Our mission is to…

increase the efficiency and effectiveness of researchers engaged in data-driven science and scholarship through **sustainable** software

Development is funded by...

U.S. DEPARTMENT OF **ENERGY**

NSF

THE UNIVERSITY OF **CHICAGO**

NIH

powered by **amazon** web services

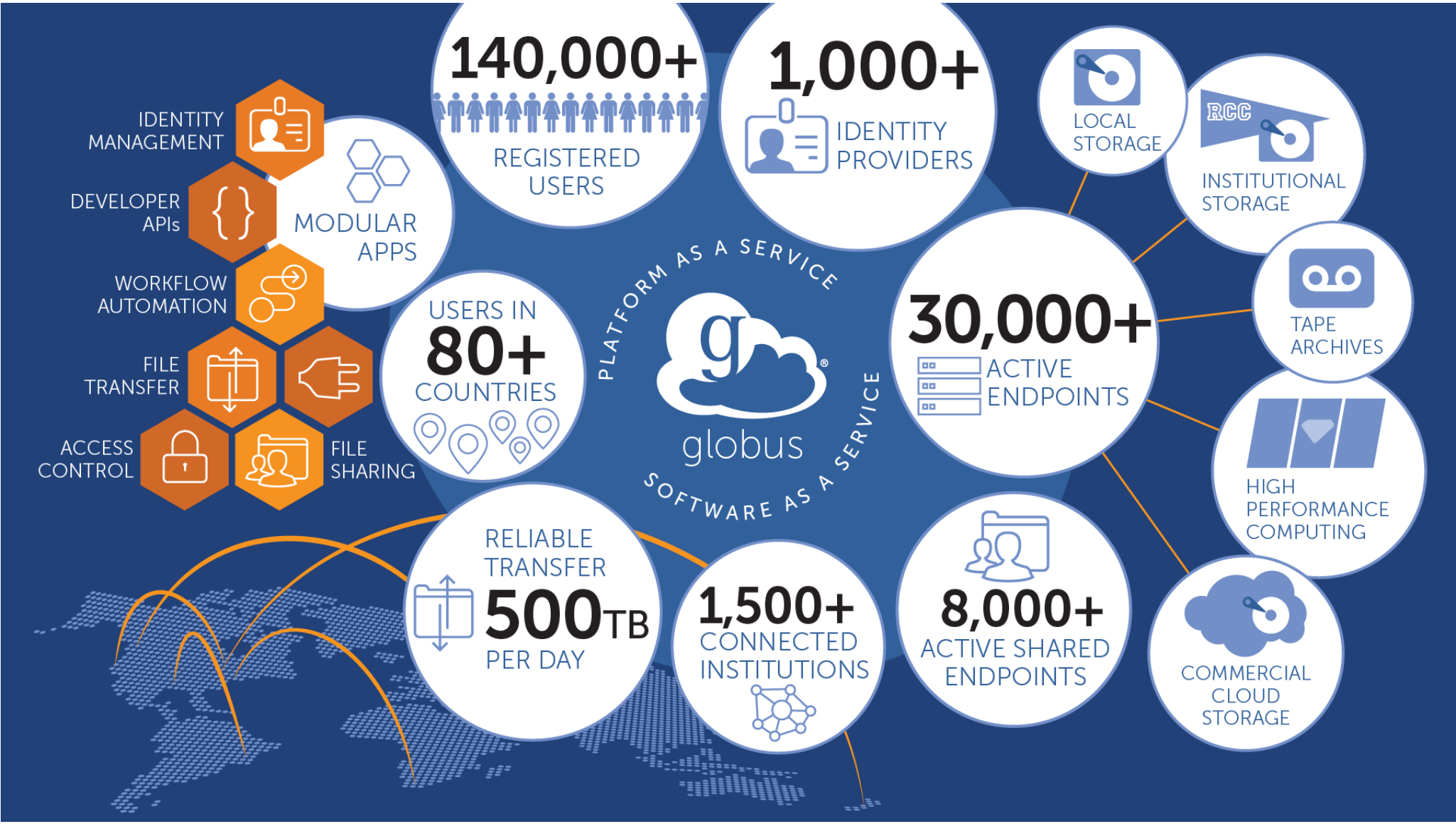ALFRED P. SLOAN FOUNDATION — S — 1934

NIST National Institute of Standards and Technology U.S. Department of Commerce

**Argonne** NATIONAL LABORATORY

# Operations are funded by subscribers

**100x ALL PRINTED MATERIAL** of the Library of Congress

**400 HUMAN BRAINS** worth of memory storage

**137 years** of observational data from the Rubin Observatory in Chile

**237,823 years** of non-stop video calls

**1,000,000,000,000,000,000**

1 EXABYTE – A QUINTILLION BYTES TRANSFERRED BY GLOBUS

**266 BILLION** human genomes worth of sequence information

Research data management simplified.

**66.7 YEARS** of the Large Hadron Collider's experimental data
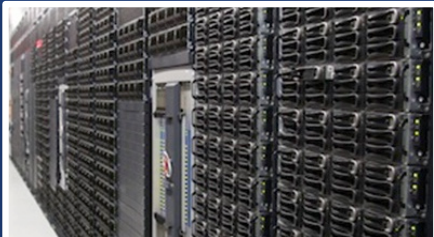
# Key capabilities

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
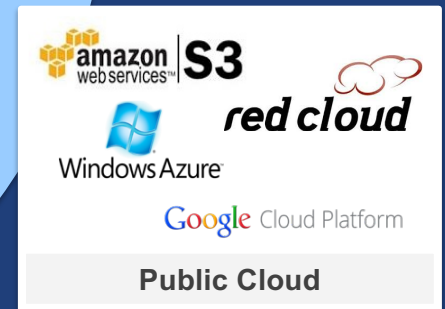
**Research Computing HPC**

**National Resources**

**Personal Resources**

**Desktop Workstations**

**Mass Storage**

**Instruments**

**Public Cloud**

# Globus Connectors

amazon web services™ | **S3**

Google Cloud

**SPECTRA**

IBM Cloud Object Storage

box

Caringo Swarm™

Quantum

ActiveScale Object Storage

iRODS®

ceph

Google Drive

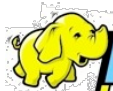IBM **Spectrum Scale**

HPSS

wasabi hot cloud storage

SCALITY

lustre™

OneDrive

hadoop HDFS

**Planned**   Microsoft Azure Blob Storage   Dropbox

# Share with collaborators/community



Project repositories, replication stores

Public repositories

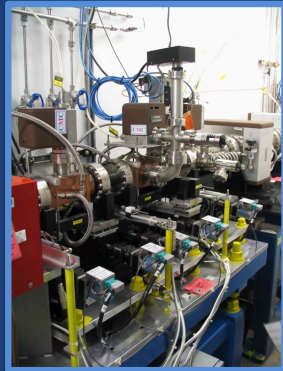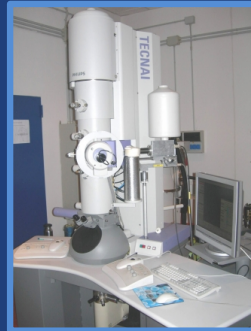External campus storage

Public / private cloud stores

# Manage data from instruments

Next-Gen Sequencer

Advanced Light Source

Cryo-EM

MRI

Light Sheet Microscope

Analysis store

High-durability, low-cost store

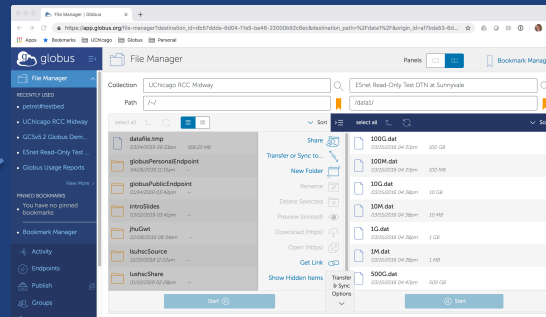Remote visualization

Personal system

# Use(r)-appropriate interfaces



**Web**

**CLI**

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
...
```

**Rest API**

# Globus SaaS / PaaS: 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!

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

**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.

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

**5**

**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**
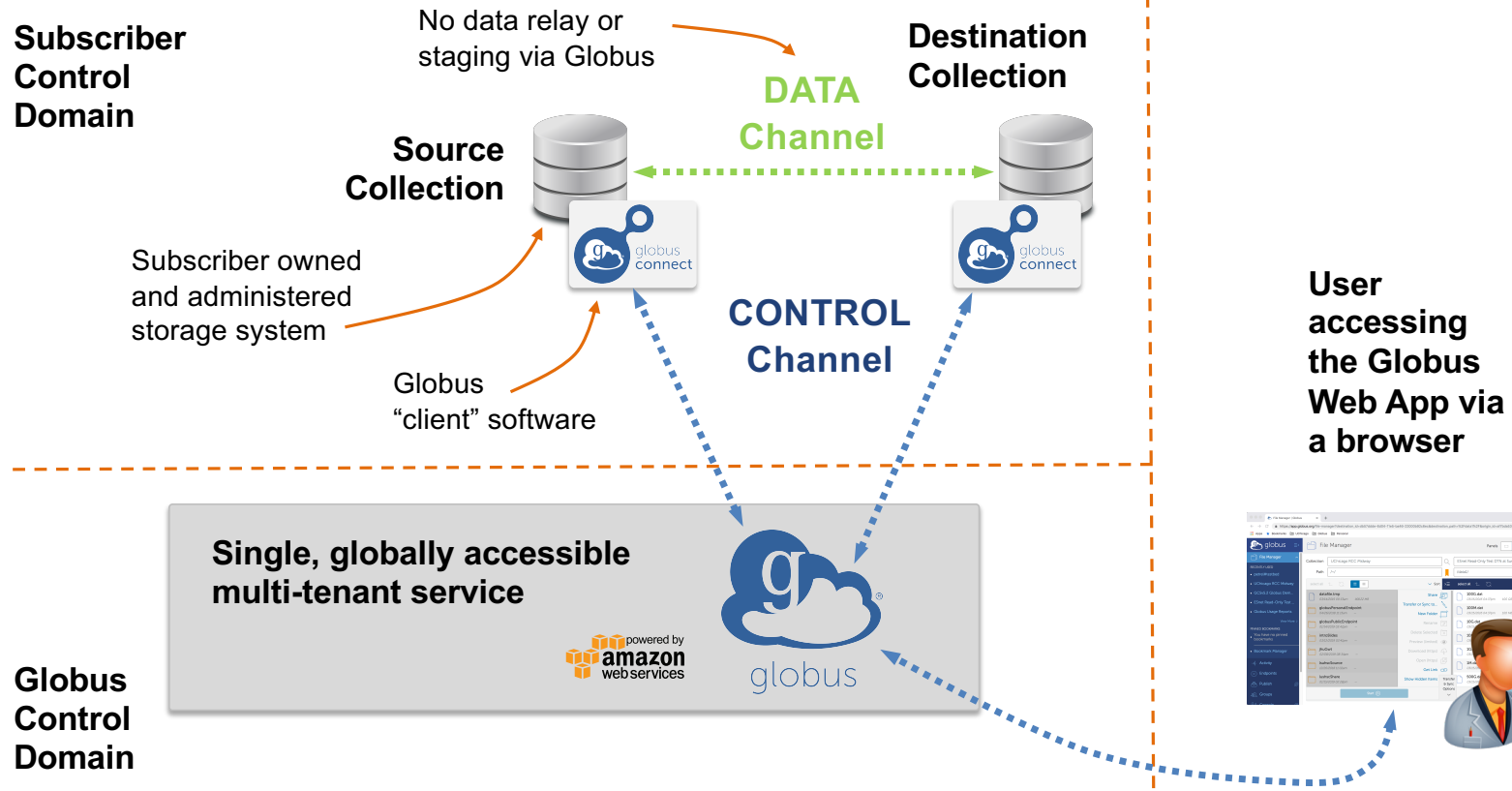- **Use an existing identity**

**Personal Computer**

# Globus core security features

- **Access Control**
  - Identities provided and managed by institution
  - Institution controls all access policies
  - Globus is identity broker; no access to/storage of user credentials
  - Fine grained access control on the collections

- **Data remain at institutions, not stored by Globus**

- **Data does not flow through the Globus Service but directly between Endpoints and their Collections**

- **Integrity checks of transferred data**

- **High availability and redundancy**

- **Encryption of user files and Globus control data**

# Conceptual architecture: Hybrid SaaS

**Subscriber Control Domain**

No data relay or staging via Globus

**Destination Collection**

**DATA Channel**

**Source Collection**

Subscriber owned and administered storage system

Globus "client" software

**CONTROL Channel**

**User accessing the Globus Web App via a browser**

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

powered by amazon web services

globus

**Globus Control Domain**

# Conceptual architecture: Sharing



Subscriber managed filesystem permissions

Globus managed "overlay" permissions

# Endpoints, Collections and Globus Connect

- **Globus Connect Server**
  - Multi user Linux Systems
  - https://docs.globus.org/globus-connect-server/

- **Globus Connect Personal**
  - Personal Workstations and Laptops
  - https://www.globus.org/globus-connect-personal
  - OS specific instructions
    o https://docs.globus.org/how-to/

# Globus Demo

- Authenticating to Globus
  - Linking in other identities

- The hamburger menu – file / folder management

- Transfer to and from ALCF resources
  - theta, cooley, thetagpu, grand
  - University of Chicago RCC Midway

- Globus Connect Personal
  - Installation
  - Transfer

# 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**

# Data centric applications leveraging Globus

# Developer references

- Globus documentation: docs.globus.org
    - Command Line Interface: docs.globus.org/cli/
    - Transfer API : docs.globus.org/api/transfer/
    - SDK: globus-sdk-python.readthedocs.io/en/stable/
- Globus GitHub: github.com/globus/
    - Jupyter Notebooks
        - Stand alone notebooks and hub integrations that walk through much of the functionality of our SDK
        - https://github.com/globus/globus-jupyter-notebooks
    - Automation Examples
        - Shell scripted CLI and Python module examples of common research data management use cases
        - https://github.com/globus/automation-examples

# Support resources

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

- **YouTube channel:** **youtube.com/user/GlobusOnline**

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

- **Mailing Lists**
  - globus.org/mailing-lists

- **Customer engagement team**