# Globus Research Data Management: Introduction and Service Overview

**XSEDE**

Extreme Science and Engineering
Discovery Environment

Steve Tuecke
Vas Vasiliadis

globus

Presentations and other useful information available at

**globus.org/events/xsede15/tutorial**

# 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

3

# Agenda

- **Research data management scenarios and challenges**

- **Introduction to Globus**

- **Demonstrations and Exercises**

  – Accessing Globus and Transferring Files

  – File sharing and Group Management

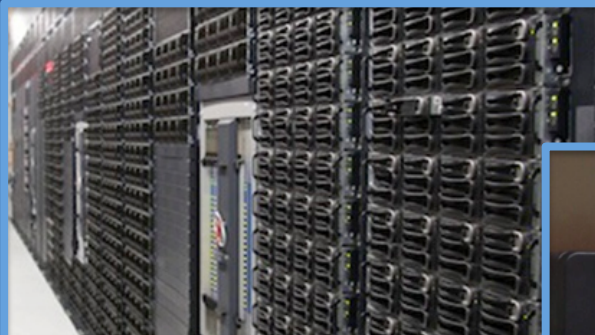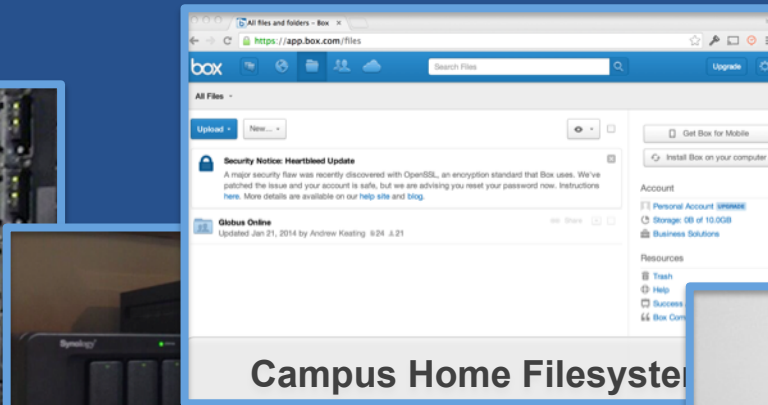  – Data publication and discovery

- **Globus: today and tomorrow**

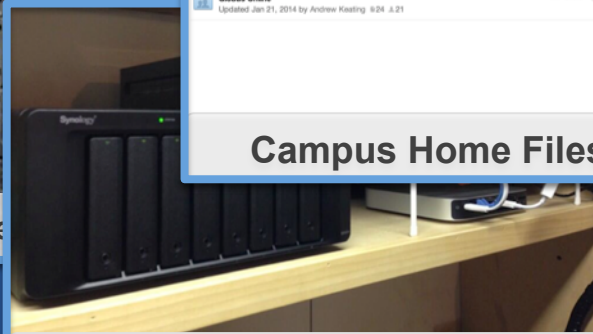# Research data management scenarios and challenges

"I need to easily, quickly, & reliably move or mirror portions of my data to other places."

Research Computing HPC Cluster

Campus Home Filesystem

Lab Server

Personal Laptop

Desktop Workstation
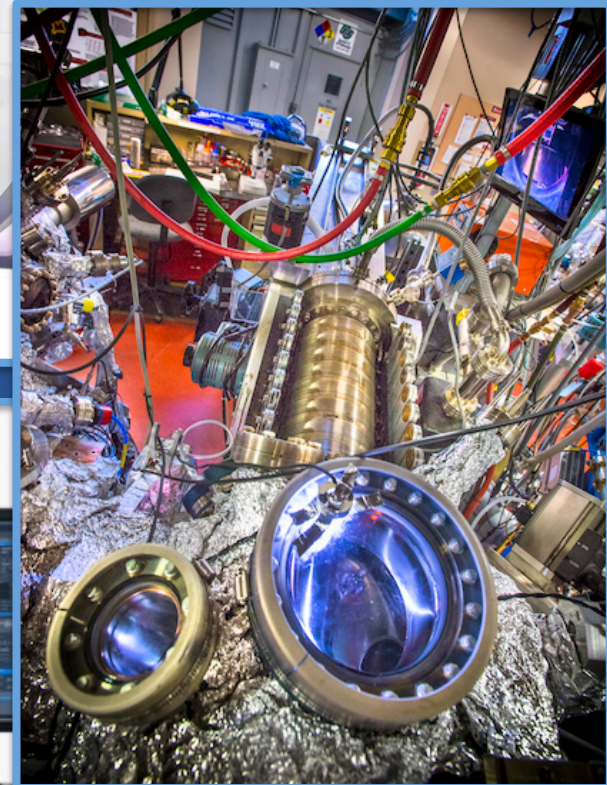
XSEDE Resource

Public Cloud

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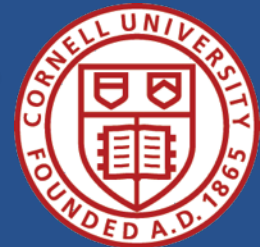
MRI

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 a good place to store / backup / archive my (big) research data, at a reasonable price."

Campus Store

Mass Store

Public Cloud Archive

# "I need to publish my data so that others can find it and use it."

Reference Dataset

Scholarly Publication

Active Research Collaboration

# Globus introduction and demonstration

Globus is...

Research data management...

...delivered via SaaS

Globus delivers...

Big data transfer, sharing, publication, and discovery...

...directly from your own storage systems

# It's about the user experience...

**flickr** ...for your photos

**Gmail** by Google ...for your e-mail

**NETFLIX** ...for your entertainment
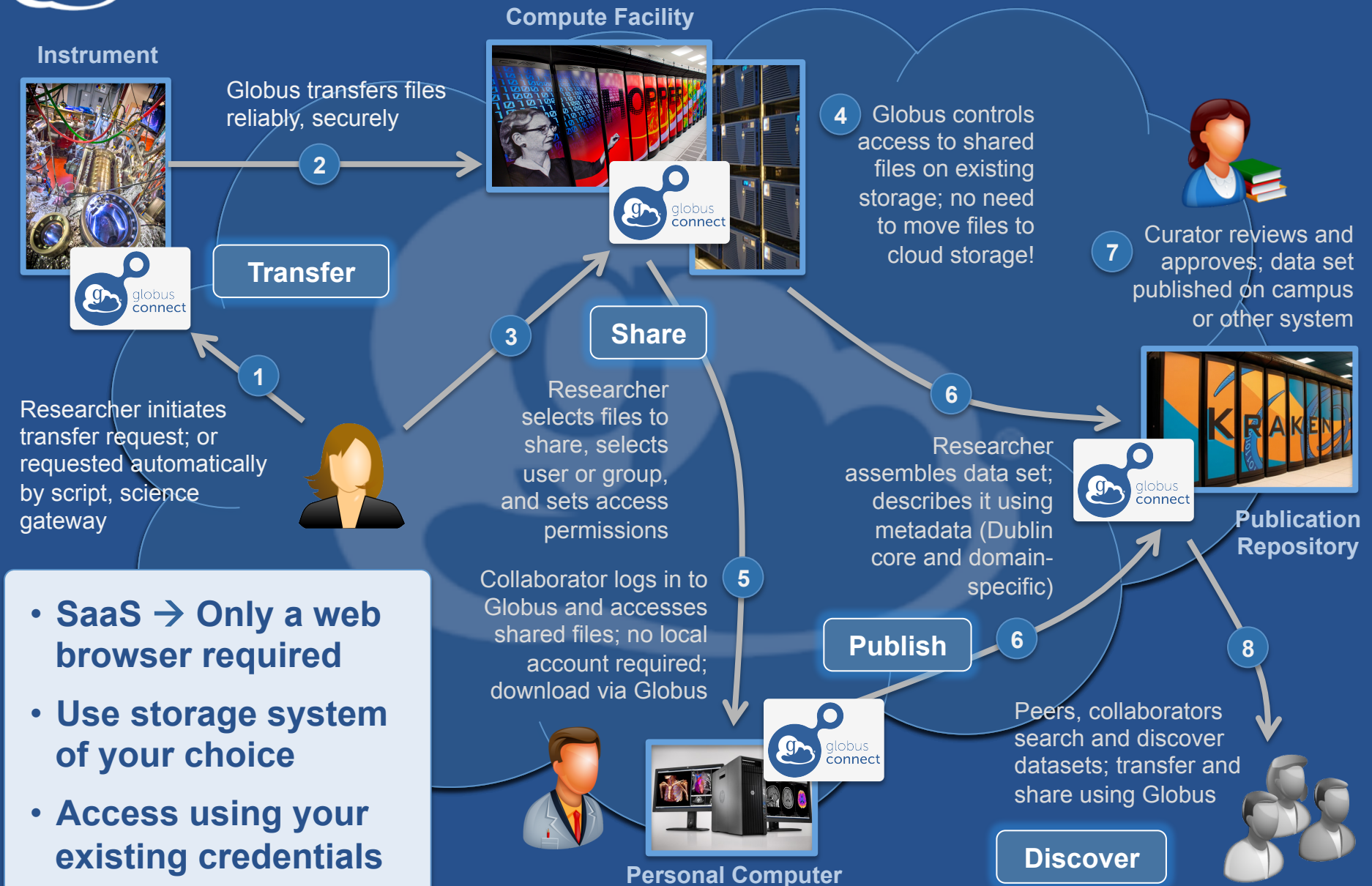
**globus** ...for your research data

# Globus is SaaS

- **Web, command line, and REST interfaces**

- **Reduced IT operational costs**

- **New features automatically available**

- **Consolidated support & troubleshooting**

- **Easy to add your laptop, server, cluster, supercomputer, etc. with Globus Connect**

# Globus and the research data lifecycle

**Instrument**

**Compute Facility**

Globus transfers files reliably, securely

**2**

**Transfer**

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

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

**1**

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

**3**

**Share**

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

**6** 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**

- **SaaS → Only a web browser required**
- **Use storage system of your choice**
- **Access using your existing credentials**

**8**

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

**Personal Computer**

**Discover**

# Globus and XSEDE

- **Globus endpoints available on all XSEDE storage resources**

- **Accessed using your XSEDE user portal username/password (see demo)**

- **First service to pass XSEDE acceptance test**
  - Available for file transfer
  - File sharing in final stages of approval

# Demonstration:
# - Accessing Globus
# - File Transfer

# Exercise: 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 file(s) from esnet#lbl-diskpt1 to your laptop**

# Demonstration:

- File Sharing

- Group Management

# Exercise: File sharing

- **Join the "Tutorial Users" group**
  - Go to "Groups" and search for "tutorial"

- **Enable sharing in Globus Connect Personal configuration**
  - e.g. on Mac OS: Preferences → Access → Shareable

- **Create a shared endpoint on your laptop**

- **Grant your neighbor permissions on your shared endpoint**

- **Access your neighbor's shared endpoint**

# Globus Data Publication and Discovery

# 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"

# Scenario 1: Genomics analysis

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
```

**Campus HPC**

Processing metadata…

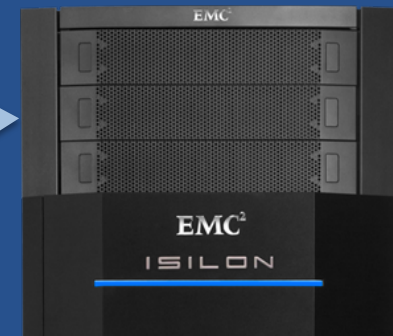- Pipeline description
- Tool parameters
- Exec environment

Moderate durability/cost

Automated curation

- Machine validated
- Exception review

Identify…

- URL

# Scenario 2: 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 – Current 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"**

# Demonstration:
# Data Publication
# and Discovery

# Globus: today and tomorrow

# Globus today…

## ~ 100PB moved

## >10,000 endpoints

## >300 active users/day

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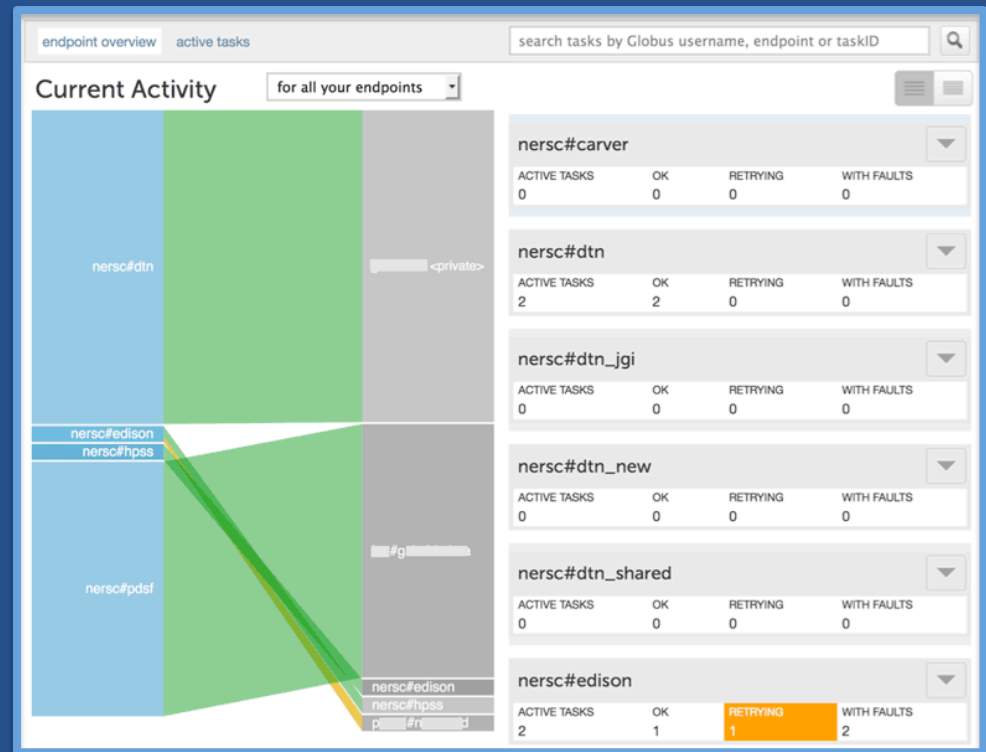
- **Globus Provider Plan**
  - Shared endpoints
  - Data publication
  - Amazon S3 endpoints
  - Management console
  - Usage reporting
  - Priority support
  - Application integration

- **Branded Web Site**

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

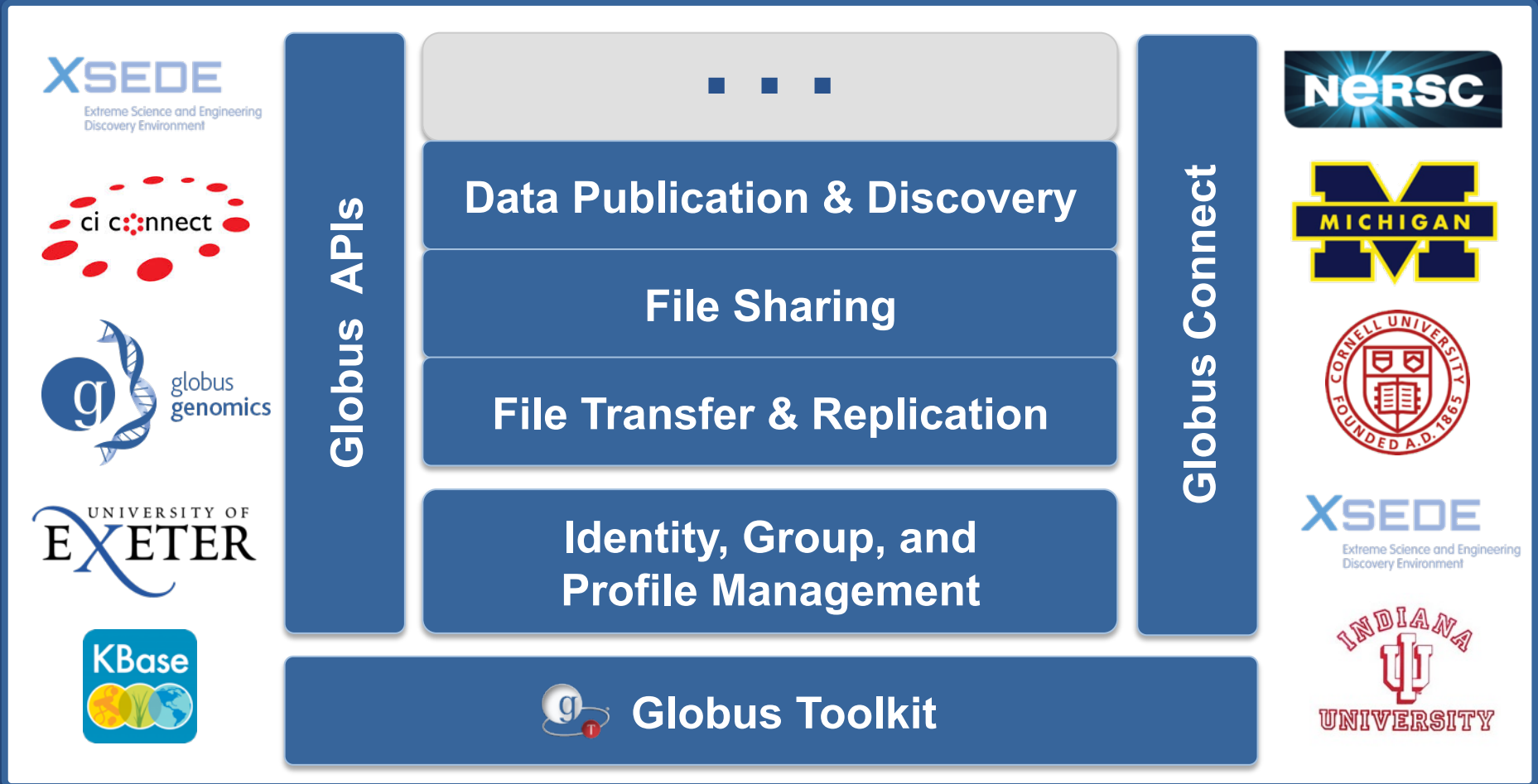- **Mass Storage System optimization**

## globus.org/provider-plans

# Demonstration:

# Globus management console

# Globus Platform-as-a-Service

**Globus APIs**

**Globus Connect**

Data Publication & Discovery

File Sharing

File Transfer & Replication

Identity, Group, and Profile Management

Globus Toolkit

# Some early Globus adopters

# End: Introduction and Service Overview