

# Materials Data Facility: A Distributed Model for the Materials Data Community

Logan Ward<sup>1</sup> (loganw@uchicago.edu)

Ben Blaiszik<sup>1,2</sup> (blaiszik@uchicago.edu),

Ian Foster (foster@uchicago.edu)<sup>1,2</sup>, Ryan Chard<sup>2</sup>

Jonathon Gaff<sup>1</sup>, Kyle Chard<sup>1</sup>, Jim Pruyne<sup>1</sup>,

Rachana Ananthakrishnan<sup>1</sup>, Steven Tuecke<sup>1</sup>

Michael Ondrejcek<sup>3</sup>, Kenton McHenry<sup>3</sup>, John Towns<sup>3</sup>

University of Chicago<sup>1</sup>, Argonne National Laboratory<sup>2</sup>, University of Illinois at Urbana-Champaign<sup>3</sup>

[materialsdatafacility.org](http://materialsdatafacility.org)  
[globus.org](http://globus.org)

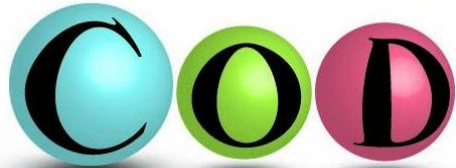


Materials Genome Initiative

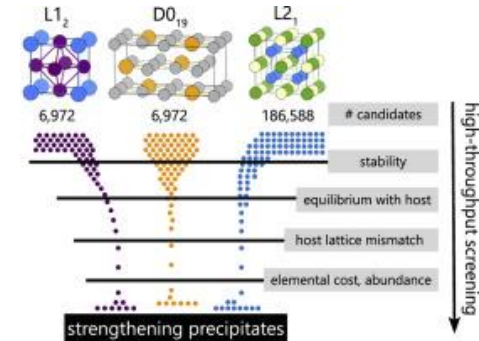


# Data-Intensive Materials Science

## Materials Databases

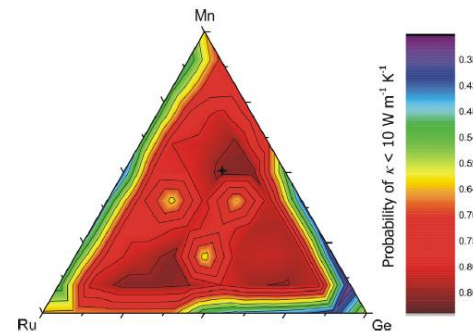
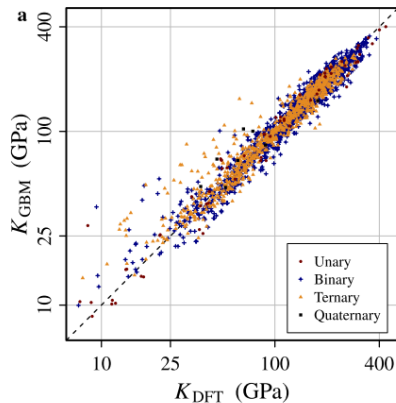


## High-Throughput Screening

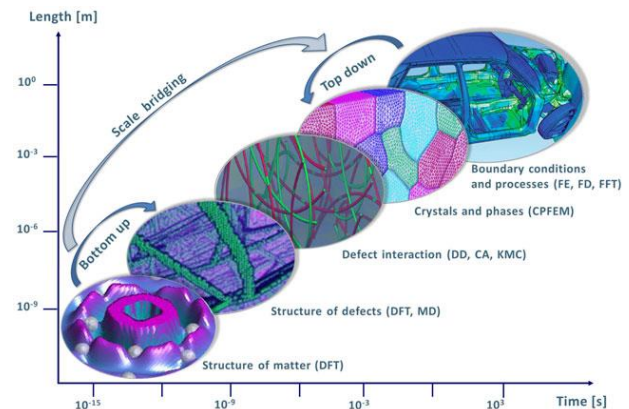


Kirklin *et al.* Acta Mat. (2016)

## Machine Learning



## Multi-scale Modeling



de Jong *et al.* Sci Rep. (2016)

Sparks *et al.* Scr. Mat. (2015)

<https://www.mpg.de/>

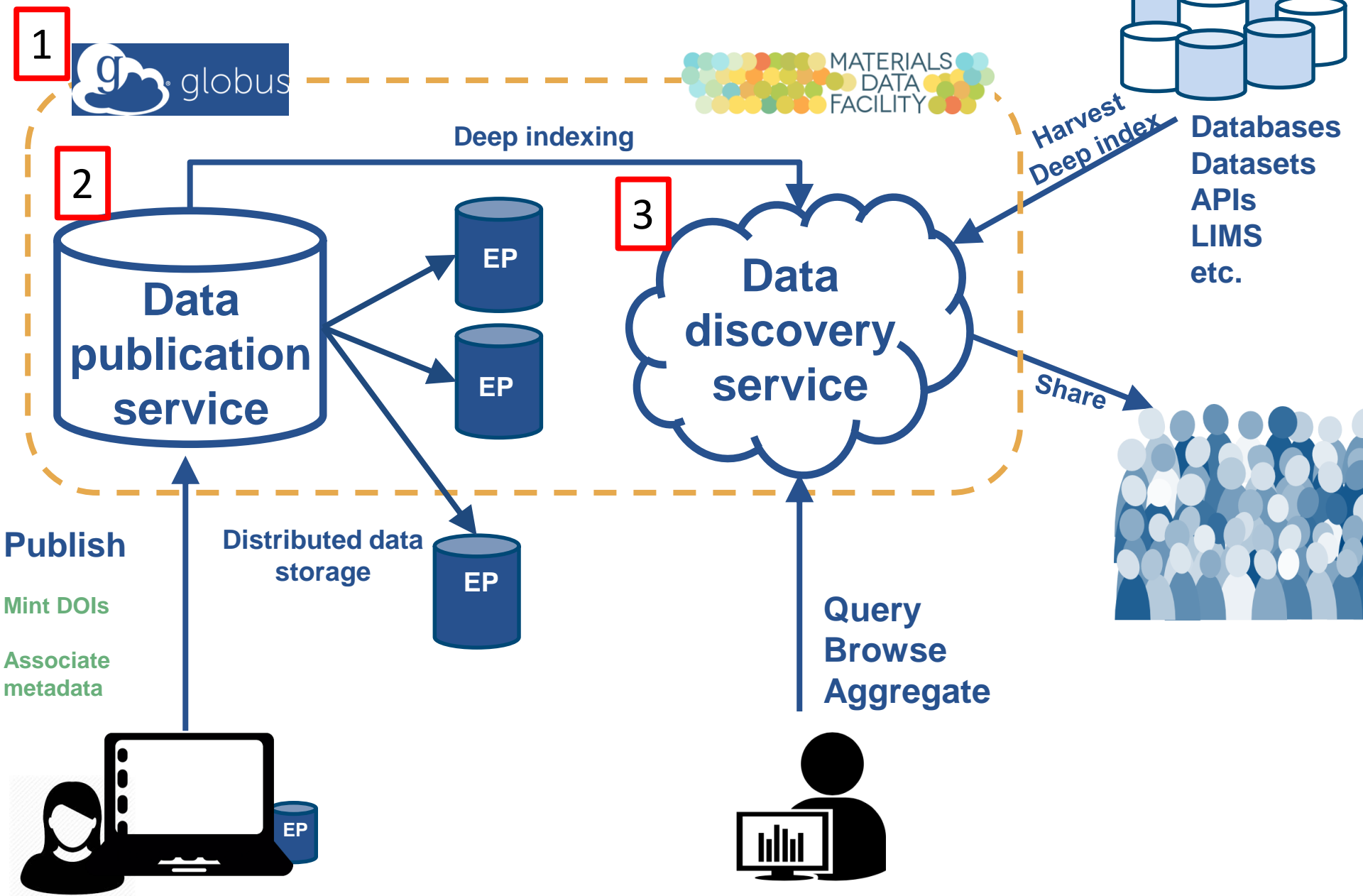
# Data-Intensive Materials Science

Science is becoming limited by the ability to handle data

- Where to get it?
- How to selectively share it?
- Where to store it?
- How to know what it is?
- How to build software that uses it?
- How to get others to share theirs?
- How to keep track of provenance?
- ....?

**Our goal is to create infrastructure that provides easy answers to these questions**

# What is the MDF?



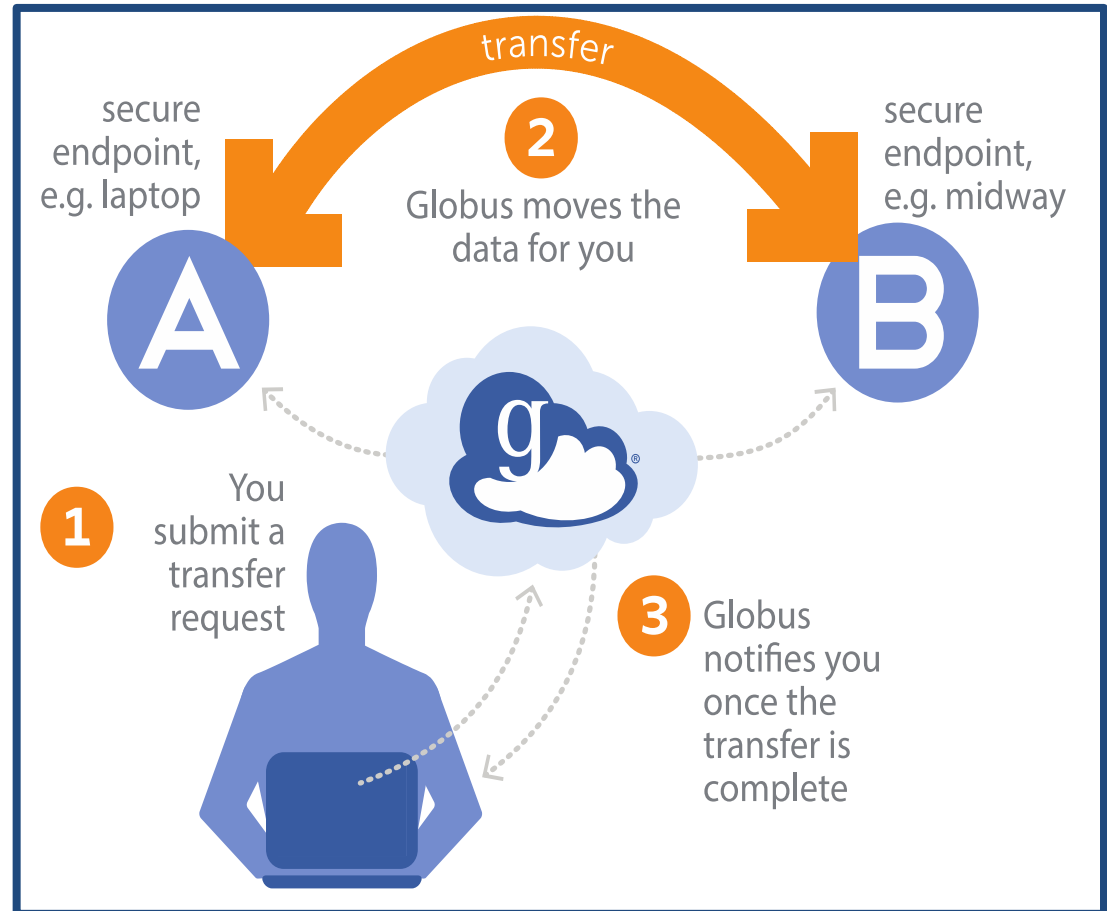
# Globus Background

## Endpoint

- E.g. laptop or server running a Globus client (e.g. Dropbox client)
- Enables advanced file transfer and sharing
- Currently GridFTP, future GridFTP + HTTP

## Some Key Features

- REST API for automation and interoperability
- Web UI for convenience
- Optimizes and verifies transfers
- Handles auto-restarts



313,101,618,927 MB TRANSFERRED

# Globus Platform-as-a-Service (PaaS)

## Identity management

- create and manage a unique identity linked to external identities for authentication

## User groups

- Manage user group creation and administration flows
- Share data with user groups

## Publication

## Discovery

## Data transfer

- High-performance data transfer from a web browser
- Optimize transfer settings and verify transfer integrity
- Add your laptop to the Globus cloud with Globus Connect Personal

## Data sharing

- Share directly from your storage device (laptop or cluster)
- File and directory-level ACLs

# Data sharing and Globus

Transfer Files | Activity | Manage Endpoints | Dashboard | Console

Transfer Files

Get Globus Connect Personal  
Turn your computer into an endpoint.

Manage Shared Endpoint

« shared endpoints list

Manage Permissions for ranantha#demo12

Host: ucrcc#midway: /~/share/gptest/

| name                               | read                                | write                               |
|------------------------------------|-------------------------------------|-------------------------------------|
| Path: /                            |                                     |                                     |
| Rachana Ananthakrishnan (ranantha) | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Ian Foster (ian1)                  | <input checked="" type="checkbox"/> | <input type="checkbox"/>            |

view link for sharing

Manage Roles (new tab) + Add Permission

**Easily control who gains access to your data:**

- Globus can use University/Laboratory credentials
- You can establish groups of authorized users

# REST APIs, Clients, and Docs

- New Python SDK available
  - <https://github.com/globusonline/globus-sdk-python>
- Jupyter Notebook Examples
  - <https://github.com/globus/globus-jupyter-notebooks>
- Sample Data Portal
  - <https://github.com/globus/globus-sample-data-portal>
- (alpha) MDF Data Publication Service API

## Endpoint search

Globus has over 8000 registered endpoints. To find endpoints of interest you can access powerful search capabilities via the SDK. For example, to search for a given string across the descriptive fields of endpoints (names, description, keywords):

```
search_str = "Globus Tutorial Endpoint"
endpoints = tc.endpoint_search(search_str)
print("==== Displaying endpoint matches for search: '{}' ===".format(search_str))
for ep in endpoints:
    print("{} ({}).format(ep["display_name"] or ep["canonical_name"], ep["id"])))
```

## Restricting search scope with filters

There are also a number of default filters to restrict the search for 'my-endpoints', 'my-gcp-endpoints', 'recently-used', 'in-use', 'shared-by-me', 'shared-with-me')

```
search_str = None
endpoints = tc.endpoint_search(
    filter_fulltext=search_str, filter_scope="recently-used")
for ep in endpoints:
    print("{} ({}).format(ep["display_name"] or ep["canonical_name"], ep["id"])))
```

## Endpoint details

You can also retrieve complete information about an endpoint, including name, owner, location, and server configurations.

```
endpoint = tc.get_endpoint(tutorial_endpoint_1)
print("Display name:", endpoint["display_name"])
print("Owner:", endpoint["owner_string"])
print("ID:", endpoint["id"])
```

## Transfer

Creating a transfer is a two stage process. First you must create a description of the data you want to transfer (which also creates a unique submission\_id), and then you can submit the request to Globus to transfer that data.

If the submit\_transfer fails, you can safely resubmit the same transfer\_data again. The submission\_id will ensure that this transfer request will be submitted once and only once.

```
# help(tc.submit_transfer)
source_endpoint_id = tutorial_endpoint_1
source_path = "/share/godata/"

dest_endpoint_id = tutorial_endpoint_2
dest_path = "/-/

label = "My tutorial transfer"

# TransferData() automatically gets a submission_id for once-and-only-once submission
tdata = globus_sdk.TransferData(tc, source_endpoint_id,
                                dest_endpoint_id,
                                label=label)

## Recursively transfer source path contents
tdata.add_item(source_path, dest_path, recursive=True)

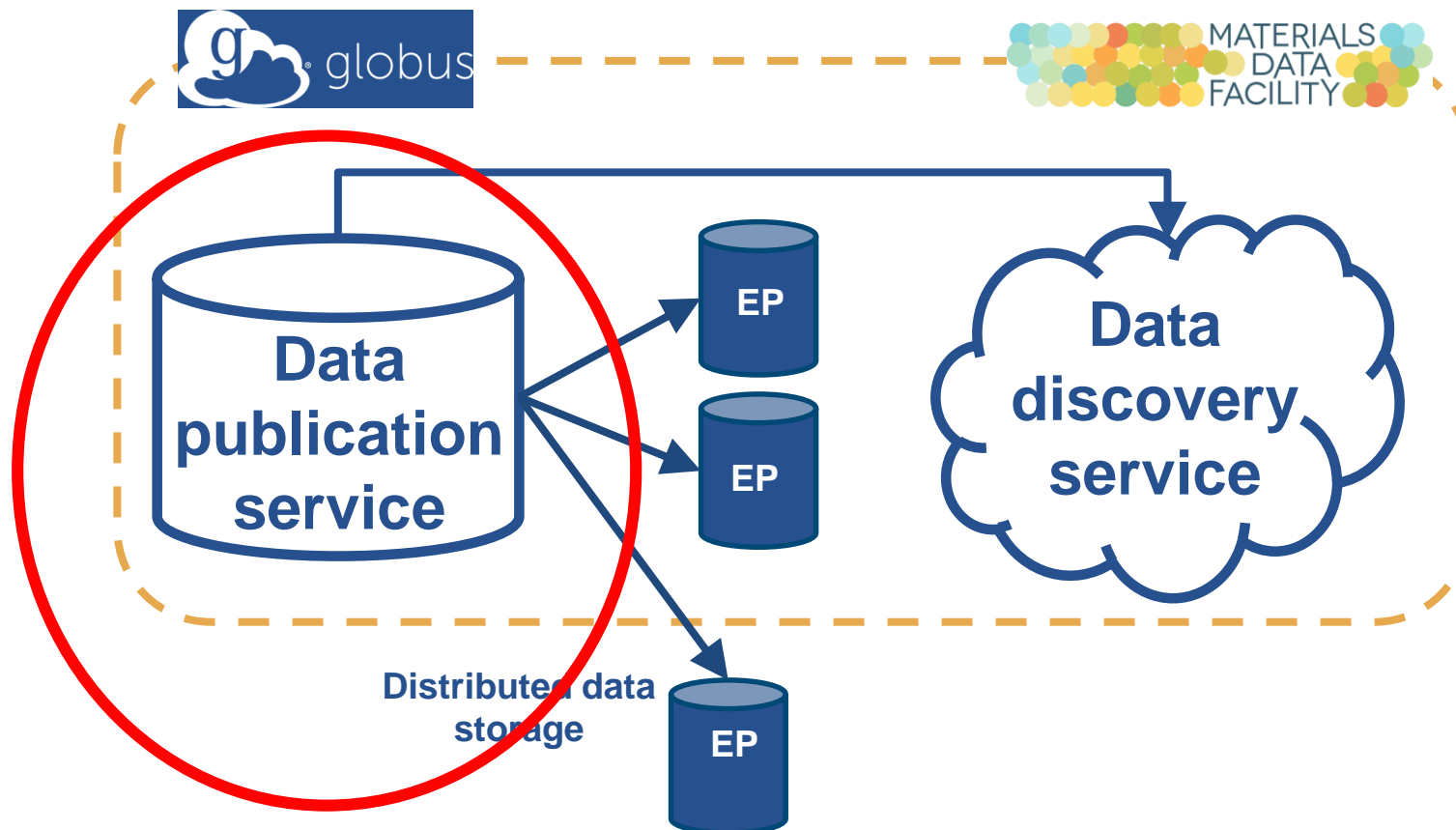
## Alternatively, transfer a specific file
# tdata.add_item("/source/path/file.txt",
#               "/dest/path/file.txt")

# Ensure endpoints are activated
tc.endpoint_autoactivate(source_endpoint_id)
tc.endpoint_autoactivate(dest_endpoint_id)

submit_result = tc.submit_transfer(tdata)
print("Task ID:", submit_result["task_id"])
```



# DATA PUBLICATION



# Materials Data Publication Service

MDF Open Collection home page

Open collection for submission of materials-related datasets

Submit to This Collection

Browse

Issue Date

Author

Title

Subject

Datasets in Collection (sorted by Submit Date in Descending order): 1 to 20 of 25

[next >](#)

| Issue Date  | Title  | Author(s)  |
|-------------|--|--|
| 22-Sep-2017 | <a href="#">Dataset for A New Generation of Effective Core Potentials for Correlated Calculations</a>              | <i>Bennett, M. Chandler; Melton, Cody A.; Annaberdiyev, Abdulgani; Wang, Guangming; Shulenburg, Luke; Mitas, Lubos</i> |
| 11-Sep-2017 | <a href="#">Probing the growth and melting pathways of a decagonal quasicrystal in real-time</a>                   | <i>Han, Insung; Xiao, Xianghui; Shahani, Ashwin J.</i>   |
| 6-Sep-2017  | <a href="#">Simulated microstructures of gamma' precipitates in cobalt-based superalloys</a>                       | <i>Jokisaari, Andrea M.; Naghavi, Shahab; Wolverton, Chris; Voorhees, Peter W.; Heinonen, Olle G.</i>                  |
| 23-Aug-2017 | <a href="#">Solute transport database in Mg using ab initio and exact diffusion theory</a>                         | <i>Agarwal, Ravi; Trinkle, Dallas R.</i>   |
| 29-Jun-2017 | <a href="#">Characterizing the Unifying Thread in High Temperature Superconductors Using Realistic Simulations</a> | <i>Narayan, Awadhesh; Busemeyer, Brian; Wagner, Lucas K.</i>   |

# Datasets Are Citable

| Title  | 1-20   | Cited by | Year |
|--|--|----------|------|
| Implications of Grain Size Variation in Magnetic Field Alignment of Block Copolymer Blends         | Y Rokhlenko, PW Majewski, SR Larson, P Gopalan, KG Yager, CO Osuji<br>American Chemical Society                            |          | 2017 |
| X-ray Scattering Image Classification Using Deep Learning  | B Wang, K Yager, D Yu, M Hoai<br>Applications of Computer Vision (WACV), 2017 IEEE Winter Conference on, 697-704           | 1        | 2017 |
| Dataset of synthetic x-ray scattering images for classification using deep learning                | KG Yager, J Lhermitte, D Yu, B Wang, Z Guan, J Liu<br>Materials Data Facility  | 1        | 2017 |
| Magnetic field alignment of coil-coil diblock copolymers and blends via intrinsic chain anisotropy | Y Rokhlenko, P Majewski, S Larson, K Yager, P Gopalan, A Avgeropoulos, ...<br>Bulletin of the American Physical Society 62 |          | 2017 |

# Publication statistics

|                     |                                      |                               |                             |
|---------------------|--------------------------------------|-------------------------------|-----------------------------|
| <b>Data Volumes</b> | <b>15.0 TB</b><br><b>13.4 TB out</b> |                               |                             |
| <b>Publication</b>  | <b>50</b><br>Total datasets          | <b>16</b><br>CHiMaD datasets  |                             |
|                     | <b>94</b><br>Authors                 | <b>14</b><br>Institutions     | <b>&gt;1000</b><br>Accesses |
| <b>Pipeline</b>     | <b>+30</b><br>Total datasets         | <b>+14</b><br>CHiMaD datasets |                             |

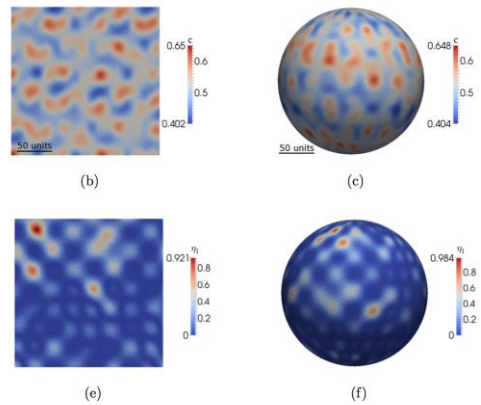
# Publication Route #1: MDF Storage

Grain Structure, Grain-averaged Lattice Strains, and Macro-scale Strain Data for Superelastic Nickel-Titanium Shape Memory Alloy Polycrystal Loaded in Tension

- Largest dataset to date (>1.5 TB). Showcases MDF unique capabilities and makes a unique dataset discoverable for code development, analysis, and benchmarking

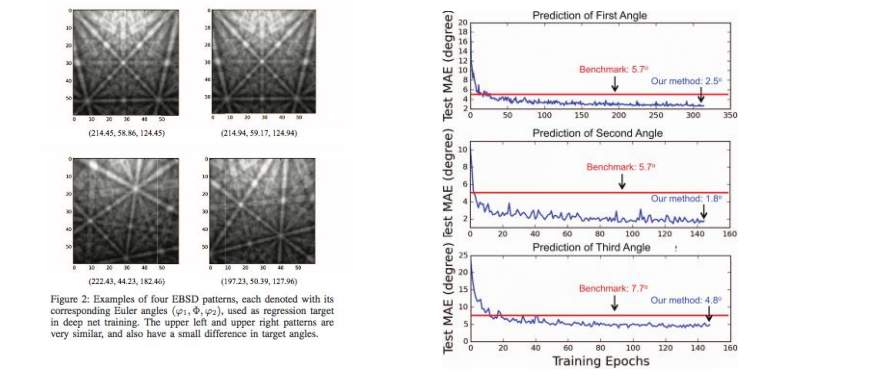
Paranjape et al.

## Phase Field Benchmark I Dataset



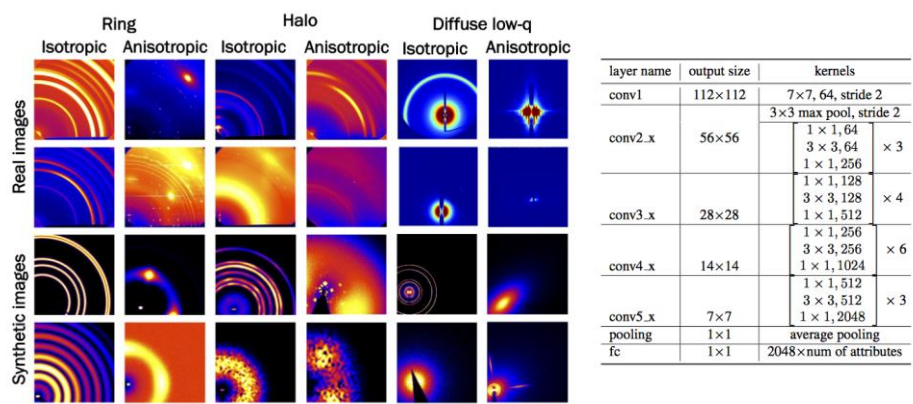
Jokisaari et al.

## Electron Backscattering and Diffraction Datasets for Ni, Mg, Fe, Si



Marc De Graef et al.

## X-ray Scattering Image Classification Using Deep Learning



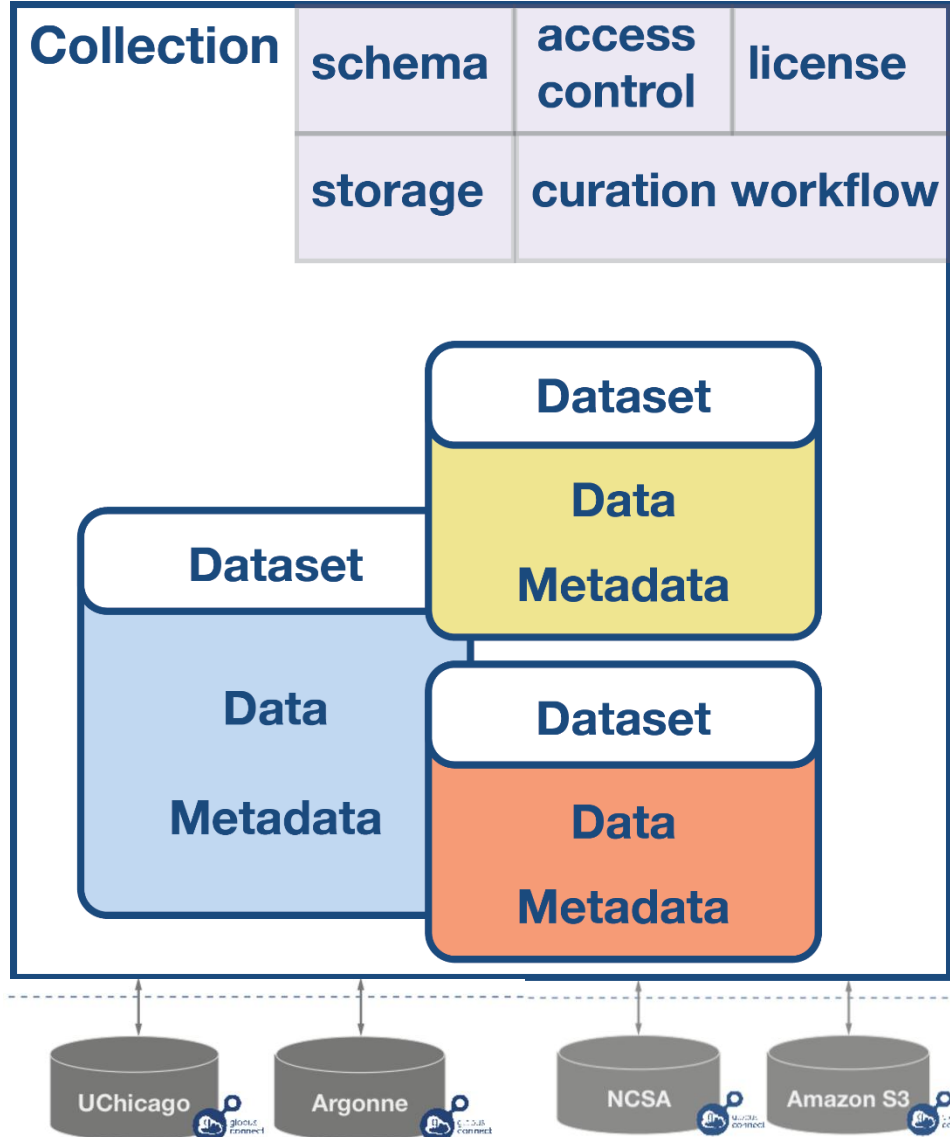
Yager et al. <http://dx.doi.org/10.18126/M2Z30Z>

# Customization: Collection Model

## Collections in this community

|   |
|---|
| <a href="#">APS Sector 1</a><br>Collection of datasets from Sector 1 at the Advanced Photon Source at Argonne National Laboratory |
| <a href="#">CHiMaD Team</a>   |
| <a href="#">Citrine Test</a>  |
| <a href="#">Hersam Group</a>  |
| <a href="#">MDF Open</a><br>MDF Open Collection   |
| <a href="#">MDF Test</a><br>Test Collection for MDF   |
| <a href="#">Voorhees Group</a>  |
|   |

# Customization: Collection Model



- Collections might be a research group or a research topic...
- Collections have specified
  - Mapping to storage endpoint
    - Currently handled as automatically created shared endpoints
  - Metadata schemas
  - Access control policies
  - Licenses
  - Curation workflows
- Collections contain
  - Datasets
    - Data
    - Metadata
- Metadata Persistence
  - Metadata log file with dataset
  - Metadata replicated in search index

# Share Data with Flexible ACLs



- Share data publicly, with a set of users, or keep data private

## Leverage Curation Workflows



- Collection administrators can specify the level of curation workflow required for a given collection e.g.
  - No curation
  - Curation of metadata only
  - Curation of metadata and files



# Example: NUCAPT Data Publication



## Goal:

- Aid metadata capture
- Simplify data publication

## Approach: Lightweight web service

- Form-based metadata capture
- Automatic file management
- “One-click” data publication

## Results:

- Beta version deployed Sept '17

**Sample Information**

Sample Title  
New Sample

Short description of sample

Sample Description  
Sample for a screenshot

Longer-form description of the sample

Sample Metadata

| Key        | Value |
|------------|-------|
| Aging time | 4 hr  |

**Data Collection Metadata**

Metadata about how a sample was collected

LEAP Model  
NUCAPT

Model of LEAP used to collect data

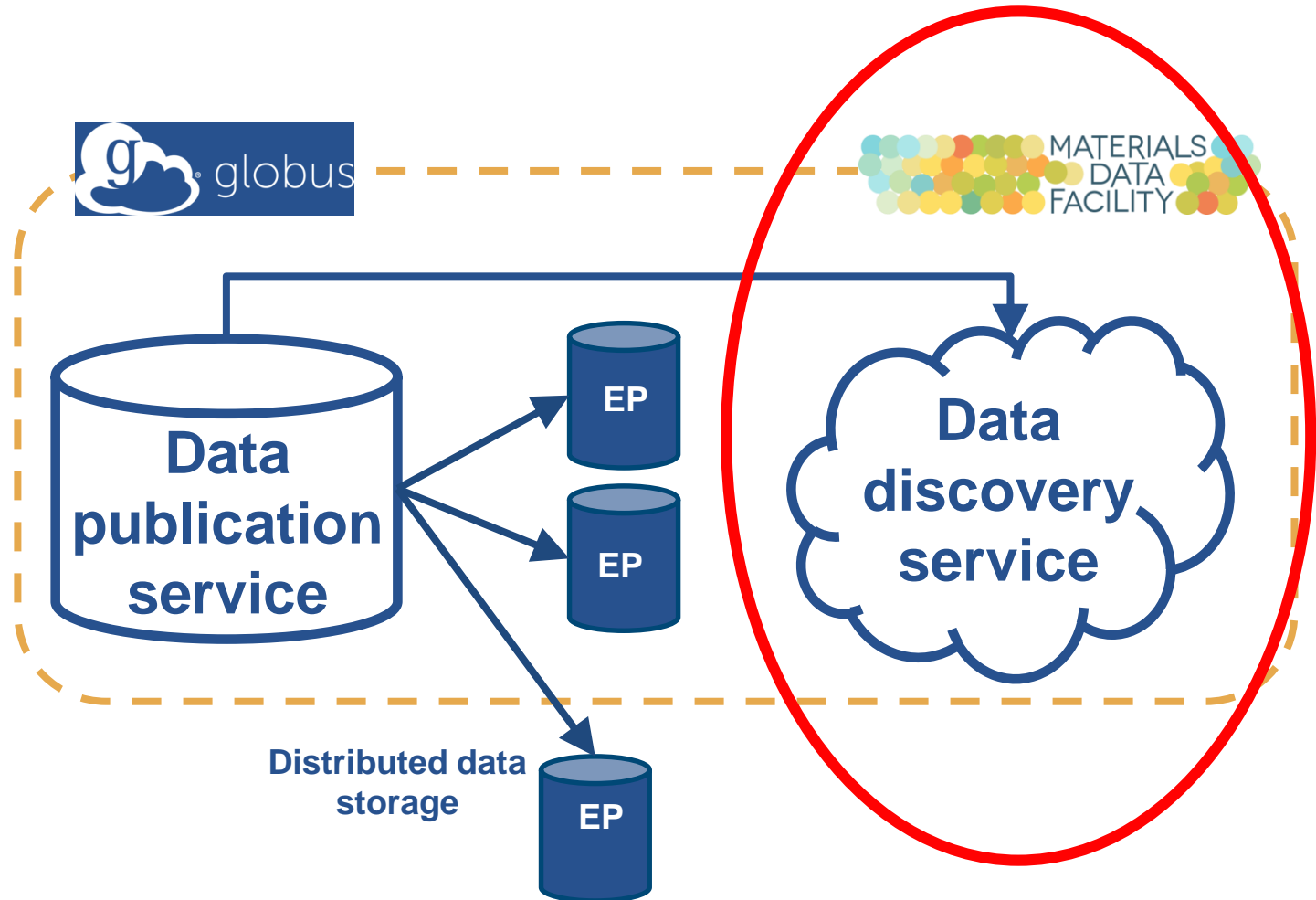
Evaporation Mode  
☒ Voltage  
☐ Laser

18Jul17\_Ward\_0

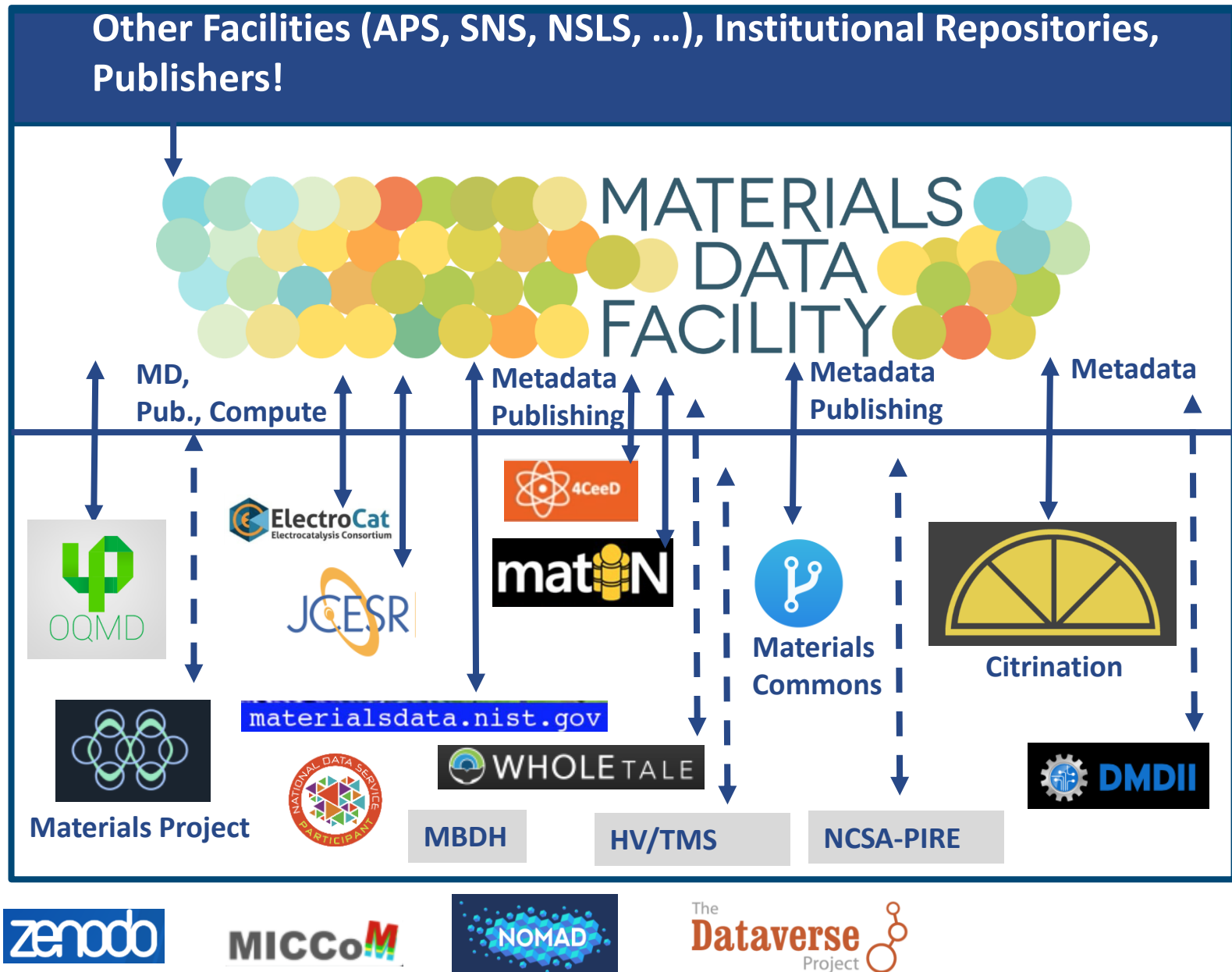
- Sample1
  - Recon1
    - Reconstruction2
      - Reconstruction3
        - 1D\_Concentration\_Profile
        - 2D\_Concentration\_Map
        - Component\_Distribution
        - Mass\_Spectrum
        - Proximity\_Histogram
        - Tip\_Composition
        - Visualization
        - example.POS

data.yaml

# DATA DISCOVERY [AND USE]



# Part 1: Linking with the Data Community



# Many Databases, Single Search

globus search Demo

Logan Ward

Log Out

Aluminum

mdf

☐ Enable Advanced Searching Options

Resource Type

☐ record

(1583)

☐ dataset

(95)

Elements

☐ Al

(1243)

☐ O

(760)

☐ C

(188)

☐ Si

(167)

☐ H

(165)

☐ N

(101)

☐ Ni

(86)

☐ S

(58)

☐ Pd

(55)

☐ r

(51)

Tags

☐ sdf

(356)

☐ alloy

(95)

☐ parent\_id

(95)

☐

(89)

☐ Computational File Repository Ca...

(6)

☐ Aluminum

(5)

☐ cif

(5)

☐ dif

(5)

☐ File Repository Categories::Chem...

(3)

☐ aluminum

(2)

You are searching as **Logan Ward** (LoganWard2012@u.northwestern.edu)

Search Results

AMCS - Aluminum

Collection: AMCS

Material Composition: Al4

AMCS - Aluminum

Collection: AMCS

Material Composition: Al4

Aluminum

Collection: NIST Material Measurement Laboratory

Description: Aluminum has many outstanding attributes that lead to a wide range of applications, including: 1) Good corrosion and oxidation resistance; 2) High electrical and thermal conductivities; 3) Low density; 4) High reflectivity; 5) High ductility and reasonably high strength; and 6) Relatively low cost. 6xxx and 6061 mentioned numerous time throughout

AMCS - Aluminum

Collection: AMCS

Material Composition: Al4

AMCS - Aluminum

Collection: AMCS

Material Composition: Al4

# MDF + NIST Database Tools



## Querying Nanomine Data

Example using the [Materials Data Facility](#) to query data from [NanoMine](#)

```
In [1]: from md_forge.forge import Forge
```

### Get All Records

Get all of the records in NanoMine

MDF automates publicizing data and provides a uniform search interface

```
In [2]: forge = Forge()
```

```
In [3]: data = forge.search('mdf.source_name=nanomine AND mdf.resource_type=record', advanced=True)
```

```
In [4]: print('Found %d records in NanoMine'%len(data))
```

Found 227 records in NanoMine

### Get Records with Olefin Matrices

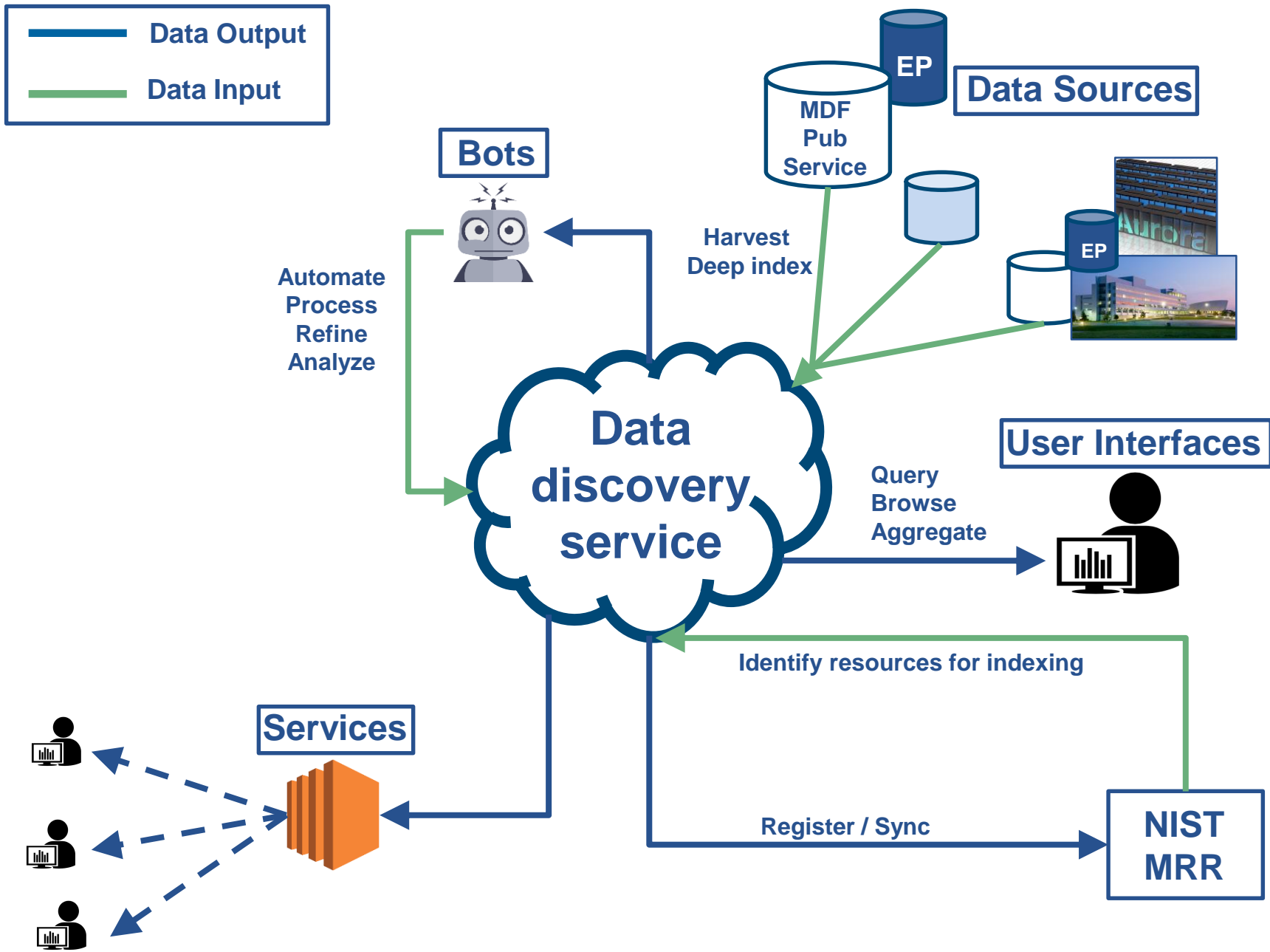
Example of a more-complex query

```
In [5]: data = forge.search('mdf.source_name=nanomine AND '
                             'content.PolymerNanocomposite.MatrixComponent.ChemicalName=olefin', advanced=True)
```

```
In [6]: print('Found %d olefin records'%len(data))
```

Found 6 olefin records

# MDF data discovery ecosystem



# Summary

## Three Major Components of Materials Data Facility

### 1. Globus

- High speed data transfer
- Easy data sharing

### 2. Data Publication Service

- Simple data publication, from your own
- Free data publication

### 3. Data Discovery Service

- Single search engine for many materials databases
- Python API for accessing these databases

# Thanks to our sponsors!



U.S. DEPARTMENT OF  
**ENERGY**



THE UNIVERSITY OF  
**CHICAGO**