validated

**Access**

Anonymous    Public    Embargoed    Collaborators

**Preservation**

Transient    Project Lifetime    Archive    "forever"

# Publish a dataset

1. **Go to trial.publish.globus.org**

2. **Log in, click "Submit a New Dataset"**

3. **Select either of the Open Trial collections and continue**

4. **Accept the license terms**

5. **Enter required metadata to describe the dataset**

6. **Assemble data set from the vas#sc15 endpoint (or your own laptop if you installed Globus Connect Personal)**

7. **Complete the workflow and submit**

8. **Curators (a.k.a. presenters) will "review" your submission and publish**

9. **Search for your published dataset and browse the data**

# Globus: today and tomorrow

# Globus today…

| | | | |
|---|---|---|---|
| **4**<br>major services | **118 PB**<br>transferred | **20 billion**<br>files processed | **31,000**<br>registered users |
| **13**<br>national labs use Globus | **10,000**<br>active endpoints | **~350**<br>active daily users | **99.95%**<br>uptime |
| **35+**<br>institutional subscribers | **1 PB**<br>largest single transfer to date | **3 months**<br>longest continuously managed transfer | **130**<br>federated campus identities |

We are a non-profit, delivering a production-grade service to the non-profit research community

We are a non-profit, delivering a production-grade service to the non-profit research community

Our challenge:

**Sustainability**

# Globus Provider Subscriptions

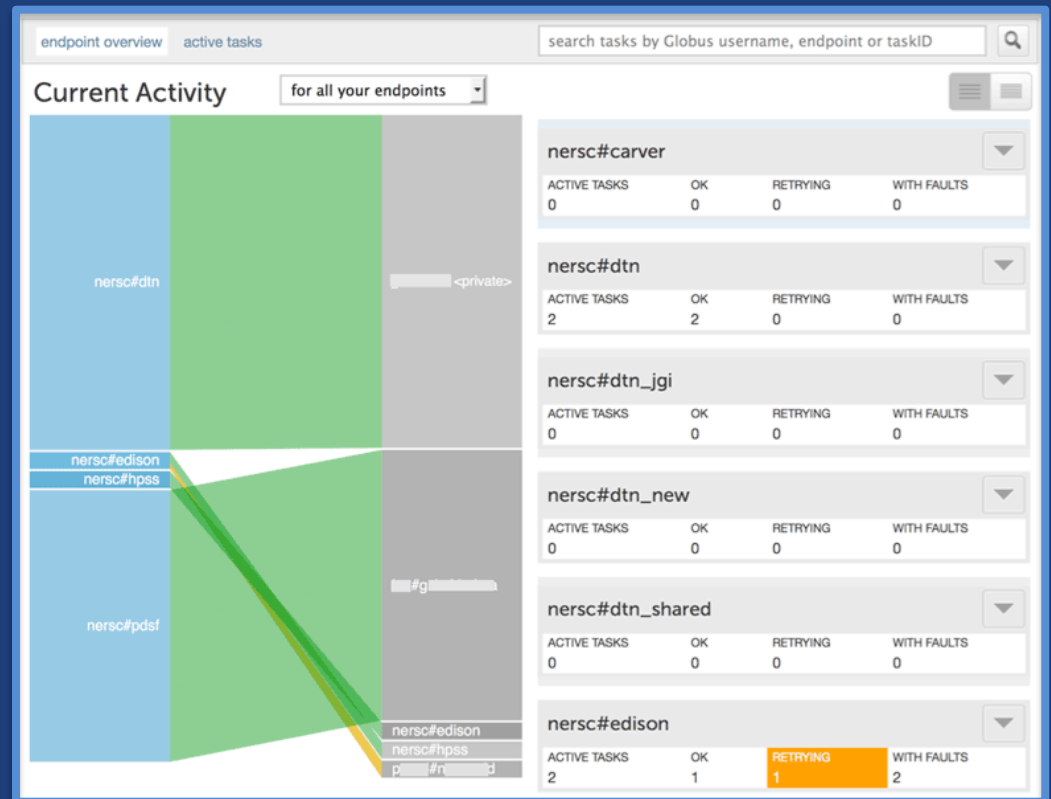- **Globus Provider Plan**
  - Shared endpoints
  - Data publication
  - Amazon S3 endpoints
  - Management console
  - Usage reporting
  - Priority support
  - Application integration

- **Branded Web Site**

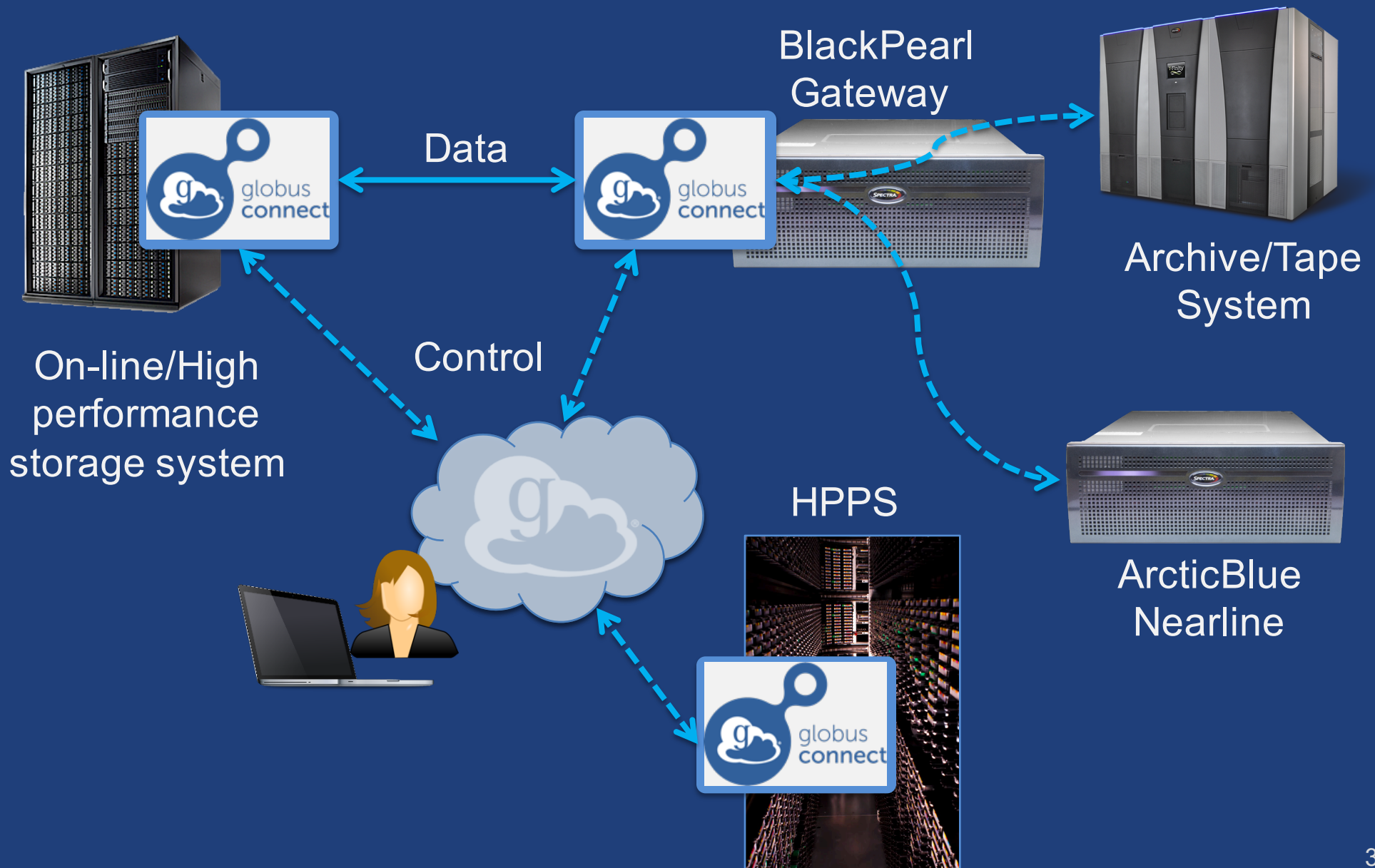- **Alternate Identity Provider (InCommon is standard)**

- **Mass Storage System optimization**



## globus.org/provider-plans

# Bridging the storage hierarchy

BlackPearl Gateway

Data

Control

On-line/High performance storage system

Archive/Tape System

HPPS

ArcticBlue Nearline

# Demonstration:
# Globus management console

# Demonstration:
# Bridging to Cloud Storage
# - Amazon S3: supported
# - Ceph: coming soon