# Introduction to Globus:
## Research Data Management Software at the ALCF

## Rick Wagner

rick@globus.org
rpwagner@uchicago.edu
rwagner@anl.gov

# Research data management today

**Index?**

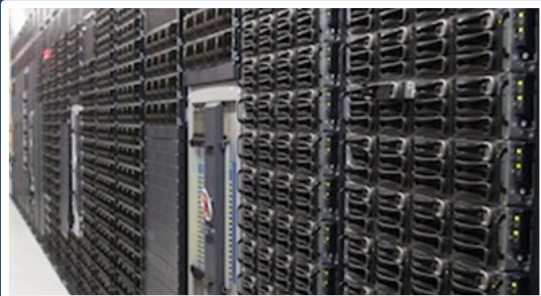How do we...
...move?
...share?
...discover?
...reproduce?

# Globus delivers...

Fast and reliable ~~big~~ data transfer, sharing, and platform services…

…directly from your own storage systems…

...via software-as-a-service using existing identities with the overarching goal of...

# …unifying access to data across tiers

**Research Computing HPC**

**National Resources**

Welcome to **compute** canada

ACCESS SERVICES AND EXPERTS

DISCOVER COMPUTE CANADA

**Personal Resources**

**Desktop Workstations**

**Mass Storage**

**Instruments**

amazon web services™ **S3**

*red cloud*

Windows Azure™

Google Cloud Platform

**Public Cloud**

# Storage Connectors - globus.org/connectors

**Current**

**Planned**

IBM Spectrum Scale

# Share with collaborators/community



External campus storage

Project repositories, replication stores

Public repositories

Public / private cloud stores

# Manage data from instruments


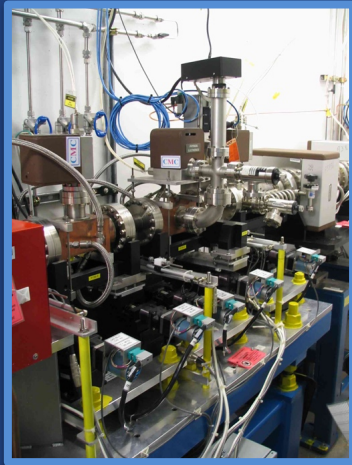Next-Gen Sequencer


MRI


Advanced Light Source


Cryo-EM


Light Sheet Microscope


Analysis store


High-durability, low-cost store


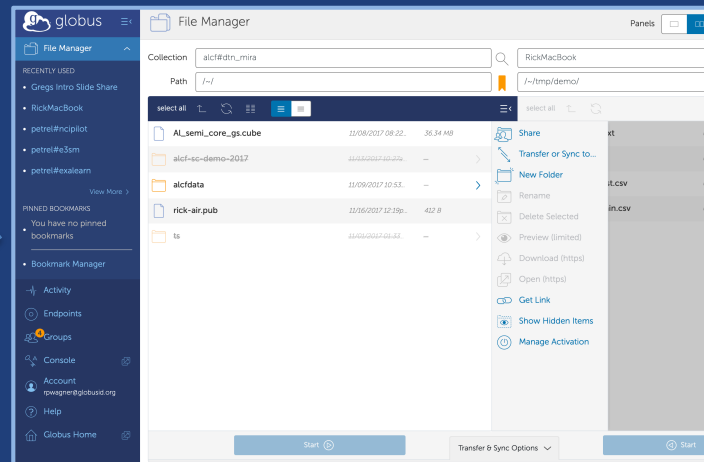Remote visualization


Personal system

# Use(r)-appropriate interfaces



**Web**

**CLI**

**Rest API**

Globus service

```
(globus-cli) jupiter:~ vas$ globus
Usage: globus [OPTIONS] COMMAND [ARGS]...

Options:
  -v, --verbose            Control level of output
  -h, --help               Show this message and exit.
  -F, --format [json|text] Output format for stdout. Defaults to text
  --map-http-status TEXT   Map HTTP statuses to any of these exit codes:
                           0,1,50-99. e.g. "404=50,403=51"

Commands:
  bookmark      Manage Endpoint Bookmarks
  config        Modify, view, and manage your Globus CLI config.
```
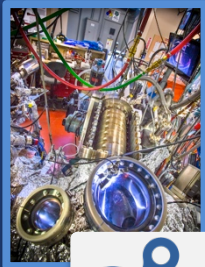
```
GET /endpoint/go%23ep1
PUT /endpoint/vas#my_endpt
200 OK
X-Transfer-API-Version: 0.10
Content-Type: application/json
...
```

# Globus SaaS / PaaS: 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!

**6** The Globus Command Line Interface, API sets, and Python SDK provide a platform…

**Transfer**

**Build**

**3**

**Share**

**1**

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

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

**7** … for building science gateways, portals and publication services.

**5**

Collaborator logs in to Globus and accesses shared files; no local account required; download via Globus
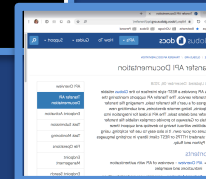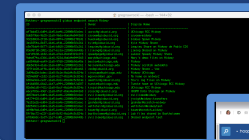
**8** Automating research workflows and ensuring those that need access to the data have it.

- **Use a Web browser or platform services**
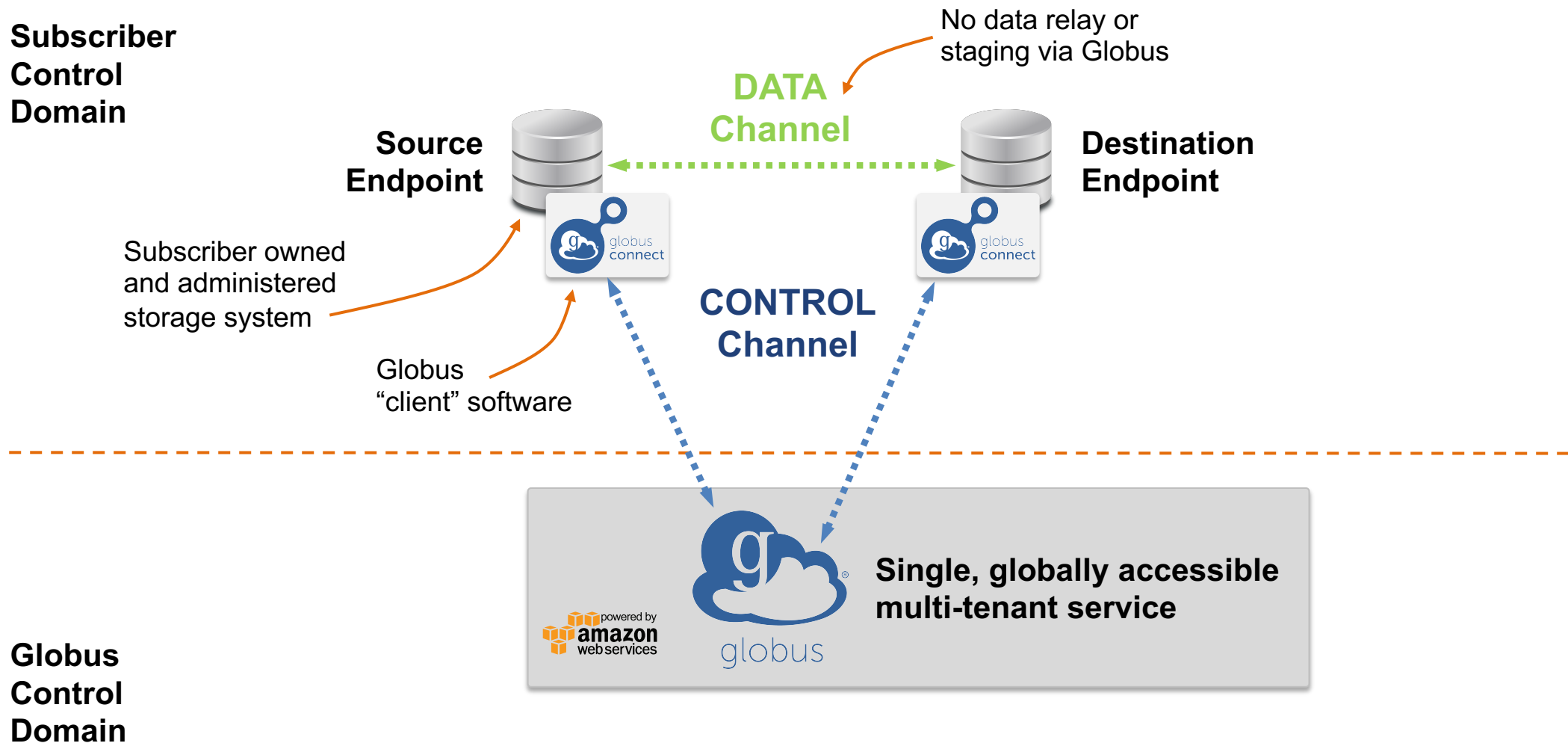- **Access any storage**
- **Use an existing identity**

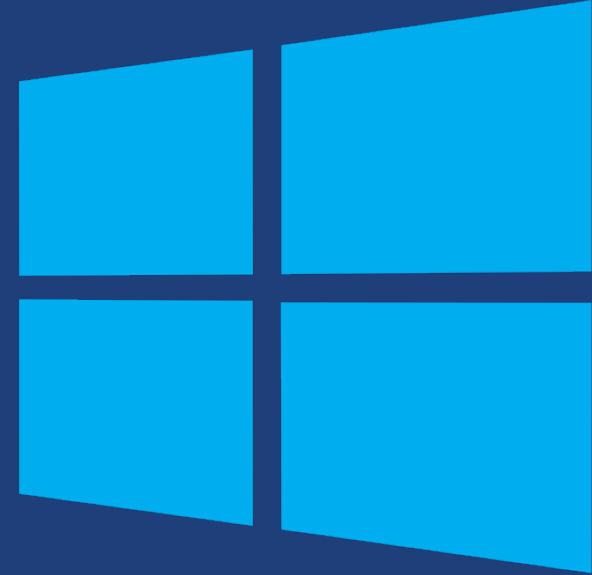**Personal Computer**

# Conceptual architecture: Hybrid SaaS

**Subscriber Control Domain**

**Source Endpoint**

Subscriber owned and administered storage system

Globus "client" software

**DATA Channel**

No data relay or staging via Globus

**Destination Endpoint**

**CONTROL Channel**

powered by **amazon** web services

**globus**

**Single, globally accessible multi-tenant service**

**Globus Control Domain**

# Conceptual architecture: Sharing

**External User Control Domain**

**Subscriber Control Domain**

DATA Channel

**Managed Endpoint**

CONTROL Channel

globus connect

globus connect

Subscriber managed filesystem permissions

**Globus Control Domain**

powered by amazon web services

globus

**Shared Endpoint**

Globus managed "overlay" permissions

...makes your storage system a Globus endpoint

# Endpoints (Collections)

- **Storage abstraction**
  - All transfers happen between two endpoints
  - Globus Connect instantiates endpoints

- **Collection ~= Endpoint**

- **Test / Demo Endpoints**
  - Globus Tutorial Endpoint 1
  - Globus Tutorial Endpoint 2
  - ESnet Test Endpoints
    - Contain file samples of various sizes

- **Globus Connect Personal**
  - Now your laptop is an endpoint
  - https://www.globus.org/globus-connect-personal

# Globus Connect Personal

- **Installers do not require admin access**
- **Zero configuration; auto updating**
- **Handles NATs**
- **Installs in seconds – easy to delete**

# The Globus Web App - Accounts

- **A Globus Account is**
  - A Primary Identity
  - Possible Linked Identities

- **Linking Identities**

- **Managing Identities**

- **Consents**

# Demonstration
**Identities**
**File Transfer**
**File Sharing**

# Activity Monitoring

- **Recent / History / Filter**

- **Drilling Down**
  - File transfer statistics
  - Overview
  - Event Log
  - Cancelling an active task

# Groups

- **What can they be used for?**
  - Sharing: Access permissions for more than one person
  - Roles: Endpoint management and monitoring

- **Groups**
  - Creating groups and setting the visibility
  - Members (invitations), Subgroups, Settings
  - Settings
    - Policies / Membership Fields / Terms & Conditions
  - Roles
    - Giving others authority over your groups

# Endpoint Sharing and Roles

- **Sharing**
  - Select the directory and create the "share"
  - A "share" is another type of endpoint
  - Share with: Users / Groups / All Globus Users

- **Roles**
  - Giving others (or groups of others) control or monitoring rights for your endpoints

# Bookmarks

- **Just like browser bookmarks – frequently used, or maybe not used frequently enough!**

- **Creating a bookmark**

- **Using a bookmark**

- **Sorting and Filtering**

- **Editing and Deleting**

# Globus Command Line Interface

```
(globus-cli) jupiter:~ vas$ globus
Usage: globus [OPTIONS] COMMAND [ARGS]...

Options:
  -v, --verbose            Control level of output
  -h, --help               Show this message and exit.
  -F, --format [json|text] Output format for stdout. Defaults to text
  --map-http-status TEXT   Map HTTP statuses to any of these exit codes:
                           0,1,50-99. e.g. "404=50,403=51"

Commands:
  bookmark        Manage Endpoint Bookmarks
  config          Modify, view, and manage your Globus CLI config.
  delete          Submit a Delete Task
  endpoint        Manage Globus Endpoint definitions
  get-identities  Lookup Globus Auth Identities
  list-commands   List all CLI Commands
  login           Login to Globus to get credentials for the Globus CLI
  logout          Logout of the Globus CLI
  ls              List Endpoint directory contents
  mkdir           Make a directory on an Endpoint
  rename          Rename a file or directory on an Endpoint
  task            Manage asynchronous Tasks
  transfer        Submit a Transfer Task
  version         Show the version and exit
  whoami          Show the currently logged-in identity.
```

**Open source, uses Python SDK**

**docs.globus.org/cli**

**github.com/globus/ globus-cli**

# The Globus CLI

- **Installation**
  - docs.globus.org/cli/installation
  - Prerequisites

- **Logging On (remember the consents?)**
  - globus login / logout

- **Getting help / list of commands**
  - globus –help
  - globus list-commands

- **Doing something**
  - It all about the UUIDs
  - Don't forget the file paths!

# The Globus CLI – Let's do a few things...

- **Find endpoints**
  - globus endpoint search Midway
  - globus endpoint search ESNet
  - globus endpoint search --filter-scope=recently-used

- **Find endpoint contents**
  - globus ls af7bda53-6d04-11e5-ba46-22000b92c6ec
  - globus ls af7bda53-6d04-11e5-ba46-22000b92c6ec:RMACC2018

- **Transfer a file**
  - From ESnet Read-Only Test DTN at CERN to Midway
  - Note the specific paths
  - globus transfer d8eb36b6-6d04-11e5-ba46-22000b92c6ec:/~/data1/1M.dat af7bda53-6d04-11e5-ba46-22000b92c6ec:/~/1M.dat

- **Transfer a directory**
  - From Globus Tutorial Endpoint 2 to Midway (create directory and contents)
  - globus transfer --recursive ddb59af0-6d04-11e5-ba46-22000b92c6ec:/~/sync-demo af7bda53-6d04-11e5-ba46-22000b92c6ec:/~/syncDemo

- **https://docs.globus.org/cli/examples/**

# Globus CLI

- **Easy install and updates**

- **It's a native application distributed by Globus**
  - https://docs.globus.org/cli/
  - https://github.com/globus/globus-cli

- **Command *globus login* gets access tokens and refresh tokens**
  - Stores the token locally (~/.globus.cfg )
  - The CLI "acts as" the logged in user

- **All interactions with the service use the tokens**
  - Tokens for Globus Auth and Transfer services

- **Command *globus logout* deletes those**

- **https://docs.globus.org/cli/examples/**

- **https://github.com/globus/automation-examples**

# Demonstration

# **Globus CLI**

# Industry software builds on platform services



cloud4scieng.org

Globus delivers… with applications and as a platform…

Fast and reliable data transfer, sharing, and file management…

…directly from your own storage systems…

…via software-as-a-service using existing identities.

# How can I integrate Globus into my research workflows?

Globus serves as…

A platform for building science gateways, portals and other web applications in support of research and education.

# Globus Platform-as-a-Service

**Globus Transfer API**

. . .

**Data Discovery**

**File Sharing**

**File Transfer & Replication**

**Globus Auth API**

**Globus Connect**

Integrate file transfer and sharing capabilities into scientific web apps, portals, gateways, etc...

Use existing institutional ID systems in external web applications

# Example web apps that leverage Globus

# Globus Transfer API Set

- **Doc**
  - https://docs.globus.org/api/transfer/

- **Sample data portal**
  - https://docs.globus.org/modern-research-data-portal/

- **Jupyter notebook**
  - https://github.com/globus/globus-jupyter-notebooks

# Globus Auth API Set

- **Doc**
  - https://docs.globus.org/api/auth/

- **Sample data portal**
  - https://docs.globus.org/modern-research-data-portal/

- **Native app examples**
  - https://github.com/globus/native-app-examples

# *Petrel:*
## A Programmatically Accessible Research Data Service

**Petrel** online store
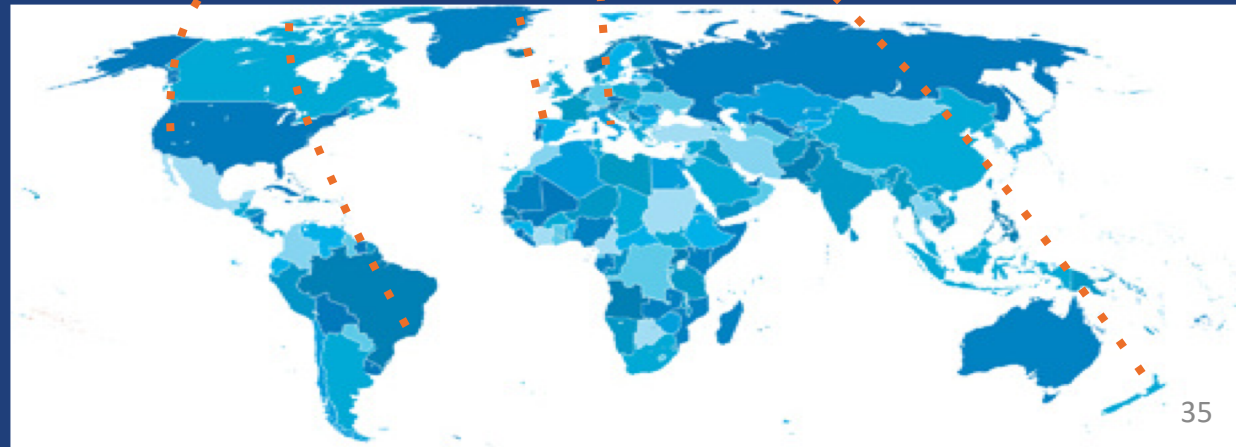
94 Gbit/s Petrel—Blue Waters



- 32 Nodes with ~~7~~ 3.2 PB usable storage
- ~~G~~ ~~P~~FS and Globus — Ceph
- 100TB allocation per project
- Transfer and sharing data with collaborators
- Federated login
- Self-managed by PIs

**PETREL**
Data Management and Sharing Pilot

3.2 petabytes
100 Gbps

## https://petrel.alcf.anl.gov

35

A bit of Globus history
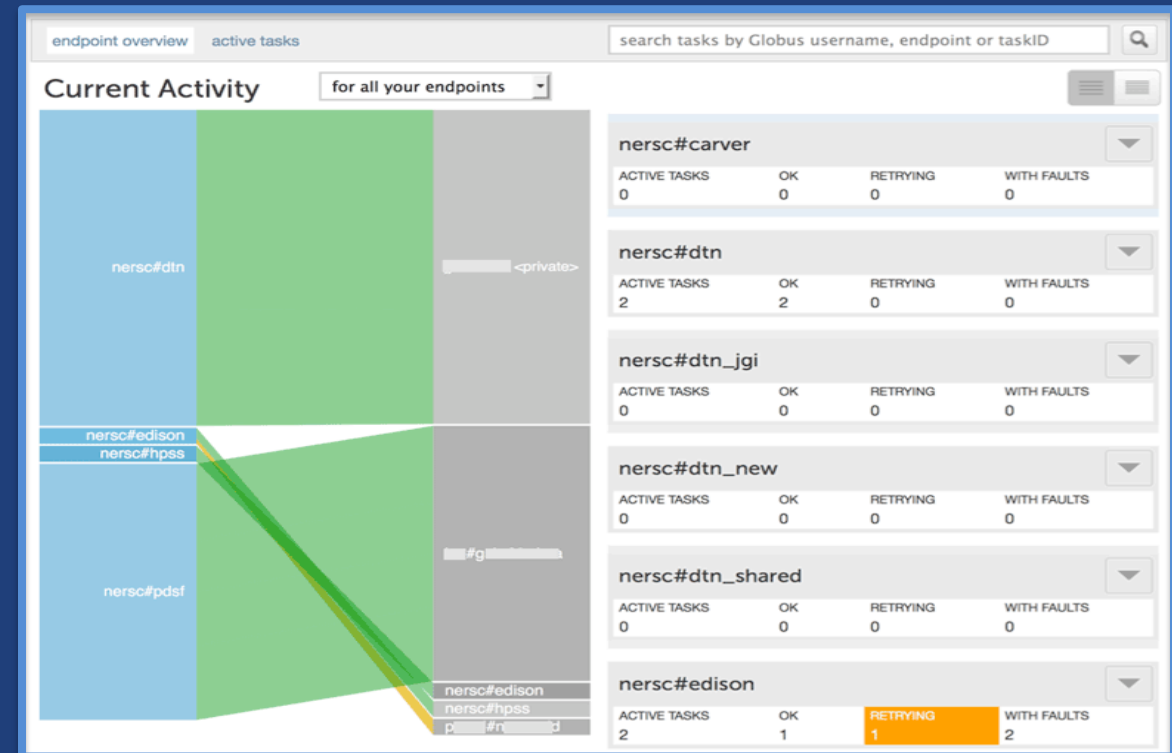
# Globus sustainability model

- **Standard Subscription**
  - Shared endpoints
  - Management console
  - Usage reporting
  - Priority support
  - Application integration
  - HTTPS support (coming soon)
- **Branded Web Site**
- **Premium Storage Connectors**
- **Alternate Identity Provider (InCommon is standard)**

# The path to sustainability

# Globus by the numbers
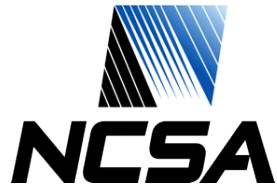
**1,042** most shared endpoints at a single institution

**635+ PB** transferred

**73 billion** files processed

**1,700** active GCS endpoints

**100+** subscribers

**100,000+** users

**3 months** longest running transfer

**18,000** active GCP endpoints

**500+** identity providers

**1 PB** largest single transfer to date

**8,000** active shared endpoints

**99.9%** availability

# ALCF Globus Resources

- **Documentation**
  - https://www.alcf.anl.gov/user-guides/data-transfer
  - https://www.alcf.anl.gov/user-guides/using-globus
- **Endpoints**
  - Theta: `alcf#dtn_theta`
  - Mira: `alcf#dtn_mira`
  - Cetus: `alcf#dtn_mira`
  - Cooley: `alcf#dtn_mira`
  - Vesta: `alcf#dtn_vesta`
  - HPSS: `alcf#dtn_hpss`

# Globus support resources

- **Globus documentation: docs.globus.org**

- **Helpdesk and issue escalation: support@globus.org**

- **Mailing Lists**
  - https://www.globus.org/mailing-lists

- **Customer engagement team**

- **Globus professional services team**
  - Assist with portal/gateway/app architecture and design
  - Develop custom applications that leverage the Globus platform
  - Advise on customized deployment and integration scenarios