# Building the
# Model Research Data Portal

## Ian Foster

foster@uchicago.edu
with Kyle Chard, Eli Dart, Steve
Tuecke, Jason Williams

globus

# The Modern Research Data Portal: A Design Pattern for Networked, Data-Intensive Science

**https://docs.globus.org/mrdp**

The Modern Research Data Portal is a new design pattern for providing secure, scalable, and high performance access to research data.

## GitHub Repo

provides code for the simple data portal that you can experiment with online

LEARN MORE ⟩

## Example Data Portal

allows you to experiment with an example implementation of the design pattern

LEARN MORE ⟩

## Code Walkthrough

provides a narrative description of the simple data portal code
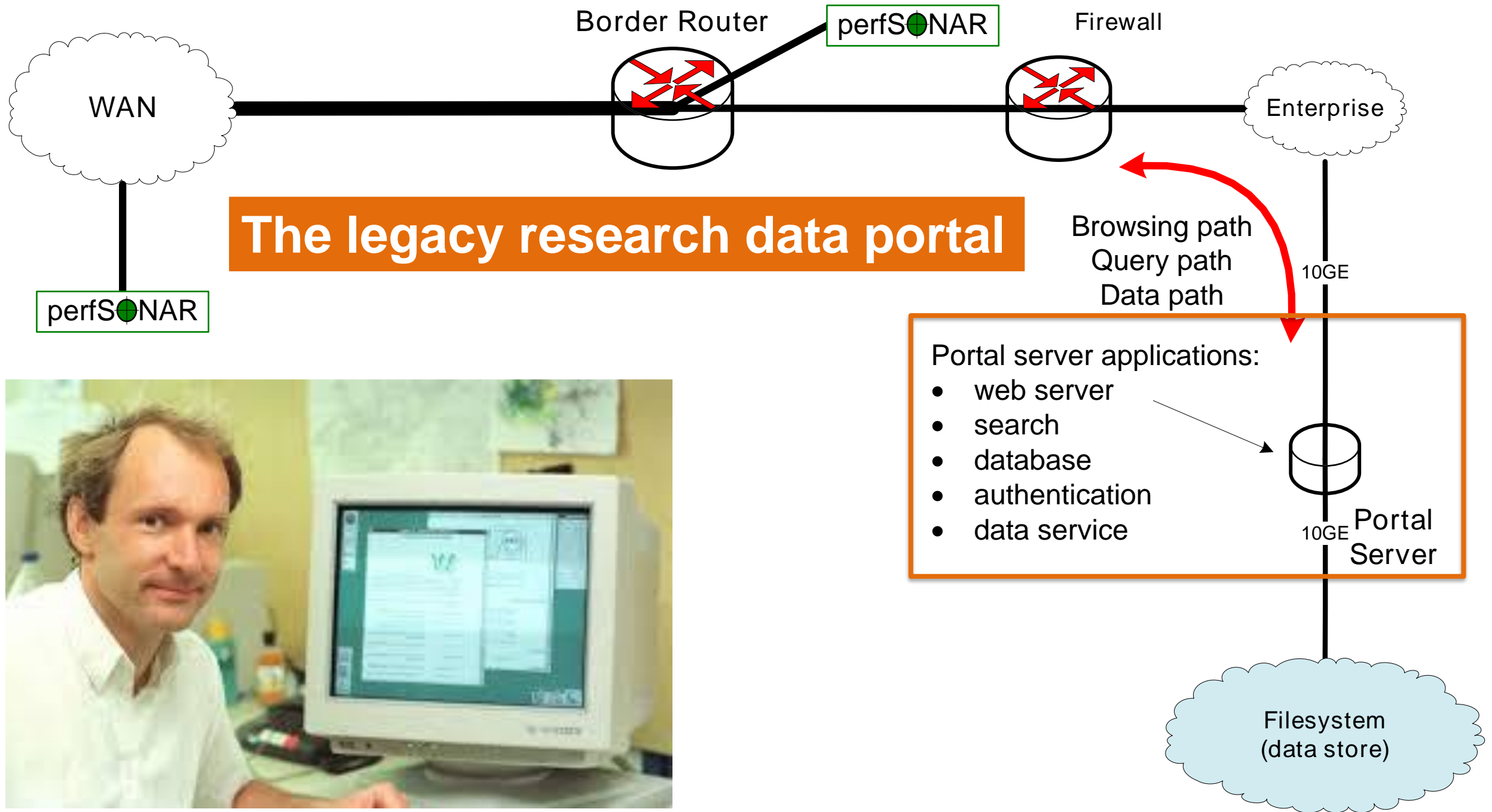
LEARN MORE ⟩

## Jupyter Notebook

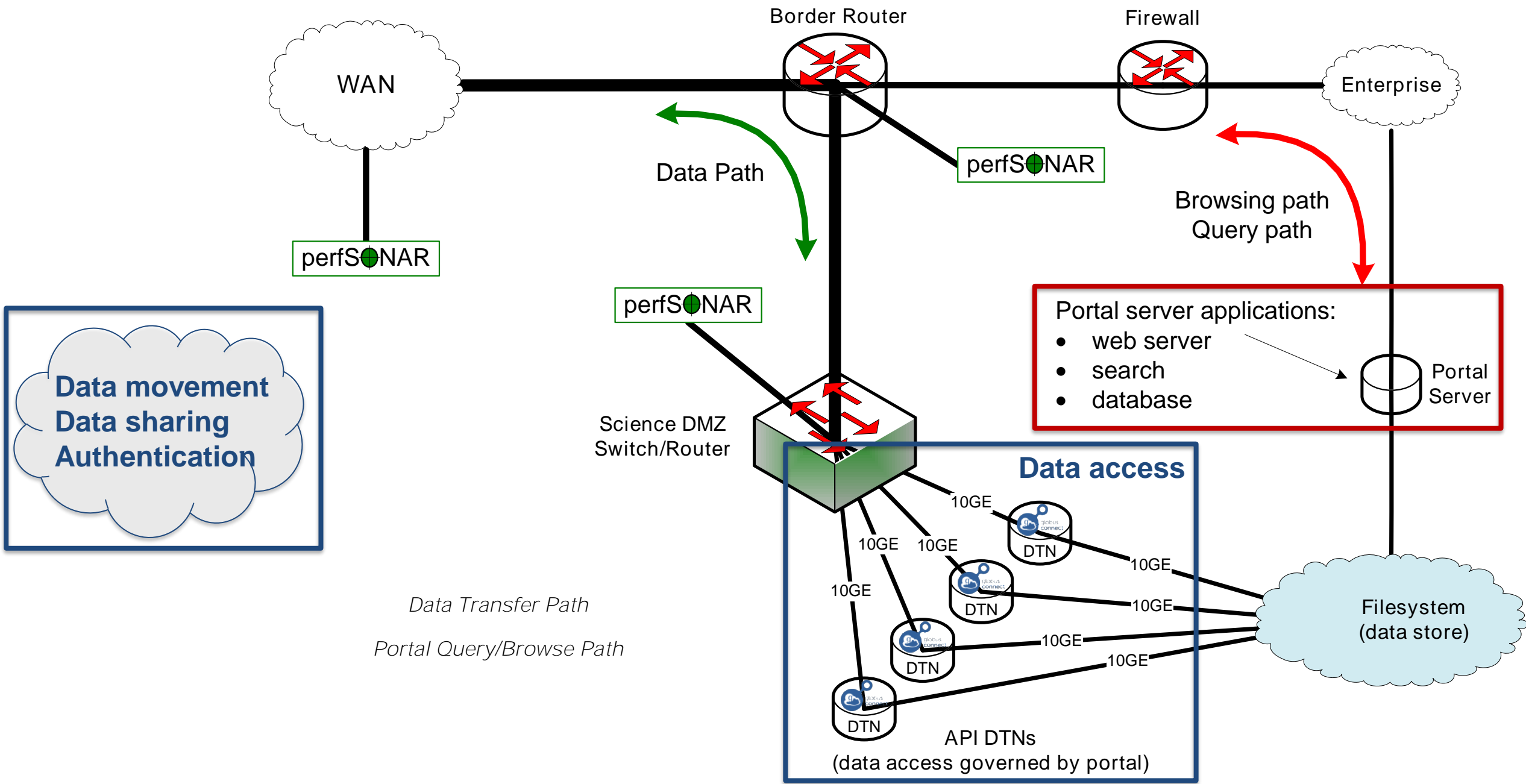demonstrates some Globus features described in the technical article

LEARN MORE ⟩

A technical article describes this design pattern, reviews representative examples at research laboratories and universities (see below), and uses coding examples to demonstrate how Globus APIs can be used to implement a range of research data portal capabilities.

Chard K, Dart E, Foster I, Shifflett D, Tuecke S, Williams J. (2017) The Modern Research Data Portal: A design pattern for networked, data-intensive science. PeerJ Preprints5:e3194v1 https://doi.org/10.7287/peerj.preprints.3194v1

LEARN MORE ⟩

Border Router    perfSONAR    Firewall

WAN

Enterprise

perfSONAR

**The legacy research data portal**

Browsing path
Query path
Data path

10GE

Portal server applications:
- web server
- search
- database
- authentication
- data service

10GE  Portal
Server

Filesystem
(data store)

# The modern research data portal



Border Router

Firewall

WAN

Enterprise

perfSONAR

Data Path

Browsing path
Query path

perfSONAR

perfSONAR

**Data movement
Data sharing
Authentication**

Portal server applications:
- web server
- search
- database

Portal
Server

Science DMZ
Switch/Router

**Data access**

10GE

10GE

10GE

10GE

DTN

10GE

10GE

DTN

10GE

Filesystem
(data store)

*Data Transfer Path*

DTN

10GE

*Portal Query/Browse Path*

DTN

10GE

DTN

API DTNs
(data access governed by portal)

# A key message: Outsource all that you can

- **Outsource responsibility for <u>determining user identities</u>**

- **Outsource <u>control over who can access</u> different data and services within the portal**

- **Outsource responsibility for <u>managing data uploads and downloads</u> between various locations and storage systems**

- **Leverage <u>standard web user interfaces</u> for common user actions**
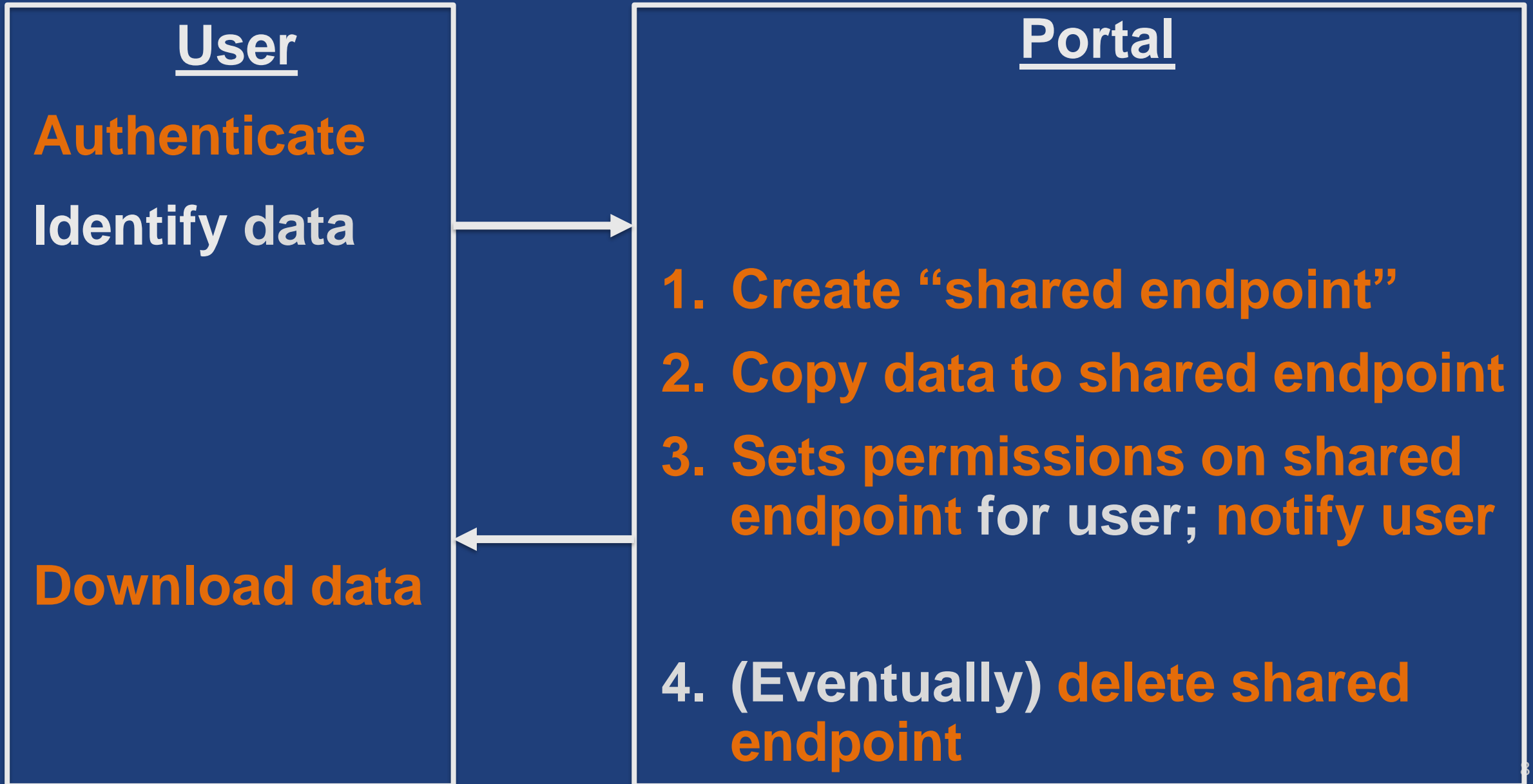
# The modern research data portal



Border Router

Firewall

WAN

Enterprise

perfS●NAR

Data Path

perfS●NAR

Browsing path
Query path

perfS●NAR

Portal server applications:
- web server
- search
- database

Portal
Server

Science DMZ
Switch/Router

**Data access**

10GE

10GE    10GE

DTN

10GE

DTN

10GE

10GE

DTN

10GE

DTN

10GE

10GE

DTN

*Data Transfer Path*

*Portal Query/Browse Path*

Filesystem
(data store)

API DTNs
(data access governed by portal)

**Globus cloud services**

Globus web widgets

Globus Auth

Globus transfer

Web browser

HTTPS

REST

Portal web server (Client)

Firewall

REST

Other services

User's endpoint (optional)

GridFTP

Portal endpoint

Science DMZ

Other endpoints

# A simple example of MDRDP logic

## User

**Authenticate**

**Identify data**

**Download data**

## Portal

1. **Create "shared endpoint"**
2. **Copy data to shared endpoint**
3. **Sets permissions on shared endpoint** for user; **notify user**

4. **(Eventually)** **delete shared endpoint**

1. **Create "shared endpoint"**

2. **Copy data to shared endpoint**

3. **Set permissions on shared endpoint for user; notify user**

**...**

4. **(Eventually) delete shared endpoint**

```python
from globus_sdk import TransferClient, TransferData
from globus_sdk import AuthClient
import sys, random, uuid


def rdp(host_id,         # Endpoint for shared endpoint
        source_path,     # Directory to copy data from
        email):          # Email address to share with
    tc = TransferClient()
    ac = AuthClient()
    tc.endpoint_autoactivate(host_id)

    # (1) Create shared endpoint:
    # (a) Create directory to be shared
    share_path = '/~/' + str(uuid.uuid4()) + '/'
    tc.operation_mkdir(host_id, path=share_path)
    # (b) Create shared endpoint on directory
    shared_ep_data = {
        'DATA_TYPE': 'shared_endpoint',
        'host_endpoint': host_id,
        'host_path': share_path,
        'display_name': 'RDP shared endpoint',
        'description': 'RDP shared endpoint'
    }
    r = tc.create_shared_endpoint(shared_ep_data)
    share_id = r['id']
```

Connect to storage system

1. **Create "shared endpoint"**

2. **Copy data to shared endpoint**

3. **Set permissions on shared endpoint for user; notify user**

…

4. **(Eventually) delete shared endpoint**

```python
# (2) Copy data into the shared endpoint
tc.endpoint_autoactivate(share_id)
tdata = TransferData(tc, host_id, share_id,
        label='RDP copy', sync_level='checksum')
tdata.add_item(source_path, '/', recursive=True)
r = tc.submit_transfer(tdata)
tc.task_wait(r['task_id'], timeout=1000,
                polling_interval=10)

# (3) Enable access by user
r = ac.get_identities(usernames=email)
user_id = r['identities'][0]['id']
rule_data = {
    'DATA_TYPE': 'access',
    'principal_type': 'identity',   # Grantee is
    'principal': user_id,           #  a user.
    'path': '/',                    # Path is /
    'permissions': 'r',             # Read-only
    'notify_email': email,          # Email invite
    'notify_message':               # Invite msg
        'Requested data are available.'
}
tc.add_endpoint_acl_rule(share_id, rule_data)

# (4) Ultimately, delete the shared endpoint
tc.delete_endpoint(share_id)
```

1. **Create "shared endpoint"**

2. **Copy data to shared endpoint**

3. **Set permissions on shared endpoint for user; notify user**

...

4. **(Eventually) delete shared endpoint**

```python
# (2) Copy data into the shared endpoint
tc.endpoint_autoactivate(share_id)
tdata = TransferData(tc, host_id, share_id,
      label='RDP copy', sync_level='checksum')
tdata.add_item(source_path, '/', recursive=True)
r = tc.submit_transfer(tdata)
tc.task_wait(r['task_id'], timeout=1000,
                polling_interval=10)

# (3) Enable access by user
r = ac.get_identities(usernames=email)
user_id = r['identities'][0]['id']
rule_data = {
    'DATA_TYPE': 'access',
    'principal_type': 'identity',  # Grantee is
    'principal': user_id,          #  a user.
    'path': '/',                   # Path is /
    'permissions': 'r',            # Read-only
    'notify_email': email,         # Email invite
    'notify_message':              # Invite msg
        'Requested data are available.'
}
tc.add_endpoint_acl_rule(share_id, rule_data)

# (4) Ultimately, delete the shared endpoint
tc.delete_endpoint(share_id)
```

1. **Create "shared endpoint"**

2. **Copy data to shared endpoint**

3. **Set permissions on shared endpoint for user; notify user**

...

4. **(Eventually) delete shared endpoint**

```python
# (2) Copy data into the shared endpoint
tc.endpoint_autoactivate(share_id)
tdata = TransferData(tc, host_id, share_id,
        label='RDP copy', sync_level='checksum')
tdata.add_item(source_path, '/', recursive=True)
r = tc.submit_transfer(tdata)
tc.task_wait(r['task_id'], timeout=1000,
                polling_interval=10)

# (3) Enable access by user
r = ac.get_identities(usernames=email)
user_id = r['identities'][0]['id']
rule_data = {
    'DATA_TYPE': 'access',
    'principal_type': 'identity',   # Grantee is
    'principal': user_id,           #  a user.
    'path': '/',                    # Path is /
    'permissions': 'r',             # Read-only
    'notify_email': email,          # Email invite
    'notify_message':               # Invite msg
        'Requested data are available.'
}
tc.add_endpoint_acl_rule(share_id, rule_data)

# (4) Ultimately, delete the shared endpoint
tc.delete_endpoint(share_id)
```

# An example MRDP
## https://docs.globus.org/mrdp

# Many variants possible

- **Manage access to data at multiple locations**

- **Manage access to data on cloud**

- **Upload data for analysis**

- **Data download from scientific instruments**

- **Data publication**

- **Transfer data to computer for analysis**
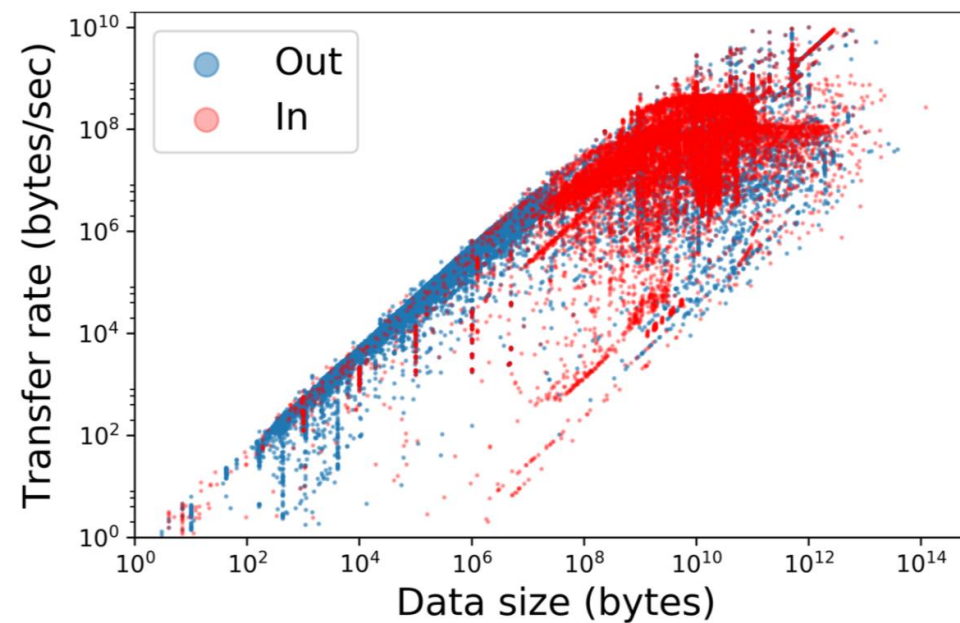
**(a)** RDA

**(b)** RDA to NERSC subset

**(c)** Sanger

**(d)** Petrel

# The Modern Research Data Portal: A Design Pattern for Networked, Data-Intensive Science

**https://docs.globus.org/mrdp**

The Modern Research Data Portal is a new design pattern for providing secure, scalable, and high performance access to research data.

## GitHub Repo

provides code for the simple data portal that you can experiment with online

LEARN MORE ⟩

## Example Data Portal

allows you to experiment with an example implementation of the design pattern

LEARN MORE ⟩

## Code Walkthrough

provides a narrative description of the simple data portal code

LEARN MORE ⟩

## Jupyter Notebook

demonstrates some Globus features described in the technical article

LEARN MORE ⟩

A technical article describes this design pattern, reviews representative examples at research laboratories and universities (see below), and uses coding examples to demonstrate how Globus APIs can be used to implement a range of research data portal capabilities.

Chard K, Dart E, Foster I, Shifflett D, Tuecke S, Williams J. (2017) The Modern Research Data Portal: A design pattern for networked, data-intensive science. PeerJ Preprints5:e3194v1
https://doi.org/10.7287/peerj.preprints.3194v1

LEARN MORE ⟩

**foster@uchicago.edu**