### Science Services and Science Platforms Using the cloud to accelerate and democratize discovery

Talk at NorduGrid annual conference, Košice, Slovakia, June 2<sup>nd</sup>, 2016







### **Thanks to co-authors and Globus team**

Globus services (globus.org)

- Foster, I. Globus Online: Accelerating and democratizing science through cloud-based services. IEEE Internet Computing(May/June):70-73, 2011.
- Chard, K., Tuecke, S. and Foster, I. Efficient and Secure Transfer, Synchronization, and Sharing of Big Data. Cloud Computing, IEEE, 1(3):46-55, 2014.
- Chard, K., Foster, I. and Tuecke, S. Globus Platform-as-a-Service for Collaborative Science Applications. Concurrency - Practice and Experience, 27(2):290-305, 2014.
- Publication (globus.org/data-publication)
- Chard, K., Pruyne, J., Blaiszik, B., Ananthakrishnan, R., Tuecke, S. and Foster, I., Globus Data Publication as a Service: Lowering Barriers to Reproducible Science. 11th IEEE International Conference on eScience Munich, Germany, 2015

**Discovery engines** 

 Foster, I., Ananthakrishnan, R., Blaiszik, B., Chard, K., Osborn, R., Tuecke, S., Wilde, M. and Wozniak, J. Networking materials data: Accelerating discovery at an experimental facility. Big Data and High Performance Computing, 2015.

### Thank you to our sponsors!



### Civilization advances by extending the number of important operations which we can perform without thinking about them

Alfred North Whitehead (1911)

### Computation may someday be organized as a public utility ... The computing utility could become the basis for a new and important industry.



### John McCarthy (1961)



### The grid vision

Accelerate discovery and innovation by providing on-demand access to computing

- "the average computing environment remains inadequate for [many] computationally sophisticated purposes"
- "if mechanisms are in place to allow reliable, transparent, and instantaneous access to high-end resources, then it is as if those resources are devoted to them" [*The Grid*, Chapter 2, 1998]

### **Another pioneer: NorduGrid**



The NorduGrid project: using Globus toolkit for building GRID infrastructure

#### M. Ellert<sup>a</sup>, A. Konstantinov<sup>b,\*</sup>, B. Kónya<sup>c</sup>, O. Smirnova<sup>c</sup>, A. Wäänänen<sup>d</sup>

<sup>a</sup> Department of Radiation Sciences, Uppsala University, Box 535, 751 21 Uppsala, Sweden <sup>b</sup> University of Oslo, Department of Physics, P.O. Box 1048, Blindern, 0316 Oslo, Norway <sup>c</sup> Elementarpartikelfysik, Fysiska Institutionen, Lunds Universitet., Box 118, 22100 Lund, Sweden <sup>d</sup> Niels Bohr Institutet for Astronomi, Fysik og Geofysik., Blegdamsvej 17, Dk-2100 Copenhagen Ø, Denmark

### **Example grid scenarios**

- "The application service providers, storage service providers, cycle providers, and consultants engaged by a car manufacturer to perform scenario evaluation during planning for a new factory"
- "Members of an industrial consortium bidding on a new aircraft"
- "A crisis management team and the databases and simulation systems that they use to plan a response to an emergency situation"
- "Members of a large, international, multiyear highenergy physics collaboration"

From: The Anatomy of the Grid, 2001

### Higgs discovery "only possible because of the extraordinary achievements of ... grid computing"—Rolf Heuer, CERN DG

10s of PB, 100s of institutions, 1000s of scientists, 100Ks of CPUs, Bs of tasks

### What has changed?

- Thousands of people learned about the joys of large-scale distributed systems
- Virtual organization concepts and technologies
- Now routine to move 100s of terabytes (e.g., GridFTP moves >2 petabyte per day)
- High throughput computing is mainstream (e.g., Condor and Globus run millions of jobs per day)
- Large Hadron Collider found the Higgs
- Earth System Grid supports >25,000 users
- Commercial cloud computing has exploded

### **Looking forward**

- Exploding data volumes and earlier successes mean that many more people face challenges of big data, big compute, big collaboration
- Networks are several orders of magnitude faster than when Grid started
- Commercial cloud providers provide a substrate on which powerful new capabilities can be built with new economies of scale

The BVP Cloudscape Top 300 Privately Held Cloud Companies

#### **Business Users**



## Cloud has transformed how software is developed and delivered





PaaS enables more rapid, cheap, and scalable delivery of powerful apps—as SaaS

# The right platform can do the same for science



We can leverage cloud to provide solutions that span the vast majority of researcher needs 14

# The right platform can do the same for science



We can leverage cloud to provide solutions that span the vast majority of researcher needs 15

# A science platform can spur a discovery cloud ecosystem



In so doing, we can slash costs, improve quality, and accelerate discovery across the sciences <sup>16</sup>

# A science platform can spur a discovery cloud ecosystem



In so doing, we can slash costs, improve quality, and accelerate discovery across the sciences <sup>17</sup>

# What can we automate and outsource in science?

Run experiment Collect data Move data Check data Annotate data Share data Find similar data Link to literature Analyze data Publish data



Automate and outsource

Discovery Cloud



### Identity and authorization challenges

DOC DB Access KERBEROS reamine wiki = NESA wiki - @ FNAL LOAP NICSA ORACLE REAPACESS JERVICE NOW @ CLOUD SVN checkout key 1"VO" MECHANSM JIRA PIOBLEN TICKET LISTSORN MAILING LIST (6) + ARCHIVE shoa-show RSA Xeys (7.) NCSA MAILING LISTE) + ARCHIVE "GROUP ACCTS FINEL ANER CLUSTER, BRAZIL ANALYSIS PORTAL, SLAC CLUSTER, BNAL CLUSTER, .... FUA GRID FTP CENTIFICATE PLODUTION: TERRAGELO MOLISS WRITE ALLESS ORACLE

# Data access: we have the highways but not the delivery service





Our highways encompass the Internet, ultra-high-speed networks, science DMZs, data transfer nodes, high-speed transport protocols A good **delivery service** automates, schedules, accelerates, adapts. It provides APIs for experts and casual users. Cuts costs and saves time.

### Globus: Research data management as a service



### 151,487,208,573 MB

## Essential research data management services

- File transfer
- Data sharing
- Data publication
- Identity and groups

Builds on 15 years of research

#### Outsourced and automated

- High availability, reliability, performance, scalability
- Convenient for
  - Casual users: Web interfaces
  - Power users: APIs
  - Administrators: Install, manage

### **Globus and the research data lifecycle**



### **Globus by the numbers**

major services

160 PB transferred

## **25 billion** files processed

## **38,000** registered users

13 national labs use Globus

35+

institutional

subscribers

**10,000** active endpoints

**1 PB** 

largest single

transfer to date

~400

active daily users

**99.9%** 

uptime

3 months

longest continuously managed transfer 130 federated campus identities

# Platforms can slash costs, simplify access, increase interoperability

- For example, by providing:
- Federated identity system with finegrained authorization
- Data management easily integrated with application workflows
   Via RESTful APIs





**Domain-independent and domain-specific services** 

AWS, Google, Azure services

Foundation services: Globus Auth, Groups, etc.

### **Globus PaaS: Ecosystem** enabler



### Globus PaaS and Open Science Grid

| Sho | w the | bookr | marks in this fo  | older 💽 h | ttps 📾 🛛 por | tal.osgconr | nect.net/ | SignIn          |      |       |       |     |         |         |      |      | Ċ       | Reader | C   |   |
|-----|-------|-------|-------------------|-----------|--------------|-------------|-----------|-----------------|------|-------|-------|-----|---------|---------|------|------|---------|--------|-----|---|
| 60  |       |       | Contacts <b>*</b> | Bonjour ▼ | MWT2 🔻       | ATLAS 🔻     | OSG ₹     | Coding <b>v</b> | IB ▼ | OSX ₹ | Mac ▼ | РМ▼ | FENIX 🔻 | Media 🔻 | GO ≖ | Biz▼ | Props v | X      | » ʃ | ÷ |



Support - Resources - O

OSG Connect -

Sign In / Register -

#### Efficiently connect your science to cycles and data



OSG Connect offers users simple access to distributed high throughput computing resources, and reliable, high-performance file transfer services.

| Sign In   | Sign Op with Globus Online |
|---|----------------------------|
| Using your InCommon / CILogon login                   | . alternate login          |
|   |                            |
| You will now be redirected to In authentication page. | Common / CILogon's         |

### Simple web app server login



KBase apps are ready-to-use workflows consisting of a set of chained methods that together perform some useful analysis.

ropagate Genome-scale Model to

Reconstruct Community Metabolic

Close Genome

Model

KBase is an open platform for comparative functional genomics and systems biology for microbes, plants and their communities, and for sharing results and methods with other scientists.

Add one or more genomes to the KBase species tree. more...

+ add another Genome

Mycoplasma capricolum su...

Genome to

species tree

Sign In

Genome

### **Jetstream cloud service**

| • • • Atmosphere ×  | Steve          |
|---|----------------|
| ← → C Attps://use.jetstream-cloud.org/application/images  | 🗶 🖓 🏠 🔒 🛈 🐵 ≡  |
| Jetstream DEMO  |                |
| H Images O Help   | Login          |
| Q SEARCH STAGS  |                |
| Search across image name, tag or description              |                |
| Showing 5 of 5 images                                     | EList III Grid |
| Featured Images   |                |
| All Images  |                |
| CentOS 7 Stock 1601<br>Feb 1st 2016 08:31 am CST by admin |                |
| Imported Application - CentOS 7 Stock 1601                |                |
| CentOS 6 Stock 1601                                       |                |
| ©2016 Jetstream C Feedback & Support                      |                |

### Serving a global community: NCAR's Research Data Archive



### Serving a global community

- 17+ PB virtual processing
- 45,000+ custom orders, 4,000 users, 380 TB served in 2014



Courtesy of Thomas Joram, NCAR (2014)

### Fully automated delivery using portal developed w/ PaaS

# PaaS enabled automated workflow

User logs in with NCAR or other campus identity

Selected dataset copied to staging area (shared endpoint)



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  | NCAR RDA Email/Username | tcram@ucar.edu | 権     |
|---|-------------------------|----------------|-------|
| Name: Globus Online<br>URL: <u>https://www.globusonline.org</u> | NCAR RDA Password       | •••••          | 楕     |
|   |                         | Sign In C      | ancel |

User granted read permission for shared endpoint User receives email with link to access files ACLs deleted after five days

### **Sanger Imputation Service**

This is a free genotype **imputation** and **phasing** service provided by the Wellcome Trust Sanger Institute. You can upload GWAS data in VCF or 23andMe format and receive imputed and phased genomes back. Click here to learn more and follow us on Twitter.

#### Before you start

Be sure to read through the instructions.

You will need to set up a free account with Globus and have Globus Connect running at your institute or on your computer to transfer files to and from the service.

#### Ready to start?

If you are ready to upload your data, please fill in the details below to **register an imputation and/or phasing job**. If you need more information, see the about page.

Full name

Organisation

Email address

#### What is this 🕑

Globus user identity



#### News



#### 11/05/2016

Thanks to EAGLE, we can now return **phased data**. The HRC panel has been updated to r1.1 to fix a known issue. See ChangeLog for more details.

#### 15/02/2016

Globus API changed, please see updated instructions.

#### 17/12/2015

New status page and reworked internals. See ChangeLog.

#### 09/11/2015

Pipeline updated to add some features requested by users. See ChangeLog.

See older news...

### **Research data portal pattern**



Move portal storage into Science DMZ, with Globus endpoint

34

- Leave Portal Web server behind firewall
- Globus handles the security and data heavy lifting

### **Research data portal pattern**



Move portal storage into Science DMZ, with Globus endpoint

35

- Leave Portal Web server behind firewall
- Globus handles the security and data heavy lifting

### https://github.com/globus/globus-sample-data-portal



#### **Globus Transfer API**

API reference for transfer and sharing functions.

#### **Globus Auth API**

API reference for authentication and authorization.

#### **Frequently Asked Questions**

When all else fails ...

Developer Workshop: Building the Modern Research Data Portal

CHICAGO

globusworld

New high-speed networks make it possible, in principle, to transfer and share research data at tremendous speeds and scales-but have also proved challenging to use in practice. Two new technologies now allow us to translate this potential into reality: Science DMZ architectures provide frictionless end-to-end network paths; and Globus APIs allow programmers to provide function. Con

research data portals that leverage these paths for data distribu synchronization, and other useful purposes.

Let us know if you'd like to participate in future workshops Introduction, Concepts, and Components IMPERIAL 2 Led by: TBD

We will introduce the Modern Research Data Portal and set the context for how Globus and the ScienceDMZ combine to deliver unique data management capabilities. This will include:

- Overview of use cases: Common patterns like data publication/distribution, orchestration of data flows, etc.
- Overview of the Globus platform: Architecture and brief overview of available services
- Introduction to the Globus Auth API: Authenticating and authorizing
   a client
- Introduction to the Globus Transfer API: Make your first call and move data with Globus
- Introduction to the Python SDK for using Globus Auth and Transfer

### **Research data portal pattern**



Move portal storage into Science DMZ, with Globus endpoint

38

- Leave Portal Web server behind firewall
- Globus handles the security and data heavy lifting

### **Globus Auth**

- Foundational identity and access management (IAM) platform service
- Simplify creation and integration of advanced apps and services
- Brokers authentication and authorization interactions between:
  - End-users
  - Identity providers: XSEDE, InCommon, web apps
  - Resource servers: services with REST APIs
  - Clients: web, mobile, desktop, command line apps
  - Resource servers acting as clients to other resource servers

### https://docs.globus.org/api/auth

# Based on widely used web standards

- OAuth 2.0 Authorization Framework
  - aka OAuth2
- OpenID Connect Core 1.0

aka OIDC

- Allows use of standard OAuth2 and OIDC libraries
  - E.g., Google OAuth Client Libraries (Java, Python, etc.), Apache mod\_auth\_openidc

### **Globus Auth uses**

- Login to web app
  - "Log in with Globus"
  - Mobile, desktop, command line apps coming
- Protect all REST API communications
  - App → Globus service
  - App → non-Globus service
  - Service  $\rightarrow$  service

### **Research data portal pattern**



Move portal storage into Science DMZ, with Globus endpoint

42

- Leave Portal Web server behind firewall
- Globus handles the security and data heavy lifting

### **Globus transfer API**

Nearly all Globus Web App functionality implemented via public Transfer API

### https://docs.globus.org/api/transfer/

HOME / GLOBUS APIS

### Transfer API Documentation

This API provides a REST-style interface to the Globus reliable file transfer service. The Transfer API supports monitoring the progress of a user's file transfer tasks, managing file transfer endpoints, listing remote directories, and submitting new transfer and delete tasks. The API is ideal for integration into

Transfer API Documentation

**API** Overview

### **Globus Python SDK**

## Python client library for the Globus Auth and Transfer REST APIs

http://globus.github.io/globus-sdk-python/

| obus-sdk-python 0.2.3 documentati  | ion » next   modules   index  |
|--|---|
| able Of Contents   | Globus SDK for Python (Beta)  |
| Globus SDK for Python (Beta)<br>Installation<br>Basic Usage<br>API Documentation | This SDK provides a convenient Pythonic interface to Globus REST APIs, including the Transfer API and the Globus Auth API. Documentation for the REST APIs is available at https://docs.globus.org. |
| License<br>ext topic   | Two interfaces are provided - a low level interface, supporting only GET, PUT, POST, and DELETE operations, and a high level interface providing helper methods for common API resources.           |
| High Level API   | Source code is available at https://github.com/globus/globus-sdk-python.  |
| his Page<br>Show Source  | Installation  |
| uick search  | The Globus SDK requires Python 2.6+ or 3.2+. If a supported version of Python is not already installed on your system, see this Python installation guide.  |
| Enter search terms or a module, class or function name.                          | The simplest way to install the Globus SDK is using the pip package manager (https://pypi.python.org/pypi/pip), which is included in most Python installations:                                     |
|  |   |

pip install globus-sdk

### Jupyter (iPython) notebooks

### https://github.com/globus/globus-jupyter-notebooks

#### **Globus SDK**

https://github.com/globus/globus-sdk-python

#### **Globus SDK Docs**

http://globus.github.io/globus-sdk-python/

#### **Requirements**

- You need to be in the tutorial users group for sharing: https://www.globus.org/app/groups/50b6a29c-63ac-11e4-8062-22000ab68755
- Installed Globus Python SDK

```
In [15]: from __future__ import print_function # for python 2
tutorial_endpoint_1 = "ddb59aef-6d04-11e5-ba46-22000b92c6ec" # endpoint "Globus Tutorial Endpoint 1"
tutorial_endpoint_2 = "ddb59af0-6d04-11e5-ba46-22000b92c6ec" # endpoint "Globus Tutorial Endpoint 2"
tutorial_users_group = "50b6a29c-63ac-11e4-8062-22000ab68755" # group "Tutorial Users"
```

#### Configuration

First you will need to configure the client with an OAuth2 access token. For the purpose of this tutorial, you can obtain access tokens via the tokens.globus.org website. Click the "Jupyter Notebook" option and copy the resulting text below, or click on "Globus CLI" and copy the resulting text into ~/.globus.cfg.

In [16]: transfer\_token = None # if None, tries to get token from ~/.globus.cfg file

45

### **Research data portal pattern**



Move portal storage into Science DMZ, with Globus endpoint

46

- Leave Portal Web server behind firewall
- Globus handles the security and data heavy lifting

### **Globus helper pages**

**Globus-provided** web pages designed for use by your web apps gm

- Browse Endpoint
- Select Group
- Logout



Account

Manage Data - Publish

Groups

### https://docs.globus.org/api/helper-pages/

### **Globus helper pages**

| 🕑 globi   | JS                                | l                | Manage D      | ata -    | Publi   | ish           | Groups          | Support -   | tuecke -      |
|---|-----------------------------------|------------------|---------------|----------|---------|---------------|-----------------|---|---------------|
|   |                                   | Browse 8         | Discover      | Data     | Publica | ation [       | Dashboard       | Communities                                       | & Collections |
|   | that datasets you sub<br>will be  | mit to this tria | l collection: | (1)      |         | Input         | Form*           | Datacite Man                                      | datory + R 🜲  |
| Jobus globus  | Manage Data                       | Publish Group    | s Support -   | tuecke - |         | Subm<br>Workf | ission<br>flow* | Default   | *             |
| Select Group  |                                   |                  |               |          |         | Curati        | ion Type*       |   |               |
| Globus Team   |                                   |                  | ✓ details     | ٩        |         | Curut         |                 | Edit Metadata                                     | •             |
| Globus Team - Development                                 |                                   |                  |               |          |         |               |                 |   |               |
| Globus Team - User Services                               |                                   |                  | ✓ details     |          |         |               |                 |   |               |
| Globus Team Plus Sponsor                                  |                                   |                  | ✓ details     |          |         | Collec        | tion Permiss    | ions  |               |
| Globus Team     Globus Publication Users                  |                                   |                  | details       |          |         |               |                 |   |               |
|   |                                   |                  | ✓ details     |          |         |               |                 |   |               |
| Wellcome Trust Sanger Institute Scien                     | tific Users                       |                  | ✓ details     |          |         | Subm          | itters          | <ul> <li>All Users</li> <li>Restricted</li> </ul> | to Group      |
| Submit  |                                   |                  |               |          | ~       | Acces         | s to Data       |   |               |
| © 2010-2016 Computation Institute, University of Chicago, | Argonne National Laboratory legal |                  |               |          |         |               |                 | <ul> <li>All Users</li> <li>Restricted</li> </ul> | to Group      |
|   |                                   |                  |               |          |         |               |                 |   |               |
|   |                                   |                  |               |          |         | Curati        | ion Group       | null<br>Char                                      | ige           |



### Can skin Globus Auth pages



### **Research data portal pattern**



Move portal storage into Science DMZ, with Globus endpoint

50

- Leave Portal Web server behind firewall
- Globus handles the security and data heavy lifting

### **Globus Connect HTTPS**

- The future of research CI is ... the web
- Globus Connect HTTPS unlocks all research storage to the web
- Globus Auth provides security glue using standard web security
- GridFTP doesn't go away async, bulk data transfer is important, but its not the end-all, be-all

### **Research data portal pattern**



Move portal storage into Science DMZ, with Globus endpoint

52

- Leave Portal Web server behind firewall
- Globus handles the security and data heavy lifting

### Why create your own services?

- Front-end / back-end within your portal
  - Remote backend for portal
  - Backend for pure Javascript browser apps

 Extend your portal with a public REST API, so that other app and service developers can integrate with and extend your portal

# Why Globus Auth for your service?

- Outsource all identity management, authentication
  - Federated identity with InCommon, Google, etc.
- Outsource your REST API security
  - Consent, token issuance, validation, revocation
  - You provide service-specific authorization
- Apps use your service like all others
  - Its standard OAuth2 and OIDC
- Your service can seamlessly leverage other services
- Other services can leverage your service

Add your service to the international science platform

# A science platform can spur a discovery cloud ecosystem



In so doing, we can slash costs, improve quality, and accelerate discovery across the sciences <sup>55</sup>

### **Enabling the Discovery Cloud**



Together we can create an integrated ecosystem of services and applications for the research and education community

Thank you! foster@uchicago.edu @ianfoster