# Research Data Management using Software-as-a-Service

*Vas Vasiliadis*
vas@uchicago.edu

THE UNIVERSITY OF CHICAGO

globus

Presentation material available at

globus.org/events/educause-2015

bit.ly/educause2015

# Thank you to our sponsors!

U.S. DEPARTMENT OF **ENERGY**

NSF

NATIONAL INSTITUTES OF HEALTH

ALFRED P. SLOAN FOUNDATION

S

1934

THE UNIVERSITY OF CHICAGO

Argonne

NATIONAL LABORATORY

powered by **amazon** web services

# Agenda

- **Research data management challenges**
- **Globus: a high-level flyover**
- **Accelerating and streamlining collaboration: transfer and sharing**
- **Enhancing reproducibility and discoverability: data publication**
- **Our sustainability challenge**
- **Campus deployment, security overview**
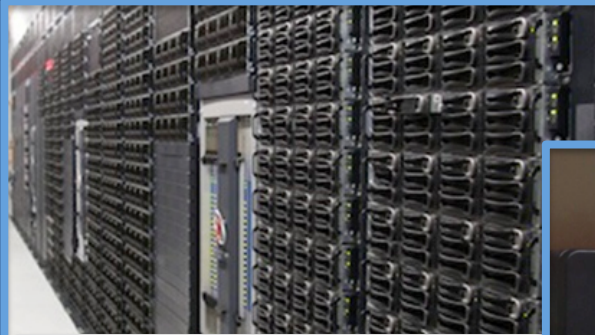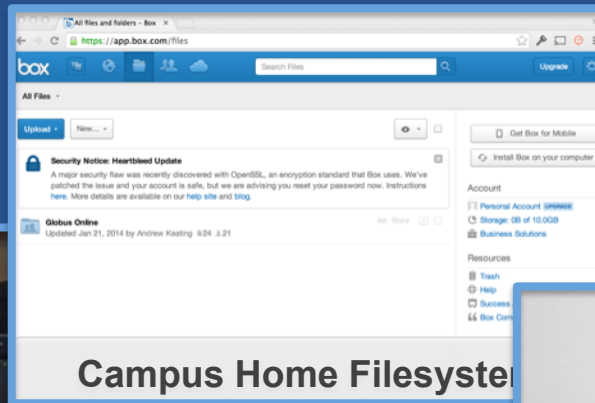- **Leveraging the Globus platform**

# Who are you?

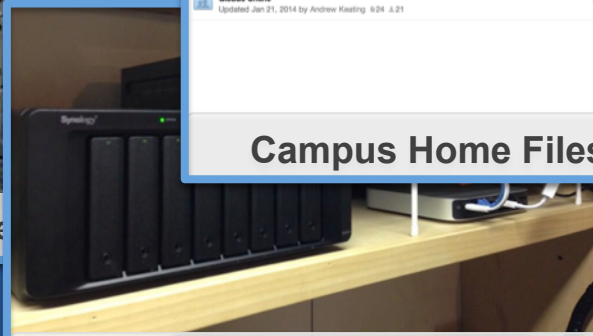# Research data management scenarios and challenges

# "I need to easily, quickly, & reliably move portions of my data to other locations."
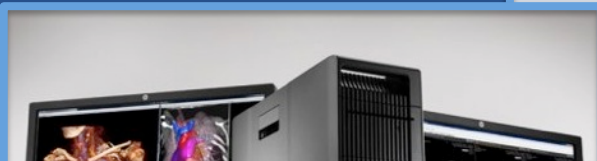
**Research Computing HPC Cluster**

**Campus Home Filesystem**

**Lab Server**

**Personal Laptop**

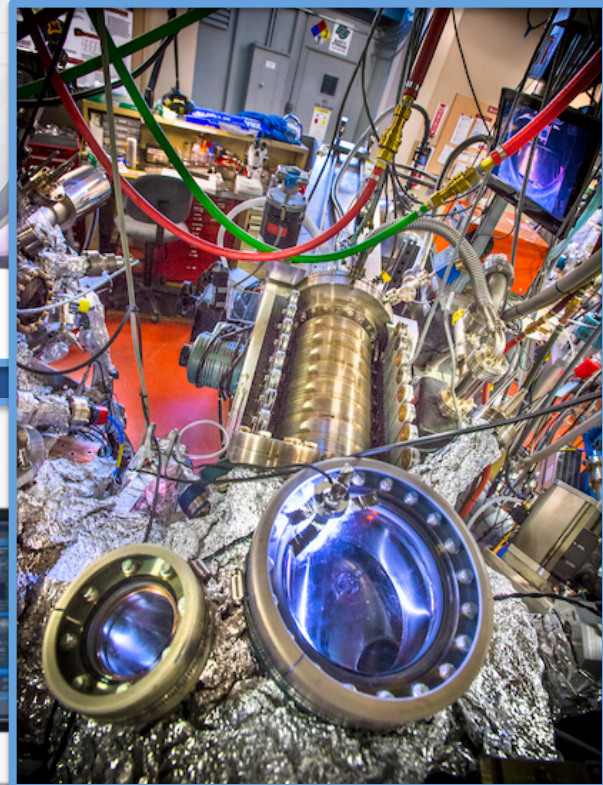**Desktop Workstation**

**XSEDE Resource**

**Public Cloud**

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

MRI

Advanced Light Source

Next Gen Sequencer

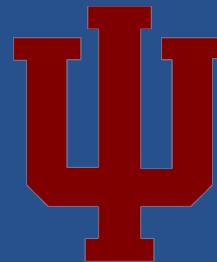Light Sheet Microscope

"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

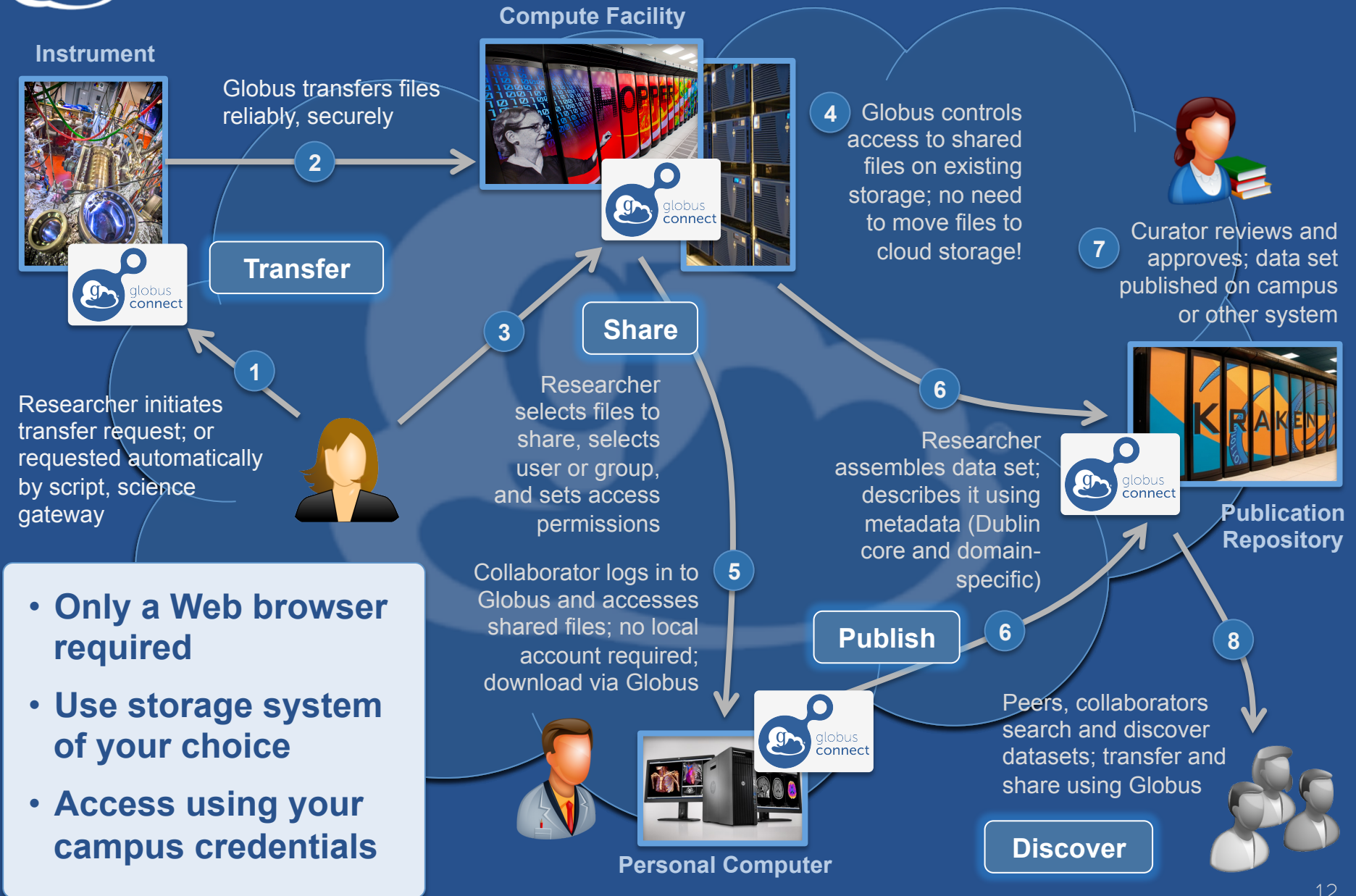Research Community Collaboration

# Research data management today



Index?

# Globus and the research data lifecycle

**Instrument**

**Compute Facility**

Globus transfers files reliably, securely

**2**

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

**Transfer**

**3**

**Share**

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

**1**

**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**

**Publish** **6**

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

**8**

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

**Discover**

**Personal Computer**

12

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

# Example: Scaling up

**Move datasets to supercomputer, national facility**

**Move results to campus (…)**

# 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#ebs 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

**Campus HPC**

globus **genomics**

```
##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
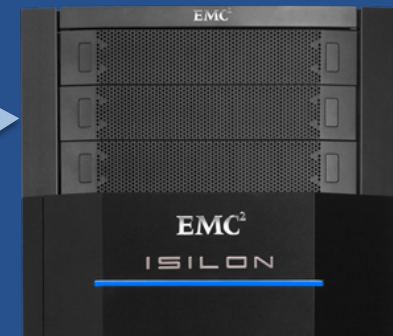- Tool parameters
- Exec environment

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**

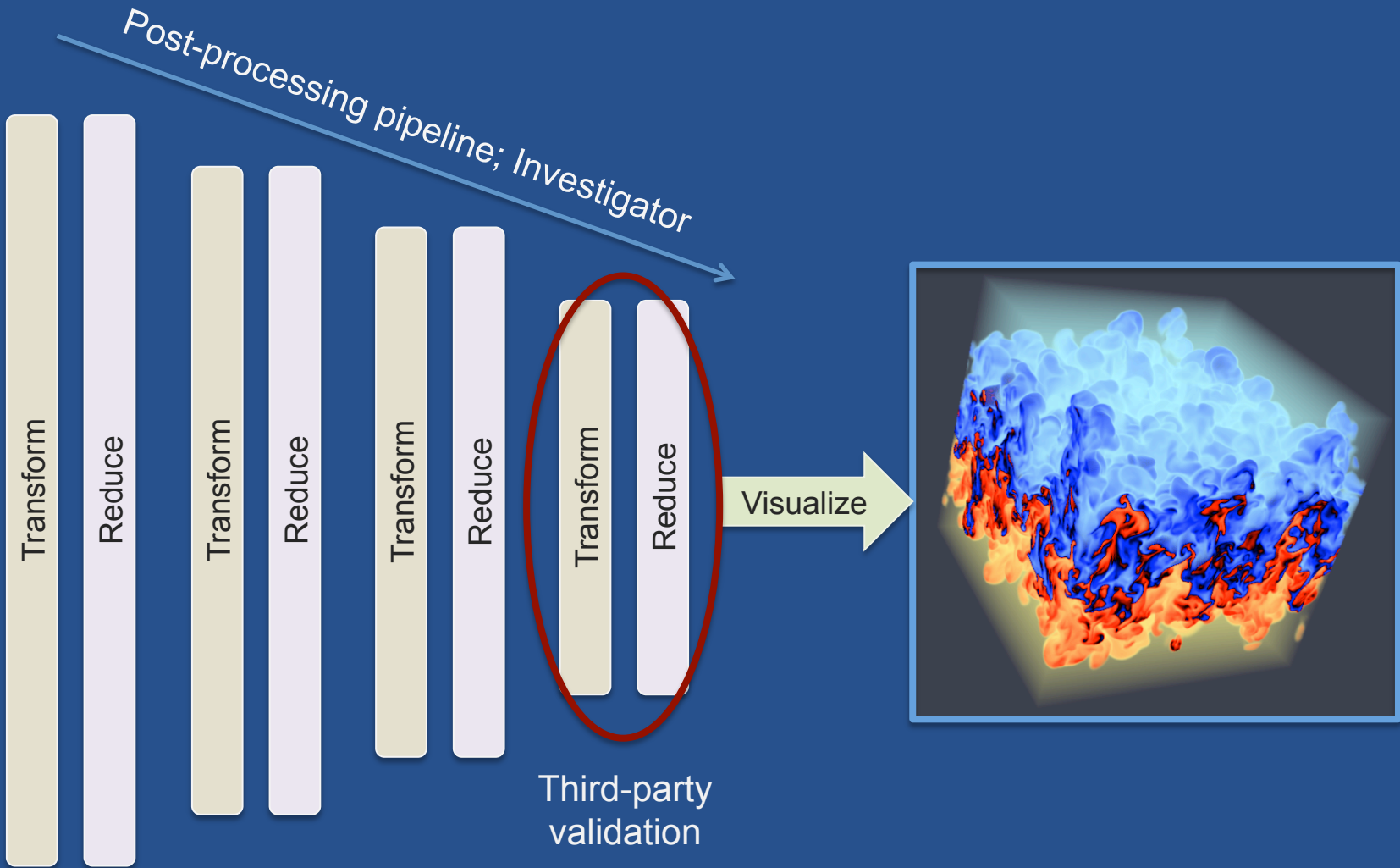Transient       Project Lifetime       Archive       "forever"

# 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#ebs 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**

# Reproducibility example

Post-processing pipeline; Investigator

Transform | Reduce | Transform | Reduce | Transform | Reduce | Transform | Reduce

Visualize

Third-party validation

# Reproducibility example

nek-workflow  /  Demo.ipynb

## Figure 1

Start by loading some boiler plate: matplotlib, numpy, scipy, json, functools, and a convenience class.

```
In [1]:  %matplotlib inline
         import matplotlib
         matplotlib.rcParams['figure.figsize'] = (10.0, 8.0)
         import matplotlib.pyplot as plt
         import numpy as np
         from scipy.interpolate import interp1d, InterpolatedUnivariateSpline
         from scipy.optimize import bisect
         import json
         from functools import partial
         class Foo: pass
```

And some more specialized dependencies:

1. Slict provides a convenient slice-able dictionary interface
2. Chest is an out-of-core dictionary that we'll hook directly to a globus remote using...
3. glopen is an open-like context manager for remote globus files

```
In [2]:  from chest import Chest
         from slict import CachedSlict
         from glopen import glopen, glopen_many
```

Configuration for this figure.

```
In [3]:  config = Foo()
         config.name     = "HighAspect/HA_visc/HA_visc"
         config.arch_end = "maxhutch#alpha-admin/~/pub/"
```

# Repository planning

**NEW YORK UNIVERSITY**

## *Mix-Model: Technology meets Pedagogy*

- **Central IT provides infrastructure**
  - Storage, Computing, Cluster, Servers...

- **Library responsible for data stewardship**
  - Collection, Acquisition, Search & Discovery, Metadata, Preservation...

- **Staff: technologists + librarians and subject specialists**

# 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<br>use Globus | **10,000**<br>active endpoints | **~350**<br>active daily users | **99.95%**<br>uptime |
| **35+**<br>institutional<br>subscribers | **1 PB**<br>largest single<br>transfer to date | **3 months**<br>longest<br>continuously<br>managed transfer | **130**<br>federated<br>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

Black Pearl Gateway

Archival/Near-line storage system

Data

On-line/High performance storage system

Control

HPPS

# Demonstration:

# Globus management console

Demonstration:

Bridging to Cloud Storage

- Amazon S3: supported

- Ceph: coming soon

# Campus Deployment Overview

# Globus Connect Server



- **Create endpoint in minutes; no complex software install**
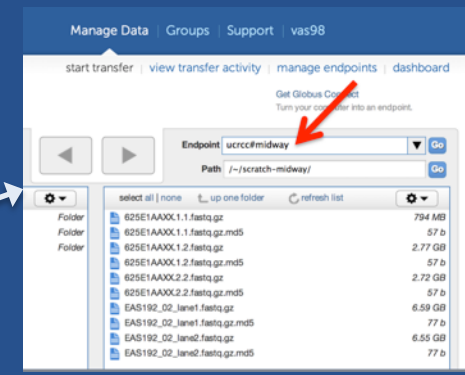- **Enable all users with local accounts to transfer files**
- **Native packages: RPMs and DEBs**

# Standard package installation

**1** **Install Globus Connect Server**
- Access server with root privileges
- Update package repos
- Install packages
- Setup Globus Connect Server

Server
(AWS EC2)

ssh

**2** **Log into Globus as end user/ researcher**

globus connect

globus

**3** **Access newly created endpoint**

**4** **Transfer a file**

Test Endpoint

Endpoint activation using MyProxy

1. Access endpoint
2. username password
3. campus username password — TLS handshake
4. username password
5. certificate
6. certificate — Transfer request
7. Authorization (resolve local user)
8. Access files
9. Control channel authorization — session certificate
10. Data transfer

Globus transfer and sharing hosted service

MyProxy Onine CA / PAM

GridFTP Server — certificate

Globus Connect Server

Local Authentication System (LDAP, RADIUS, Kerberos, ...)

Local Storage

Campus Cluster

GridFTP Server — Remote cluster with Globus Connect Server or laptop/PC with Globus Connect Personal

**Endpoint activation using MyProxy OAuth**

Campus Cluster

1 — Access endpoint → globus — Globus transfer and sharing hosted service

2 — OAuth redirect

3 — campus username password

7 — certificate

8 — certificate — Transfer request

11 — Control channel authorization — session certificate

OAuth Server

globus connect

4 — username password → MyProxy Onine CA / PAM

6 — certificate

GridFTP Server — certificate

5 — username password

9 — Authorization (resolve local user)

10 — Access files

12 — Data transfer

Local Authentication System (LDAP, RADIUS, Kerberos, …)

Local Storage

GridFTP Server — Remote cluster with Globus Connect Server or laptop/PC with Globus Connect Personal

# Common Configurations

- **Change endpoint name**

- **Customize filesystem access**

- **Enable sharing; set path restrictions**

- **Integrate with campus identity system**

- **Scale your campus deployment**
  - Data node(s)
  - Science DMZ

# Deployment best practice

## Science DMZ + Globus

**Border Router**

**Enterprise Border Router/Firewall**

perfSONAR

perfSONAR

**WAN**

10G

10GE

*Clean, High-bandwidth WAN path*

*Site / Campus access to Science DMZ resources*

10GE

perfSONAR

10GE

**Site / Campus LAN**

10GE

**Science DMZ Switch/Router**

*Per-service security policy control points*

perfSONAR

10GE

**High performance Data Transfer Node with high-speed storage**

*Details at: **fasterdata.es.net***

# Use(r)-appropriate interfaces

**Web**

**CLI**

```
laptop:~ ssh vas@cli.globusonline.org
$ Welcome to globusonline.org, vas.  Type 'help' for help.
$ endpoint-modify vas#ebs --organization="University of Chicago"
$
```

Globus service

```
GET /endpoint/go%23ep1
PUT /endpoint/vas#my_endpt
200 OK
X-Transfer-API-Version: 0.10
Content-Type: application/json
…
```

**Rest API**

# Quick look:

# Globus Command Line Interface (CLI)

# Globus Platform-as-a-Service

**Globus APIs**

**Globus Connect**

- Data Publication & Discovery
- File Sharing
- File Transfer & Replication

Identity, Group, and
Profile Management

Globus Toolkit

XSEDE
Extreme Science and Engineering
Discovery Environment

ci connect

globus genomics

UNIVERSITY OF EXETER

KBase

NeRSC

MICHIGAN

NCAR

XSEDE
Extreme Science and Engineering
Discovery Environment

INDIANA UNIVERSITY

# Building bridges to global communities

# What is the RDA?

- **Free and open access to 600+ datasets for climate and weather research**

- **Worldwide usage**

- **Multiple data access pathways**
  - HTTP (wget, cURL, etc.)
  - OPeNDAP, WCS, WMS
  - Web services (CLI, API)
  - Analysis on HPC systems (NCAR users)

# RDA Usage



- ## 2014
  - 17+ PB virtual processing
  - Web downloads: 7300 users, 750 TB served
  - 45,000 custom orders, 4000 users, 380 TB served

# Globus @ RDA

- **Single shared endpoint**
- **Data copied to subdirectories under endpoint source path**
- **Allow read permission to subdirectories under the shared endpoint**
- **ACLs managed programatically via Globus CLI**

# RDA Alternate Identity login

Courtesy of Thomas Cram, NCAR

# RDA Alternate Identity login

## NCAR Research Data Archive (RDA) MyProxy Client Authorization

Welcome to the NCAR RDA OAuth for MyProxy Client Authorization Page. The Client below is requesting access to your account. If you approve, please sign in with your RDA email/username and RDA password.

**Client Information**

Name: Globus Online
URL: https://www.globusonline.org

**NCAR RDA Email/Username** tcram@ucar.edu

**NCAR RDA Password** ••••••••••••

Sign In    Cancel

# Some early Globus supporters

# Enable your campus

- Signup: **globus.org/signup**

- Enable your resource: **globus.org/globus-connect-server**

- Need help? **support.globus.org**

- Subscribe to help make Globus self-sustaining
  **globus.org/provider-plans**

- Follow us: **@globusonline**

# Thank you