# Beyond Data Management with Globus

iRODS UGM 2023
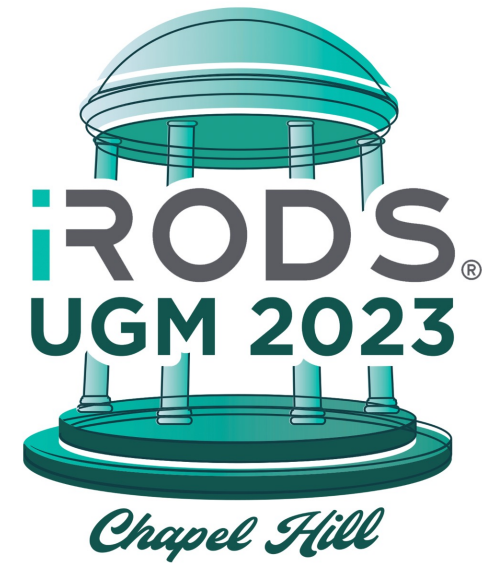Chapel Hill

**Vas Vasiliadis**

University of Chicago – Globus
Adjunct Associate Professor, Masters Program in Computer Science

**vas@uchicago.edu, vasv@anl.gov**

THE UNIVERSITY OF **CHICAGO**

Argonne
NATIONAL LABORATORY

globus

# Reimagining research IT with Globus

**Managed transfer & sync**

**Publication & discovery**

**Platform-as-a-Service**

fast
secure transfer
reliable

DOI 456

your app

**Collaborative data sharing**

collaborators

colleagues

leadership
class
computing

personal
computing

institutional
computing

commercial
computing

research
computing

**Managed remote execution**

**Software-as-a-Service**

local
storage

HPC
systems

institutional
storage

commercial
clouds

tape
archives

**Unified data access**

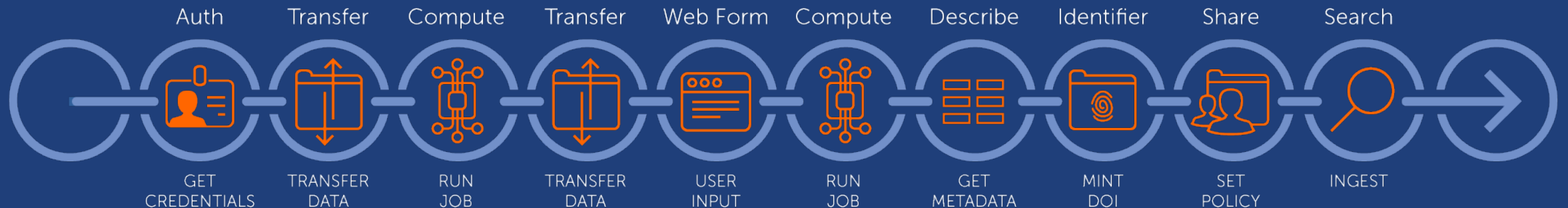Auth        Transfer        Search

GET
CREDENTIALS

TRANSFER
DATA

INGEST

**Reliable automation**

# Automating Research with Globus Flows

- A platform for orchestrating distributed research tasks

- Flows comprise **Actions**

- **Action Providers**: Called by Flows to perform tasks

- **Triggers**: Start flows based on events

- Extensible via **Action Provider API**

# How to best enable distributed, remote compute?

**Borrow page from data management playbook**

➔ **"Fire-and-forget" computation**

➔ **Uniform access interface**

➔ **Federated access control**

➔ **Move closer to researchers' environments**

    ➔ Researchers primarily work in high level languages

    ➔ Functions are a natural unit of computation

**Globus Compute**

**Managed, federated Function-as-a-Service for reliably, scaleably and securely executing functions on remote endpoints from laptops to supercomputers**

THE UNIVERSITY OF CHICAGO

ILLINOIS

Argonne
NATIONAL LABORATORY

# Globus Compute transforms any computing resource into a function serving endpoint

- **Python `pip` installable agent**

- **Elastic resource provisioning from local, cluster, or cloud system (via Parsl)**

- **Parallel execution using local fork or via common schedulers**
  - Slurm, PBS, LSF, Cobalt, K8s

Compute Service

# Executing functions with Globus Compute

- **Users invoke functions as tasks**
  - Register Python function
  - Pass input arguments
  - Select endpoint(s)

- **Service stores tasks in the cloud**

- **Endpoints fetch waiting tasks (when online), run tasks, and return results**

- **Results stored in the cloud and on Globus storage endpoints**

- **Users retrieve results asynchronously**

# User interaction with Globus Compute

**A**

**B**

**2** Globus Compute manages the reliable and secure execution on these endpoints

Compute Service

globus

**1** You request a function be executed on endpoints A and B

**3** Globus Compute returns results or stores them until requested

**Executing a bag of tasks**, e.g., running simulations with different parameters, executing ML inferences, on multiple remote computers directly from your environment, e.g., Jupyter notebook

**Constructing and running automated analysis pipelines with data processing steps that need to be executed in different locations**
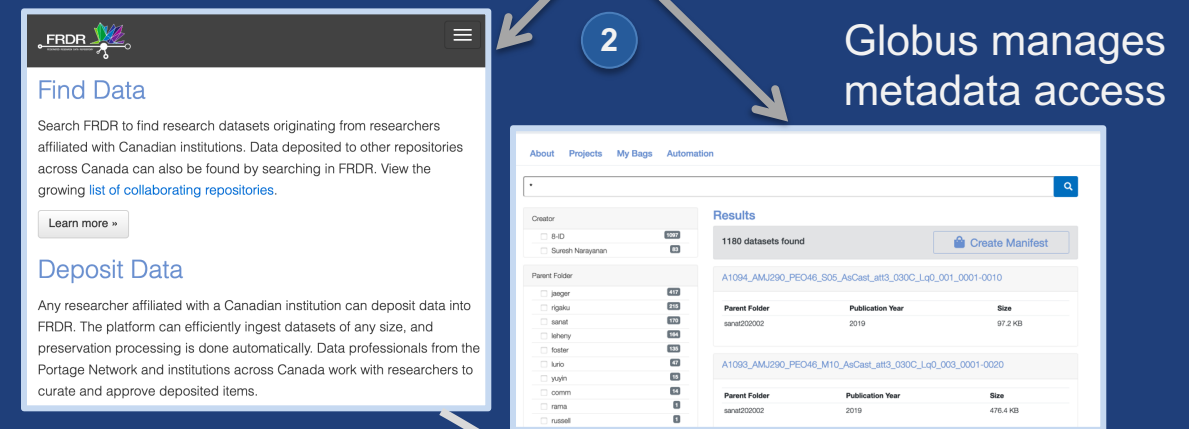
**Building new applications and services** that seamlessly execute application components or user workloads on remote resources

# Scalable data discovery with Globus Search

- **Scalable metadata store**

- **Fine-grained visibility controls**

- **Schema agnostic**
  → **dynamic schemas**

- **Federated auth integration**

- **Robust query API**
  - **GET with URL parameters**
  - **POST with facets**

**docs.globus.org/api/search**



User publishes metadata into search index

Globus manages metadata access

Users query and discover data of interest

# CR3 Portal (simulated data)



Federated logon using Globus Auth with 1,800+ identity providers

Google-like text search with facets for filtering

Variable facets based on source registry index
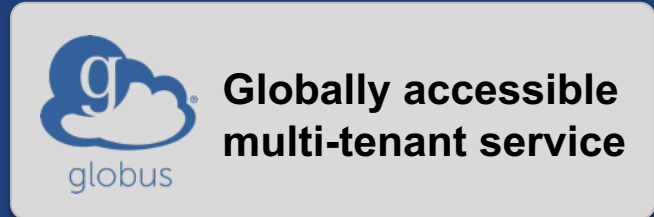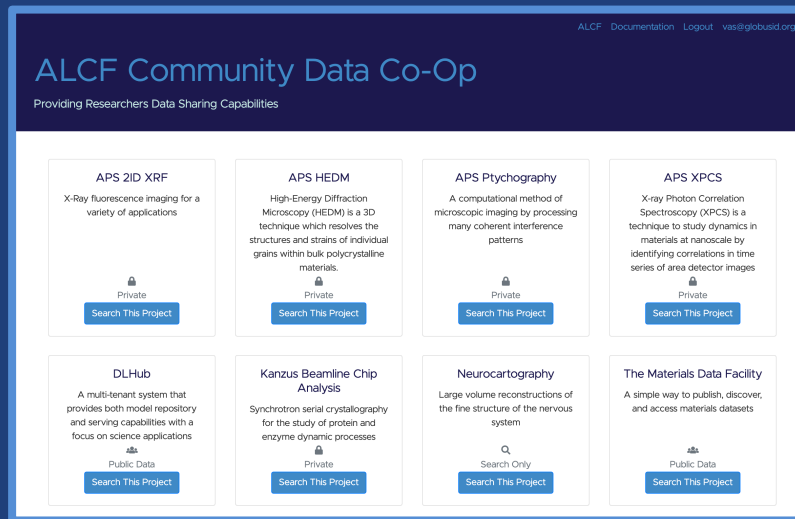
Dynamically updating charts as facets change

Developed using a framework based on the Globus Modern Research Data Portal* design pattern (**docs.globus.org/mrdp**)

* PeerJ Articles:cs-144 https://peerj.com/articles/cs-144/

# Repository data distribution

- **Faceted search via data portal**
- **Enforces fine-grained authZ**
- **HTTPS download for "small" data**
- **Managed file transfer for larger data sets**

Example: **acdc.alcf.anl.gov**



**Globally accessible multi-tenant service**

**ALCF Community Data Co-Op**

Providing Researchers Data Sharing Capabilities

| APS 2ID XRF | APS HEDM | APS Ptychography | APS XPCS |
|---|---|---|---|
| X-Ray fluorescence imaging for a variety of applications | High-Energy Diffraction Microscopy (HEDM) is a 3D technique which resolves the structures and strains of individual grains within bulk polycrystalline materials. | A computational method of microscopic imaging by processing many coherent interference patterns | X-ray Photon Correlation Spectroscopy (XPCS) is a technique to study dynamics in materials at nanoscale by identifying correlations in time series of area detector images |
| Private | Private | Private | Private |
| Search This Project | Search This Project | Search This Project | Search This Project |

| DLHub | Kanzus Beamline Chip Analysis | Neurocartography | The Materials Data Facility |
|---|---|---|---|
| A multi-tenant system that provides both model repository and serving capabilities with a focus on science applications | Synchrotron serial crystallography for the study of protein and enzyme dynamic processes | Large volume reconstructions of the fine structure of the nervous system | A simple way to publish, discover, and access materials datasets |
| Public Data | Private | Search Only | Public Data |
| Search This Project | Search This Project | Search This Project | Search This Project |

Bulk data transfer

Browser based download

Search, request data of interest

# Resources

- **Web app access: app.globus.org**

- **Documentation: docs.globus.org**

- **Helpdesk: support@globus.org**

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