# Globus 2018:
# **Beyond File Transfer**

**GlobusWorld 2018 Keynote**

**Steve Tuecke & Ian Foster**

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

400,075,759,997 MB TRANSFERRED

# Globus by the numbers

**1,042** most shared endpoints at a single institution

**400 PB** transferred

**66 billion** files processed

**100,000** users

**24** Petabyte+ institutions

**15,000** active transfer users

**3 months** longest running transfer

**20,000** active endpoints

**500+** identity providers

**1 PB** largest single transfer to date

**8,000** active shared endpoints

**99.9%+** availability

# Globus Auth adoption

- **100,000 users, 37k new in last year**

- **500+ identity providers**

- **1,100 registered apps and services**

- **103,000 user consents, 35% non-Globus apps**

- **99.994% uptime since Feb 2016**

# Progress on sustainability

- **90 subscribers, including 1/3 of R1 universities**
- **½ of our product funding is from subscriptions**
- **Need most R1 and many R2 universities**
- **You can help by encouraging others to subscribe**

# Help us get the word out!

- **Do you rely on Globus for your work?**

- **If so, please share your experiences!**
  - **Contribute** to our Usage Brief Library
    globus.org/usage-brief-library
  - **Add a slide** or logo in event talks (we can help!)
  - **Mention Globus** in news articles or interviews
  - **Tag us** in posts about projects that use Globus
  - **Acknowledge Globus** in your journal articles
    globus.org/publications

- **Why?**
  - Give your peers new ideas on how to use Globus
  - Help us grow the user community

"…, and file sharing with Globus."

"…with Globus for data transfer."

"We used Globus for…"

"I needed Globus to…"

"...and Globus."

"#ALCF #ORNL #theNCI #CANDLE #globusonline"

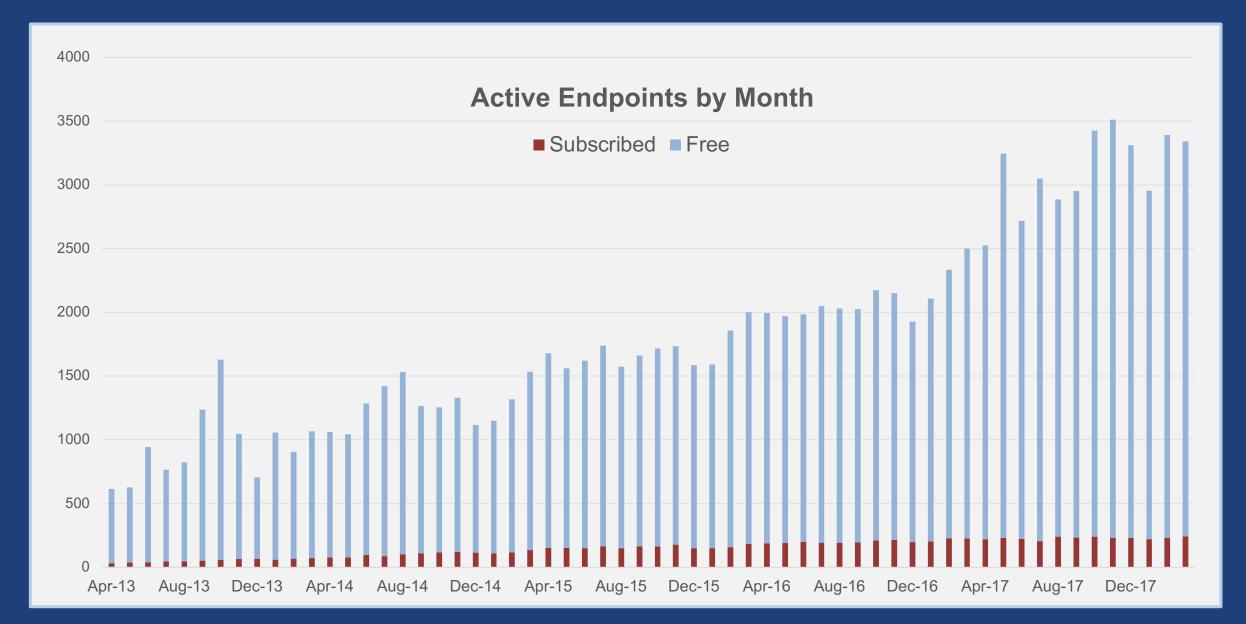"…using tool x, tool y, Globus, technology z…"

# Why subscribe? Go (way) beyond file transfer…

- **Remove friction for external collaborators**

- **Automate/scale research data flows**

- **Diversify research storage options—with a unified interface**

- **Gain visibility into research storage utilization**

- **Integrate robust data management into research apps**

- **Optimize data transfer performance**

- **Access expert support resources**

# Five-year Subscription Growth



Active Endpoints by Month

■ Subscribed  ■ Free

# THANK YOU, subscribers!

# Thank you to our sponsors

U.S. DEPARTMENT OF
**ENERGY**

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

**NSF**

NATIONAL INSTITUTES OF HEALTH

THE UNIVERSITY OF
CHICAGO

ALFRED P. SLOAN FOUNDATION
1934

**Argonne**
NATIONAL LABORATORY

powered by
**amazon**
web services

# New Features and Enhancements

# Transfer performance improvements



**August 2017**
**L380 Data Set**

alcf#dtn_mira ALCF

21.3 Gbps / 24.4 Gbps

19.2 Gbps / 21.1 Gbps / 16.4 Gbps

nersc#dtn NERSC — 17.9 Gbps — 18.4 Gbps — olcf#dtn_atlas OLCF

19.5 Gbps / 16.6 Gbps / 21.8 Gbps

9.5 Gbps / 14.3 Gbps

ncsa#BlueWaters NCSA

**Petascale DTN Project**

**November 2017**
**L380 Data Set**

**Gigabits per second (min/avg/max), three transfers**

alcf#dtn_mira ALCF

33.0/35.0/37.8 Gbps   44.1/46.8/48.4 Gbps

41.0/42.2/43.9 Gbps   43.0/50.0/56.3 Gbps   34.6/47.5/56.8 Gbps

nersc#dtn NERSC   35.9/39.0/40.7 Gbps   29.9/33.1/35.5 Gbps   olcf#dtn_atlas OLCF

23.1/33.7/39.7 Gbps   33.2/43.4/50.3 Gbps

55.4/56.7/57.4 Gbps   26.7/34.7/39.9 Gbps

21.2/22.6/24.5 Gbps

ncsa#BlueWaters NCSA

```
Data set: L380
Files: 19260
Directories: 211
Other files: 0
Total bytes: 4442781786482 (4.4T bytes)
Smallest file: 0 bytes (0 bytes)
Largest file: 11313896248 bytes (11G bytes)
Size distribution:
    1 - 10 bytes: 7 files
    10 - 100 bytes: 1 files
    100 - 1K bytes: 59 files
    1K - 10K bytes: 3170 files
    10K - 100K bytes: 1560 files
    100K - 1M bytes: 2817 files
    1M - 10M bytes: 3901 files
    10M - 100M bytes: 3800 files
    100M - 1G bytes: 2295 files
    1G - 10G bytes: 1647 files
    10G - 100G bytes: 3 files
```

**2x**



ncsa#BlueWaters-TO-nersc#dtn

- old transfer with c=64, p=4, pp=10
- new transfer with c=64, p=4, pp=10



L380 alcf-TO-ncsa stat

Maximum — Minimum

Throughput (Gbps) vs Test case: prod, sandbox-default, 1k, 10k, 20k

# Storage connectors - globus.org/connectors

# HGST ActiveScale

**Western Digital®**

- **Turnkey on-premise object storage**

- **Globus connector using S3 API**

- **Low TCO:**
  - Manufactures own drives
  - Erasure coding
  - BitDynamics: background data integrity checks with self-healing
  - Cloud-based systems management tools
  - Data Forever: automatic migration to new tech

https://docs.globus.org/premium-storage-connectors/wd-activescale/

# Connectors for S3 "compatible" systems

- **S3 API is de-facto standard API for object storage**

- **Make it easier to validate and support connectors for S3 "compatible" object storage systems**
  - Functionality and performance test suite
  - Improving connector robustness and performance
  - E.g., Ceph, ActiveScale, SwiftStack, Wasabi, IBM Cloud Object Storage System (CleverSafe)

- **Also requires vendor engagement and market interest**

# Upcoming Webinar: May 22

**Simplifying large-scale data management and lowering total cost of storage with Globus and Spectra**

- May 22, 2018 at 11 a.m. EDT / 8 a.m. PDT
- Guest speaker from UMN / MSI
- Topics include:
  - MSI's use of the Spectra® BlackPearl® solution with Globus premium connector
  - Cost model for Spectra® BlackPearl®

**https://globus.org/events/webinar-tco-spectra-msi**

*"Spectra Logic's T950 and BlackPearl are important components in our strategic and comprehensive storage plan for hundreds of terabytes of critical research data."*

*--Jeff McDonald, Assistant Director for HPC Operations, MSI*

# HPSS

- **Community has agreed on sustainability model**

- **NERSC & ORNL investing in enhancements**

- **Premium storage connector subscription**

# Globus Connect Server v5 update

- **Limited production releases:**
  - V5.0: Google Drive (Fall 2017)
  - V5.1: HTTPS and Posix shared endpoints (next week)
  - Subsequent v5.x building up to full functionality

- **Year-end target for full release**
  - End-to-end Globus Auth
  - Multi-DTN
  - Remaining connectors

- Globus Connect Community Source License

COLLECTIONS

MAPPED   GUEST   MAPPED   MAPPED   MAPPED   GUEST   GUEST   GUEST

Data access
interface
GridFTP & HTTPS

POSIX 1
STORAGE GATEWAY

mapped & guest
collections
/project

POSIX 2
STORAGE GATEWAY

only mapped
collections
/scratch

GOOGLE DRIVE
STORAGE GATEWAY

only @domain.edu

STORAGE
GATEWAYS

Policies &
configuration

ENDPOINT

ENDPOINT
Management &
config interface

POSIX
CONNECTOR

globus
connect
server
NODE 1

Google
Drive
CONNECTOR

DATA
TRANSFER
NODE 1

POSIX
CONNECTOR

globus
connect
server
NODE 2

Google
Drive
CONNECTOR

DATA
TRANSFER
NODE 2

DATA
TRANSFER
NODES

Network & storage
connected servers
in ScienceDMZ

# Globus Toolkit end of support

- **General support for open source Globus Toolkit has ended**
  - Does not effect Globus service or Globus Connect

- **Customers using the GT GridFTP with Globus service will continue to be supported until GCSv4 is discontinued**
  - Security patches continue for GridFTP, MyProxy, GSI-OpenSSH

- **GCSv4 and GT GridFTP will be discontinued 6 months after GCSv5 full release**

https://github.com/globus/globus-toolkit/blob/globus_6_branch/support-changes.md

# Command Line Interface

- **New Globus CLI is generally available**
  - Fully functional
  - Many enhancements
  - Simple updater

- **Deprecating old hosted SSH CLI**
  - Will be turned off August 1

https://docs.globus.org/cli

```
$ 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
  --jmespath, --jq TEXT   A JMESPath expression to apply to json output.
                          Takes precedence over any specified '--format' and
                          forces the format to be json processed by this
                          expression
  --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.
```
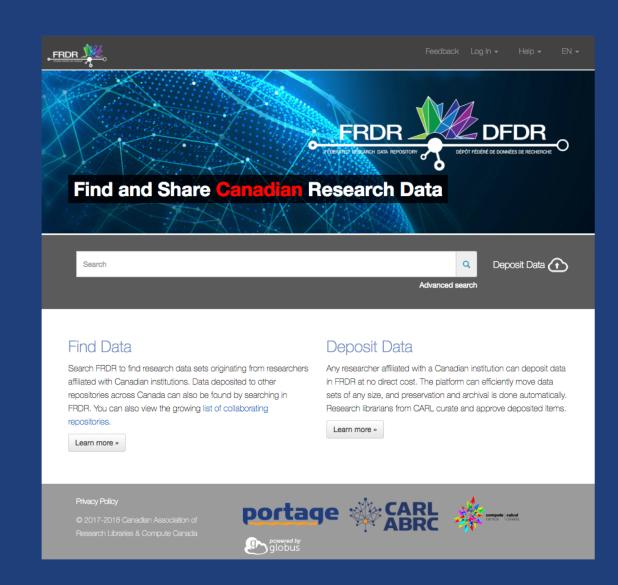
# Symbolic links delayed

- **On transfer and sync**

- **Options:**
  - Ignore
  - Keep as symlinks
  - Copy as files

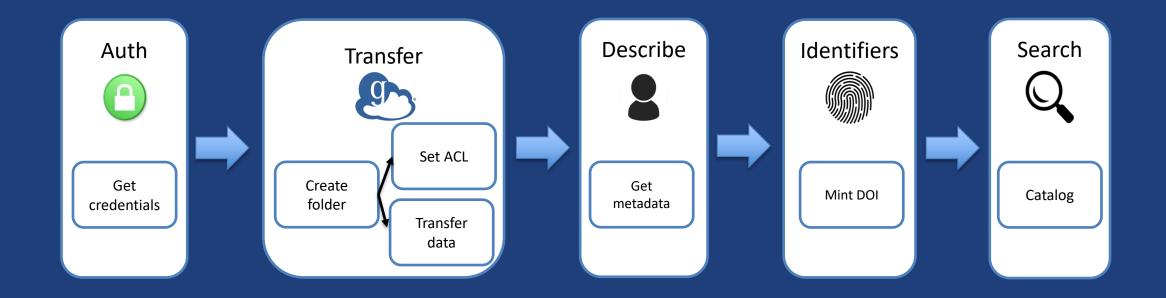- **Target later this year**

# Publication v1

- **Publication v1 app**
  - Publish datasets to Globus Search
  - Internationalization

- **Canadian Federated Research Data Repository**
  - https://frdr.ca/
  - Uses v1 open source and Globus Search

# Publication v2 platform

- **Decompose Publication v1 into platform components**
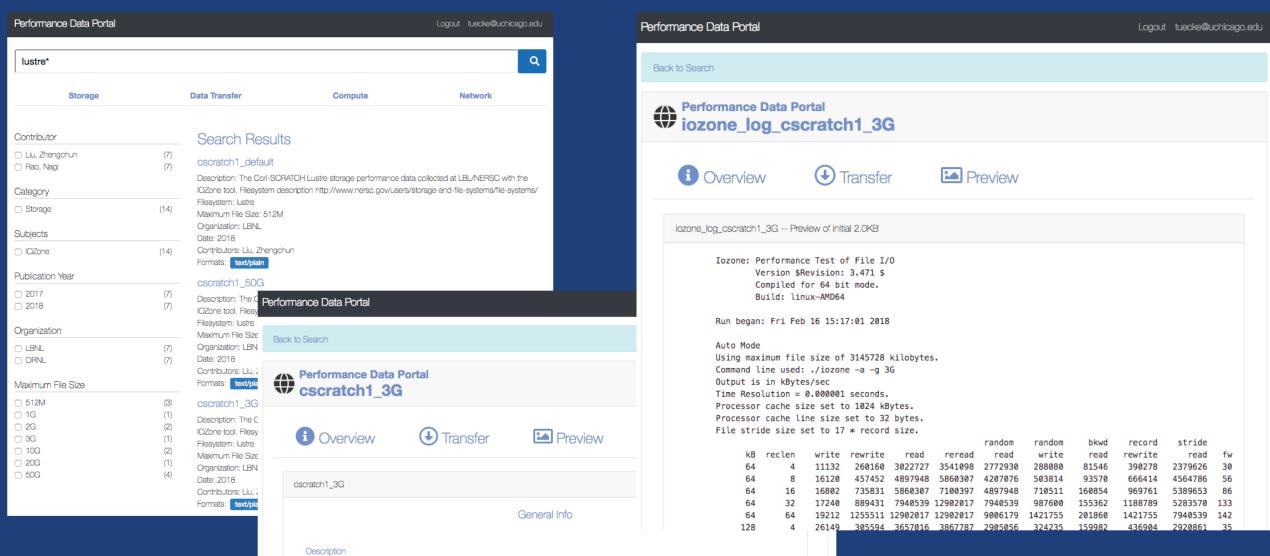- **Allow flexible re-composition & adaptation by customers**

# Globus Search platform service

- **Search service:**
  - **Scalable**: to billions of entries
  - **Schema agnostic**: can use standard (e.g., DataCite) or custom metadata
  - **Fine grain access control**: only returns results that are visible to user
  - **Plain text search**: ranked results
  - **Faceted search**: for data discovery
  - **Rich query language**: ranges, expressions, regex, fuzzy, stemming, etc.
- **Limited production, generally available target year end**
- **Tutorial: Data Publication and Discovery with Globus**

# Django Globus Portal App

Django Globus Portal App Demo

# Globus Identifiers platform service

- **Issue persistent identifiers**
  - DOI, ARK, Handle, Globus
  - E.g., https://identifiers.globus.org/doi:10.1145/2076450.2076468

- **Within a namespace**
  - E.g., Your University's DataCite namespace
  - Control which identities and groups can create identifiers in your namespace

- **Each identifier has:**
  - **Link to data**: one or more https URLs, to file, folder or manifest
  - **Landing page**: provided by service, or by user
  - **Visibility**: which identities and groups can see identifier
  - **Checksum**: of the file or manifest
  - **Metadata**: as required by identifier (e.g., DataCite), extensible
  - **Replaces / Replaced-by**: for versioning

- **Limited beta available now, generally available year end**

- **Tutorial: Data Publication and Discovery with Globus**

30

# Jupyter Integration



- **Authenticate to JupyterHub with Globus Auth**
  - Passes tokens into notebooks as environment variable

- **Use Globus data management platform from notebooks**
  - With Globus Python SDK

Jupyter Integration Demo

What's Coming Next

# Protected data

- **High assurance endpoints**
  - User must authenticate with specific identity within a specified time period, with browser session and native app device instance isolation
  - Audit logging
  - Multi-factor authentication

- **For data that requires additional security**
  - HIPAA Personal Health Information (PHI) w/ BAA
  - Personally Identifiable Information (PII)
  - Sensitive but unclassified

- **NIST 800-171 Low**

- **Two additional subscription tiers**
  - **High assurance tier**: for all added security features
  - **BAA tier**: high assurance features plus BAA with Uchicago

- **Available this Summer**
  - Transfer, sharing, web app, CLI only. Excludes publish, search, identifiers, hosted CLI, GlobusID

# SSH with Globus Auth

- **Securely access resource using SSH with federated identity**
  - Leverage same security model as rest of data infrastructure
  - Facilitates automation
  - Eliminate need to manage SSH key lifecycle and provisioning

- **Replaces GSI SSH**

- **Client side wrapper around local SSH client (globus-ssh …)**

- **No changes to the SSH server (PAM module)**

- **Status:**
  - Prototype complete, early customer feedback
  - GA by end of year

- **Lightning talk: SSH with Globus Auth**

# New Storage Connectors

- **We continue to grow our connector set**

- **On near-term radar**
  - Box
  - Google Cloud Storage

- **What else do you want?**
  - Microsoft Azure Blob Storage
  - Wasabi
  - …

# Groups

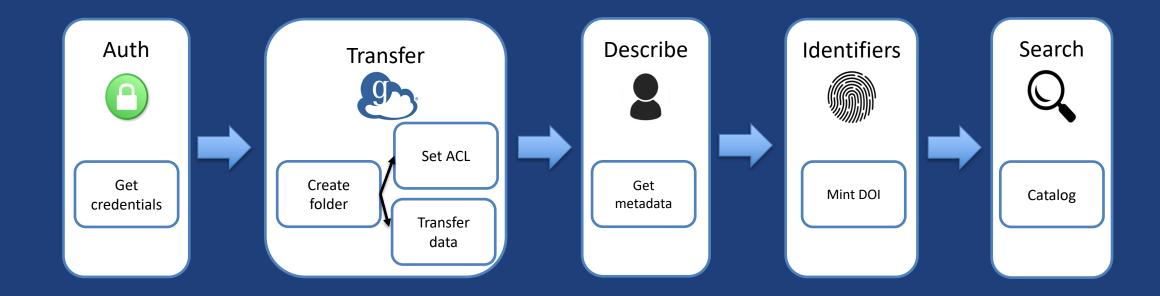- **Generally available in web app**

- **REST API has been in limited production**

- **Plan on opening some portion to general availability**
  – Please tell us your use cases
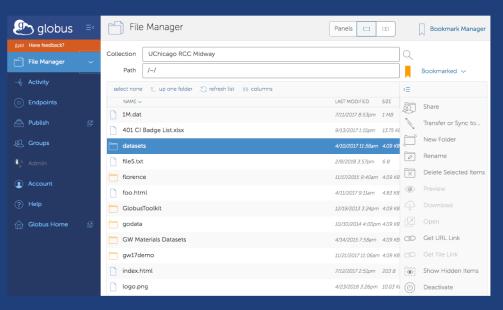
# Publication v2 platform

- **Form-based metadata entry**
- **Automate**

# New web app

- **Complete file manager for any research storage**

- **Improved browser experience**
  - **Accessibility**: WCAG 2.0 AA
  - **Responsiveness**: from large desktop to small phone
  - **Touch support**: for phones and pads

- **Leverage Globus Connect HTTPS**
  - E.g., Preview, download

- **Beta available now:**
  **https://app.globus.org**

# New Web App Demo
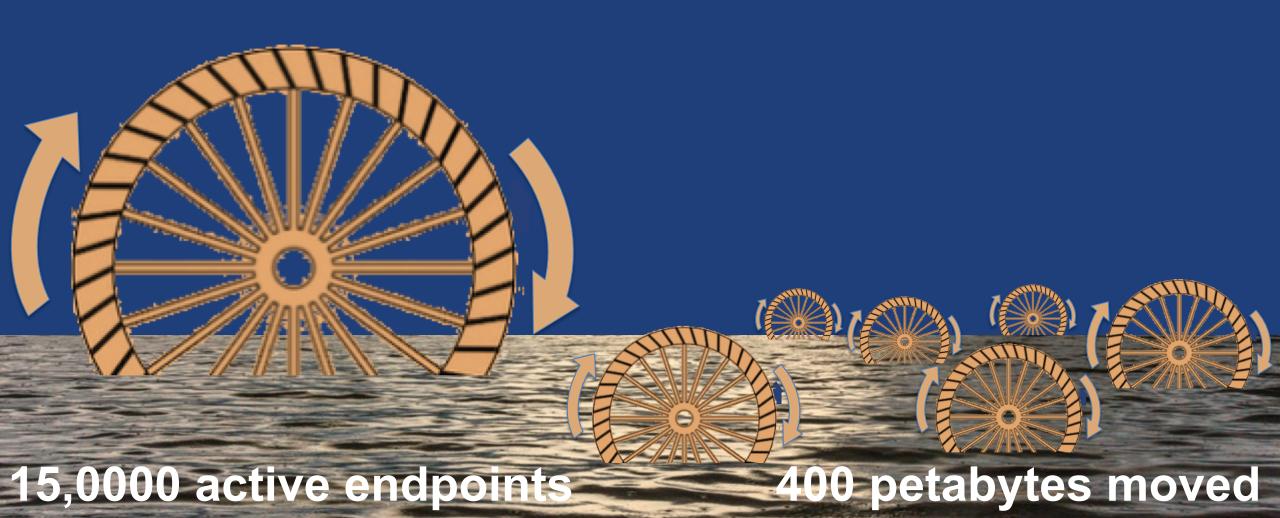
# Globus Labs

# globus △ labs -- labs.globus.org



Ian Foster

Ryan Chard

Tyler Skluzacek

Zhi Hong

Logan Ward

Yulie Zamora

Ben Blaiszik

Ricardo Lourenco

Roselyne Tchoua

Anna Woodward

Kyle Chard

Sam Nickolay

Jonathan Gaff

Steve Tuecke

# "Make all research data reliably, rapidly, and securely accessible, discoverable, and usable"

- Address computer science & domain science challenges

- Contribute to Globus with improved performance, new features, product directions, exploratory prototypes
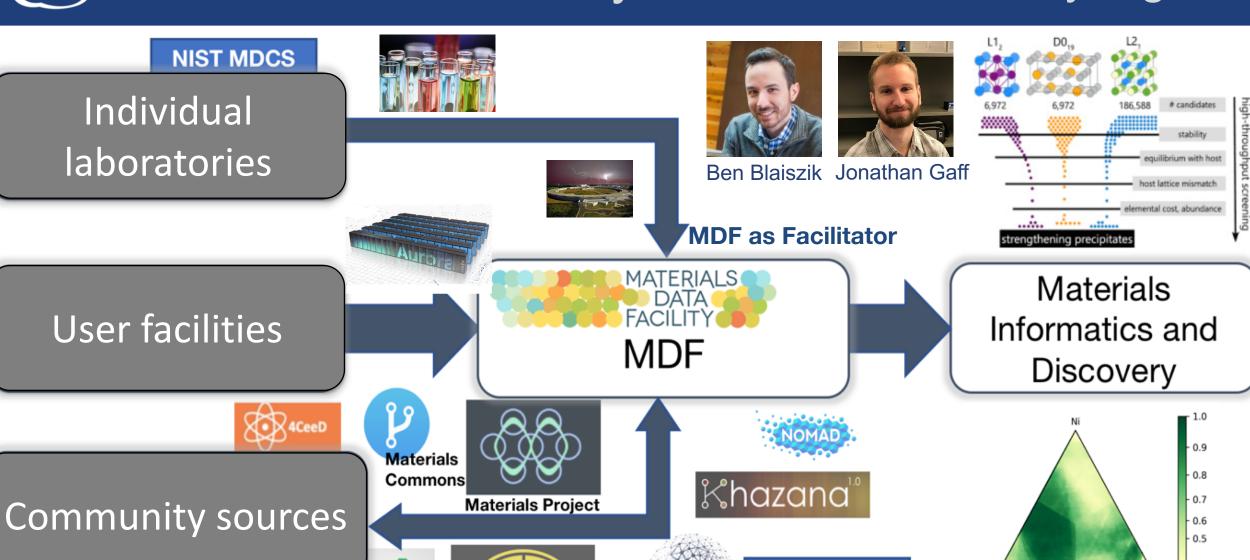
- Leverage advanced Globus features

# Building "data turbines"?

**15,0000 active endpoints**

**400 petabytes moved**

Materials Data Facility: materialsdatafacility.org

# Materials Data Facility adoption



**Publication**

| | | |
|---|---|---|
| **61** Total datasets | **29** Institutions | **22** CHiMaD datasets |
| **150** Authors | | **>18 TB** Data Volume |

**MDF Index**

| | |
|---|---|
| **117** Data resources indexed | **>3.4M** Records |
| **8** Repositories harvested | **~ 200** Datasets |
| | **~ 300 TB** Made discoverable |

**Kevin G. Yager**

Center for Functional Nanomaterials, Brookhaven National Laboratory
Verified email at bnl.gov - Homepage

scattering    SAXS    GISAXS    block-copolymers    self-assembly

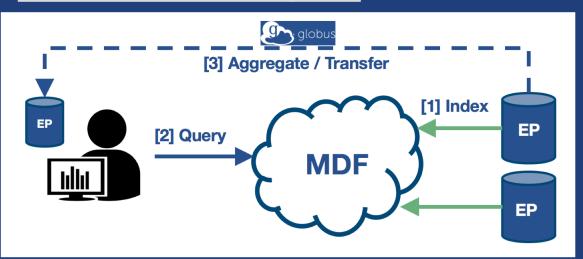| TITLE | CITED BY | YEAR |
|---|---|---|
| X-ray scattering image classification using deep learning<br>B Wang, K Yager, D Yu, M Hoai<br>Applications of Computer Vision (WACV), 2017 IEEE Winter Conference on, 697-704 | 4 | 2017 |
| Dataset of synthetic x-ray scattering images for classification using deep learning<br>KG Yager, J Lhermitte, D Yu, B Wang, Z Guan, J Liu<br>Materials Data Facility | 1 | 2017 |
| Operando grazing incidence small-angle X-ray scattering/X-ray diffraction of model ordered mesoporous Lithium-ion battery anodes<br>SM Bhaway, Z Qiang, Y Xia, X Xia, B Lee, KG Yager, L Zhang, ...<br>ACS nano 11 (2), 1443-1454 | 6 | 2017 |
| Nanoconfinement platform for nanostructure quantification via grazing-transmission X-ray scattering<br>CT Black, KG Yager<br>US Patent 9,557,283 | | 2017 |
| Beyond native block copolymer morphologies<br>GS Doerk, KG Yager<br>Molecular Systems Design & Engineering 2 (5), 518-538 | 3 | 2017 |
| Rapid assessment of crystal orientation in semi-crystalline polymer films using rotational zone annealing and impact of orientation on mechanical properties<br>C Ye, C Wang, J Wang, CG Wiener, X Xia, SZD Cheng, R Li, KG Yager, ...<br>Soft matter 13 (39), 7074-7084 | | 2017 |

# MDF is enabling new science: Dataset mixing

## Assemble Training Set





```
elements = ["Al"]
sources = ["khazana_vasp", "sluschi", "ab_initio_solute_database"]
my_ep = "c8ee7e5c-6d04-11e5-ba46-22000b92c6ec"
my_path = "/Users/ben/Desktop/blaiszik-macbookpro/dft_training_set"

mdf = Forge()
res = mdf.search_by_elements(elements=elements, sources=sources, limit=9999)
mdf.get_globus(res, dest=my_path,
               local_ep=my_ep, preserve_dir=True)
```

```
Processing records: 100%|██████████| 10/10 [00:00<00:00, 19.05it/s]
Submitting transfers: 100%|██████████| 1/1 [00:00<00:00,  3.60it/s]

All transfers submitted
Submission IDs: 3fbfc637-7181-11e7-a9fd-22000bf2d287
```
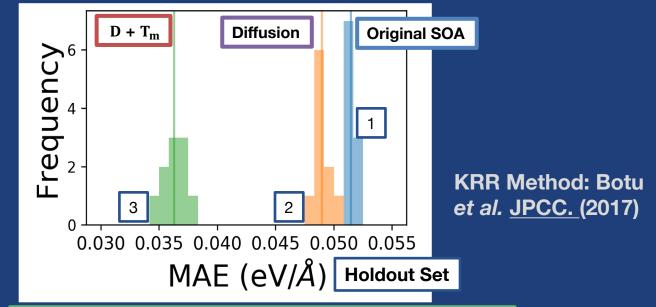
## Dataset Mixing

Build force-field potentials from different datasets

**Data resources:** 3 Aluminum DFT datasets

1 dataset from khazana.uconn.edu, 2 from materialsdata.nist.gov

**Result:** Improved performance by integrating data sources

Logan Ward



**KRR Method: Botu** *et al.* JPCC. (2017)

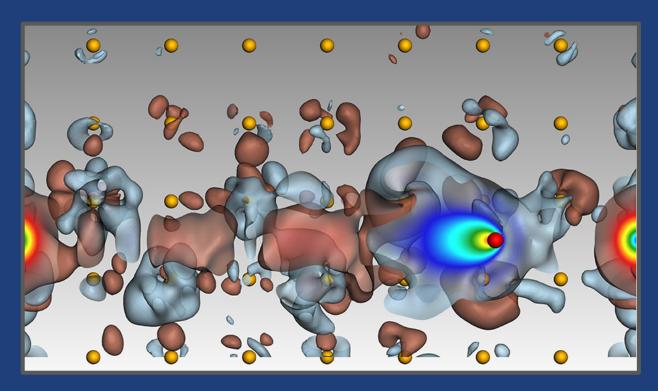Better performance in original application: No new DFT calculations

# MDF enables new science: Stopping power

**Stopping Power**: A "drag" force experienced by high speed protons, electrons, or positrons in a material

**Areas of Application**

- **Nuclear reactor safety**

- **Magnetic confinement / inertial containment for nuclear fusion**

- **Solar cell surface adsorption**

- **Medicine (e.g., proton therapy treatment)**

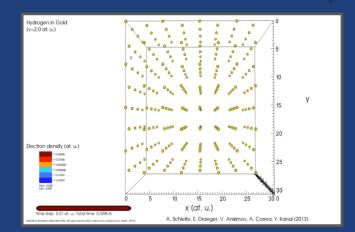- **Critical to understanding material radiation damage**



**André Schleife and Cheng-Wei Lee (UIUC) 2016 ALCF INCITE Project "Electronic Response to Particle Radiation in Condensed Matter"**
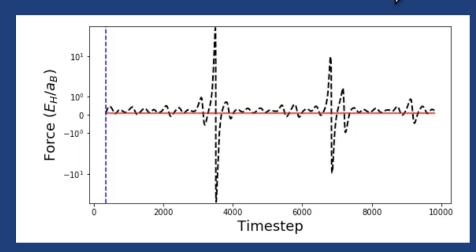
# Connecting MDF and HPC
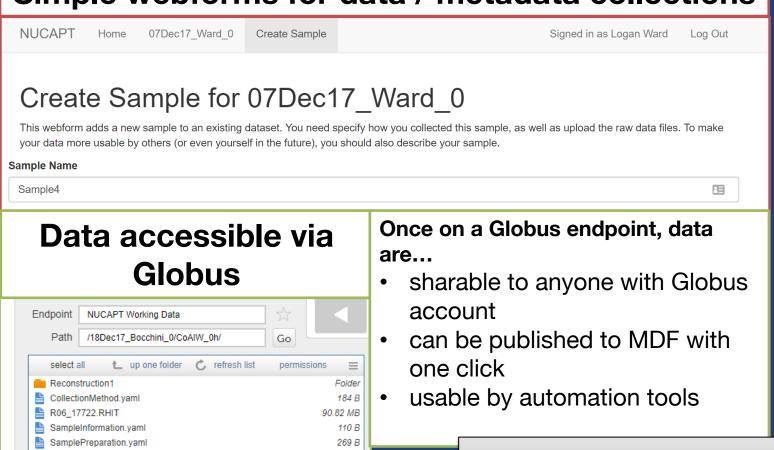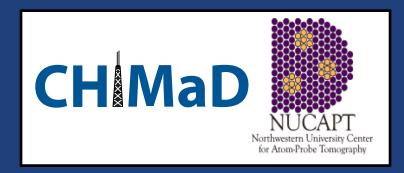
# Northwestern University Center for Atom-Probe Tomography (NUCAPT) Integration
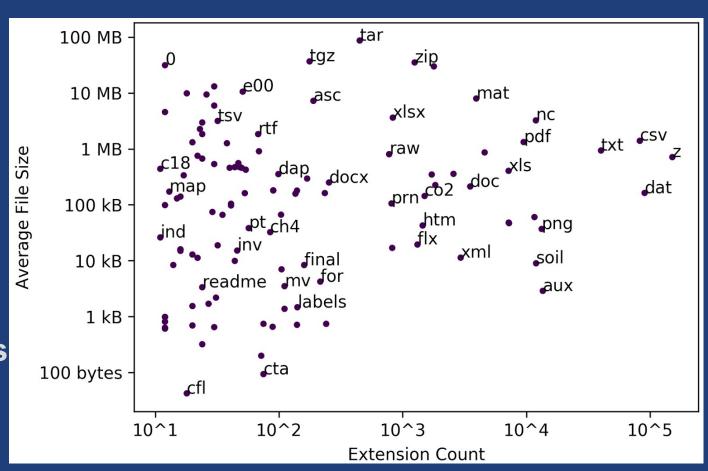
**Simple webforms for data / metadata collections**



NUCAPT    Home    07Dec17_Ward_0    Create Sample                    Signed in as Logan Ward    Log Out

Create Sample for 07Dec17_Ward_0

This webform adds a new sample to an existing dataset. You need specify how you collected this sample, as well as upload the raw data files. To make your data more usable by others (or even yourself in the future), you should also describe your sample.

**Sample Name**

Sample4

**Data accessible via Globus**

Endpoint    NUCAPT Working Data

Path    /18Dec17_Bocchini_0/CoAlW_0h/    Go

select all    up one folder    refresh list    permissions

Reconstruction1                Folder
CollectionMethod.yaml          184 B
R06_17722.RHIT                 90.82 MB
SampleInformation.yaml         110 B
SamplePreparation.yaml         269 B

**Once on a Globus endpoint, data are...**

- sharable to anyone with Globus account
- can be published to MDF with one click
- usable by automation tools

**Publish**



CHiMaD
NUCAPT
Northwestern University Center for Atom-Probe Tomography



MATERIALS DATA FACILITY

| Title | Author(s) |
|---|---|
| Influence of ruthenium in a model Co-Al-W superalloy | Sauza, Daniel; Bocchini, Peter; Chung, Ding-Wen; Dunand, David; Seidman, David |
| Atom Probe Tomography Reconstruction and Analysis for the Temporal Evolution of Co-Al-W Superalloys at 650˚C | Bocchini, Peter; Chung, Ding-Wen; Dunand, David; Seidman, David |
| γ+γ' Microstructures in the Co-Ta-V Ternary System | Reyes Tirado, Fernando; Perrin Toinin, Jacques; Dunand, David |
| γ+γ' Microstructures in the Co-Nb-V Ternary System | Reyes Tirado, Fernando; Perrin Toinin, Jacques; Dunand, David |
| Atom Probe Tomography Reconstruction and Analysis for the Temporal Evolution of Co-Al-W Superalloys at 750˚C | Bocchini, Peter; Chung, Ding-Wen; Dunand, David; Seidman, David |

# "Draining the data swamp"

**File systems and data repositories often inconsistent & disorganized**

**Many file types:**

- **Tabular**

- **Structured (JSON/XML)**

- **Images: photos, maps, plots**

- **Text: READMEs, papers, abstracts**

- **Not useful (?): Hadoop error logs, Windows installers, desktop shortcuts**
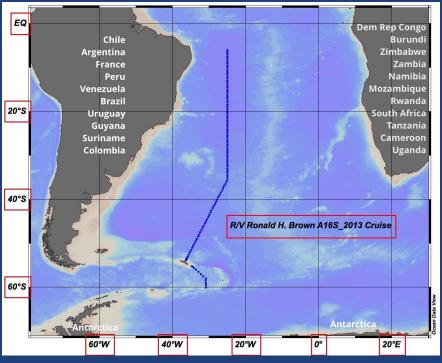
**Assumption: extension =/= content**



Carbon Dioxide Information Analysis Center (CDIAC) file extension frequency vs size for 500,001 files

# Skluma: automated metadata extraction pipelines

- Automatically crawls arbitrary file systems
- Constructs individualized extraction pipelines for each file
- Uses ML methods to determine extractors
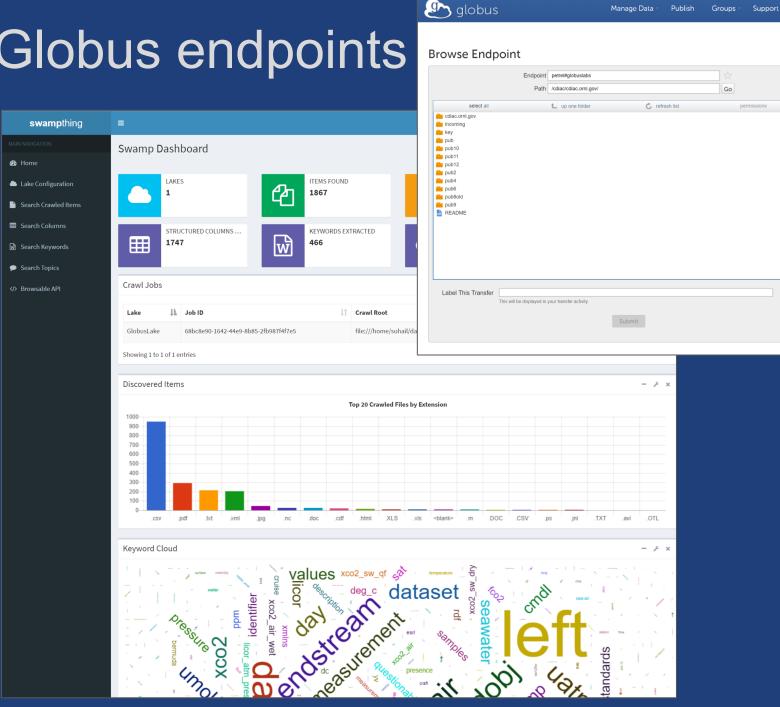- Handles *hybrid* files

# Skluma and Globus endpoints



- **Input: Globus endpoint on a file system or repository**

- **Output: Globus search index of extracted metadata**

- **User interface for visualizing and managing the data swamp**
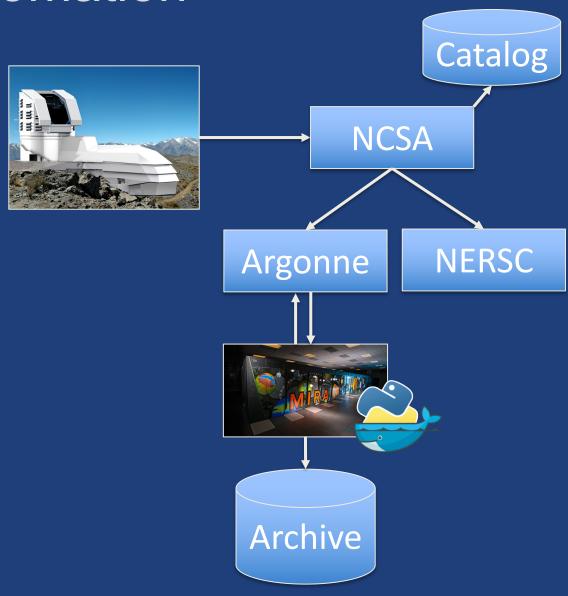
# Distributed research automation

Automation of scientific lifecycles:

- Data acquisition at different locations/times/instruments

- Analysis execution on distributed/heterogeneous resources

- Cataloging of descriptive metadata and provenance

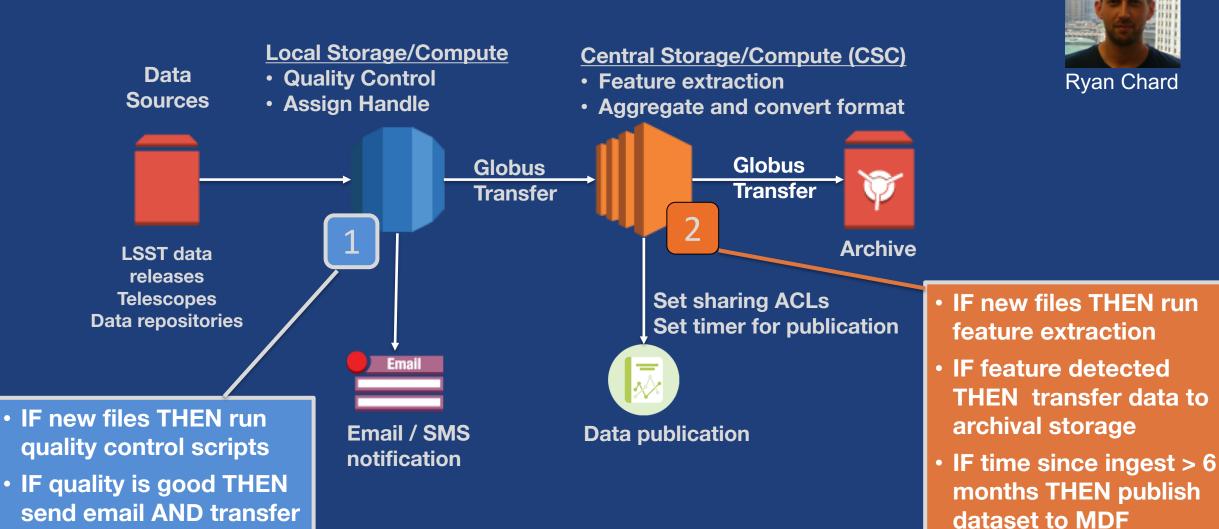- Dynamic collaborations around data and analysis



LSST data distribution and analysis pipeline

# Encoding automation flows using trigger-action programming

Ryan Chard

**Data Sources**

**Local Storage/Compute**
- **Quality Control**
- **Assign Handle**

**Central Storage/Compute (CSC)**
- **Feature extraction**
- **Aggregate and convert format**

LSST data releases
Telescopes
Data repositories

**Globus Transfer**

**Globus Transfer**

**Archive**

**1**

**2**

Email

Email / SMS notification

Set sharing ACLs
Set timer for publication

Data publication

- **IF new files THEN run quality control scripts**
- **IF quality is good THEN send email AND transfer data to CSC**

- **IF new files THEN run feature extraction**
- **IF feature detected THEN transfer data to archival storage**
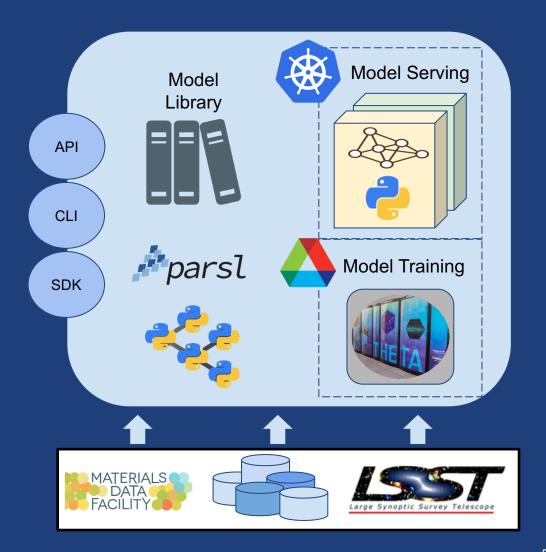- **IF time since ingest > 6 months THEN publish dataset to MDF**

# Deep Learning Hub (DLHub)

**Framework for finding, training, and serving machine learning and deep learning models**

**Scalable, container-based execution model**

**Globus integration to connect external repositories**

# Parallel Scripting in Python

**Parsl**

Parallel programming library for Python

> **pip install parsl**

Annotate functions to define Parsl *apps*
- Bash apps call external applications
- Python apps call Python functions

Apps run concurrently respecting data dependencies via futures. Natural parallel programming!

Parsl scripts are independent of where they run. Write once run anywhere!

```python
@App('python', dfk)
def hello ():
    return 'Hello World!'

print (hello().result())
```
```
Hello World!
```

```python
@App('bash', dfk)
def echo_hello(stdout='echo-hello.stdout'):
    return 'echo "Hello World!"'

echo_hello().result()

with open('echo-hello.stdout', 'r') as f:
    print(f.read())
```
```
Hello World!
```

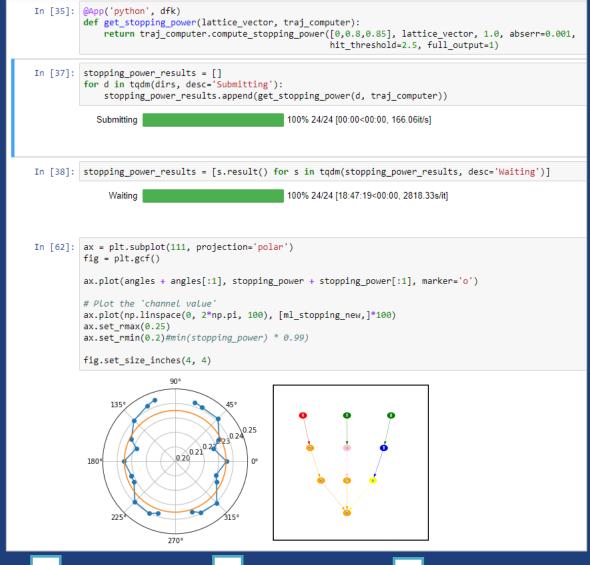# Interactive scalable computing in Jupyter



```python
In [35]: @App('python', dfk)
         def get_stopping_power(lattice_vector, traj_computer):
             return traj_computer.compute_stopping_power([0,0.8,0.85], lattice_vector, 1.0, abserr=0.001,
                                                         hit_threshold=2.5, full_output=1)
```

```python
In [37]: stopping_power_results = []
         for d in tqdm(dirs, desc='Submitting'):
             stopping_power_results.append(get_stopping_power(d, traj_computer))
```
Submitting ████████████████ 100% 24/24 [00:00<00:00, 166.06it/s]

```python
In [38]: stopping_power_results = [s.result() for s in tqdm(stopping_power_results, desc='Waiting')]
```
Waiting ████████████████ 100% 24/24 [18:47:19<00:00, 2818.33s/it]

```python
In [62]: ax = plt.subplot(111, projection='polar')
         fig = plt.gcf()

         ax.plot(angles + angles[:1], stopping_power + stopping_power[:1], marker='o')

         # Plot the 'channel value'
         ax.plot(np.linspace(0, 2*np.pi, 100), [ml_stopping_new,]*100)
         ax.set_rmax(0.25)
         ax.set_rmin(0.2)#min(stopping_power) * 0.99)

         fig.set_size_inches(4, 4)
```

Globus integration (Auth and transfer)

Multi-site execution

Exec provider/model independent

Automated elasticity

Container support
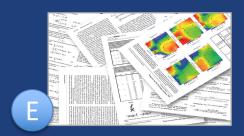
Fault tolerance

Python/Jupyter

Data servers

amazon web services

XSEDE
Extreme Science and Engineering Discovery Environment

59

# Parsl adoption across sciences



A — Simulating galaxy formation using sky surveys

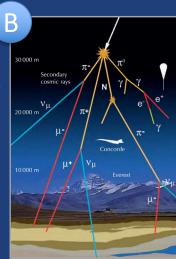B — Analysis of Cosmic ray showers

C — Real-time computing (e.g., APS image reconstruction)

D — Protein docking

E — Information extraction

F — Biomolecule modelling

G — Machine learning and data analytics

H — Computational calculation of stopping power

I — Education programs in Jupyter notebooks

# Parsl + Globus + Jupyter

**Native app integration** to provide embedded access to Globus (and other) services

**Transparent SSH-based authentication** to compute resources (soon)

**High-speed data access** and **reliable data movement** to/from repositories, laptops, supercomputers, …

**Staging to/from DTNs**



```
In [1]: import globus_sdk

CLIENT_ID = '4790b51f-7c6b-4727-8d85-a761a417b8ac'

native_auth_client = globus_sdk.NativeAppAuthClient(CLIENT_ID)

native_auth_client.oauth2_start_flow(requested_scopes="urn:globus:auth:scope:data.materialsdatafacility.org:all urn:globus:auth:s

print("Login Here:\n\n{0}".format(native_auth_client.oauth2_get_authorize_url()))

print(("\n\nNote that this link can only be used once! "
       "If login or a later step in the flow fails, you must restart it."))
```

Login Here:

https://auth.globus.org/v2/oauth2/authorize?client_id=4790b51f-7c6b-4727-8d85-a761a417b8ac&redirect_uri=https%3A%2F%2Fauth.globus.org%2Fv2%2Fweb%2Fauth-code&scope=urn%3Aglobus%3Aauth%3Ascope%3Adata.materialsdatafacility.org%3Aall+urn%3Aglobus%3Aauth%3Ascope%3Atransfer.api.globus.org%3Aall+urn%3Aglobus%3Aauth%3Ascope%3Aauth.globus.org%3Aview_identities+openid+email+profile+urn%3Aglobus%3Aauth%3Ascope%3Asearch.api.globus.org%3Aall&state=_default&response_type=code&code_challenge=6087u8mbP4JAcf1Mgfk8TewLE_-4F1RzRjByKunanE8%3D&code_challenge_method=S256&access_type=online

Log in to use SDK / Jupyter client

Use your existing organizational login

e.g., university, national lab, facility, project

University of Chicago

Didn't find your organization? Then use **Globus ID** to sign in. (What's this?)

Continue

SDK / Jupyter client would like to:

✓ HTTPS Server data.materialsdatafacility.org ⓘ
✓ Transfer files using Globus Transfer ⓘ
✓ View your identities on Globus Auth ⓘ
✓ Know who you are in Globus. ⓘ
✓ Know some details about you. ⓘ
✓ Know your email address. ⓘ
✓ Access the Globus Search API ⓘ

To work, the above will need to:

✓ View your identities on Globus Auth ⓘ
✓ Manage your Globus Groups ⓘ

```
sorted_file = File(
    "globus://ddb59aef-6d04-11e5-ba46-22000b92c6ec/~/sorted.txt")

dfu = unsorted_file.stage_in()
dfu.result()

f = sort_strings(inputs=[dfu], outputs=[sorted_file])
f.result()
```

# Program Preview

**globusworld.org/conf/program**

- **Today**
  - Lightning talks
  - Guest keynote: Alex Szalay, Building the Open Storage Network
  - Reception

- **Tomorrow**
  - Tutorials
  - Office Hours

- **Friday morning**
  - Customer forum

# #globus2018

# @globus