# The Data Mobility Exhibition

Eli Dart, Science Engagement

Energy Sciences Network (ESnet)

Lawrence Berkeley National Laboratory

CC* PI meeting

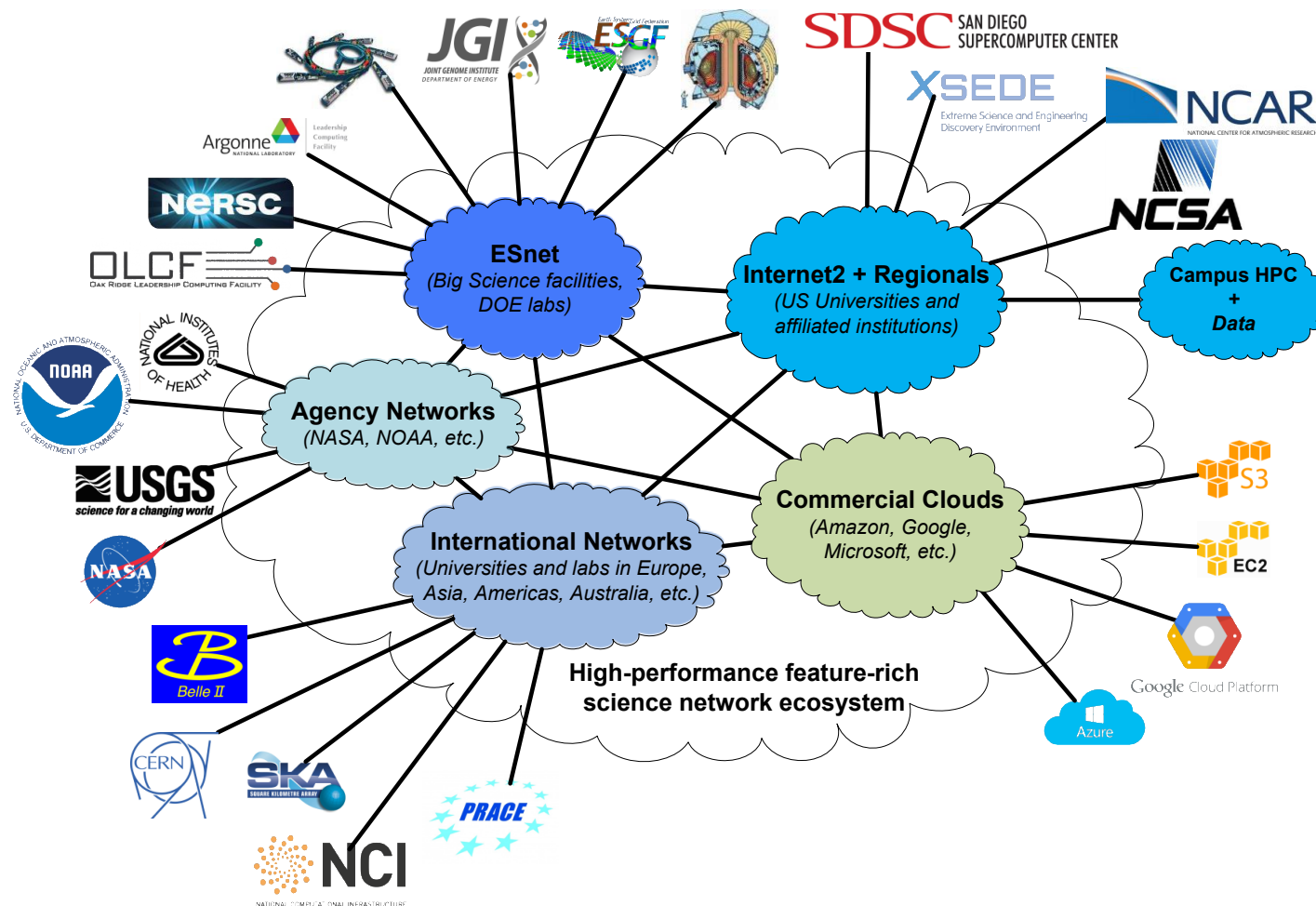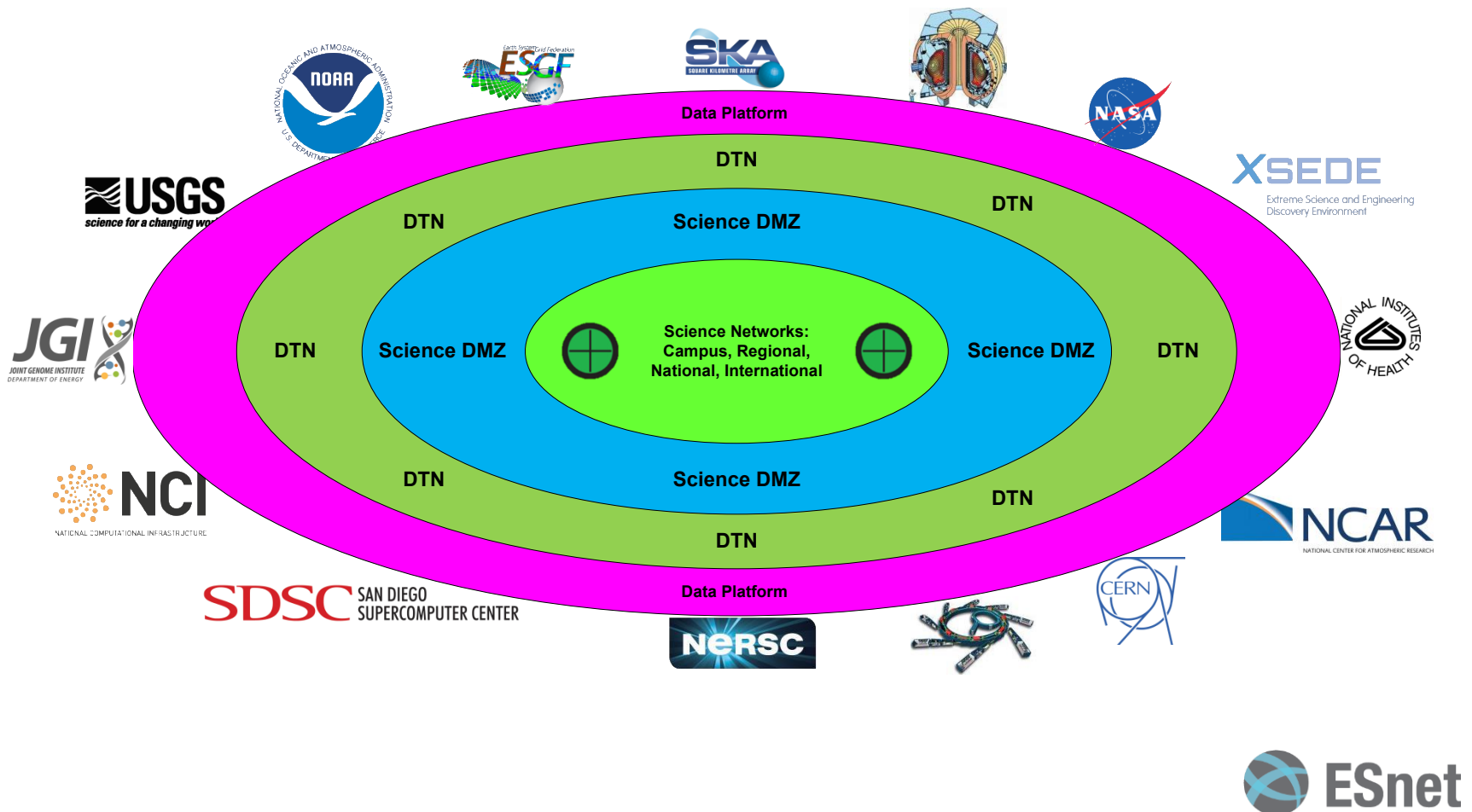Minneapolis, MN

September 24, 2019

# Our Community Has Made Great Progress

- Over the past 8-10 years or so, we have made great strides forward

- Science networks are big, fast, and clean

- Science networks are instrumented for performance

- The Science DMZ model is widely deployed

- DTNs in the Science DMZs

- Data orchestration platforms running on the DTNs

**ESnet**

# Data Ecosystem - Abstract Network Diagram

# Data Ecosystem – Non-Network View

# What Remains To Be Done?

- We're getting close to completing a transformation of Cyberinfrastructure
- Three major tasks remain
  - Data orchestration across Science DMZs
    - This includes test, verification, and performance engineering
    - Partially complete today
  - Upgrading the data portals - Modern Research Data Portal
    - This has begun – lots left to do
  - Onboarding scientists and collaborations
    - Science Engagement
    - We understand it, but we need to scale it
- Remember – this has to be useful to scientists, so it has to work for them

ESnet

# Data Movement Exhibition

- Current and previous CC* Awardees, along with the greater R&E community, are encouraged to participate

- Will be highlighted at the 2019 and 2020 CC* PI Meetings

- Using reference data sets, and existing campus CI components, participants will work on scientific data movement capabilities:
  - Download/transfer data
  - Measure performance
  - Potentially improve

- This event has begun (August 2019), and will extend for a full year (through the CC* 2020 PI meeting)

**ESnet**

# Data Movement Exhibition

- Basic Idea:
  - Create a brief (1-2 page) description of the network and data architecture for your campus environment
  - Prepare a local data transfer machine, and get Globus working
    - Hopefully this is a production DTN in your Science DMZ
    - Even better: front-end for campus HPC or other production CI
    - Even better: DTNs for a data portal – can we download from you?
  - Download the reference data sets (specific instructions forthcoming)
  - Share the description, and your results (specific instructions forthcoming)
  - For those that want to accelerate their results, 1:1 assistance, via the Engagement and Performance Operations Center (EPOC), is available: epoc@iu.edu

ESnet

# National Facility Endpoints (Petascale DTN)
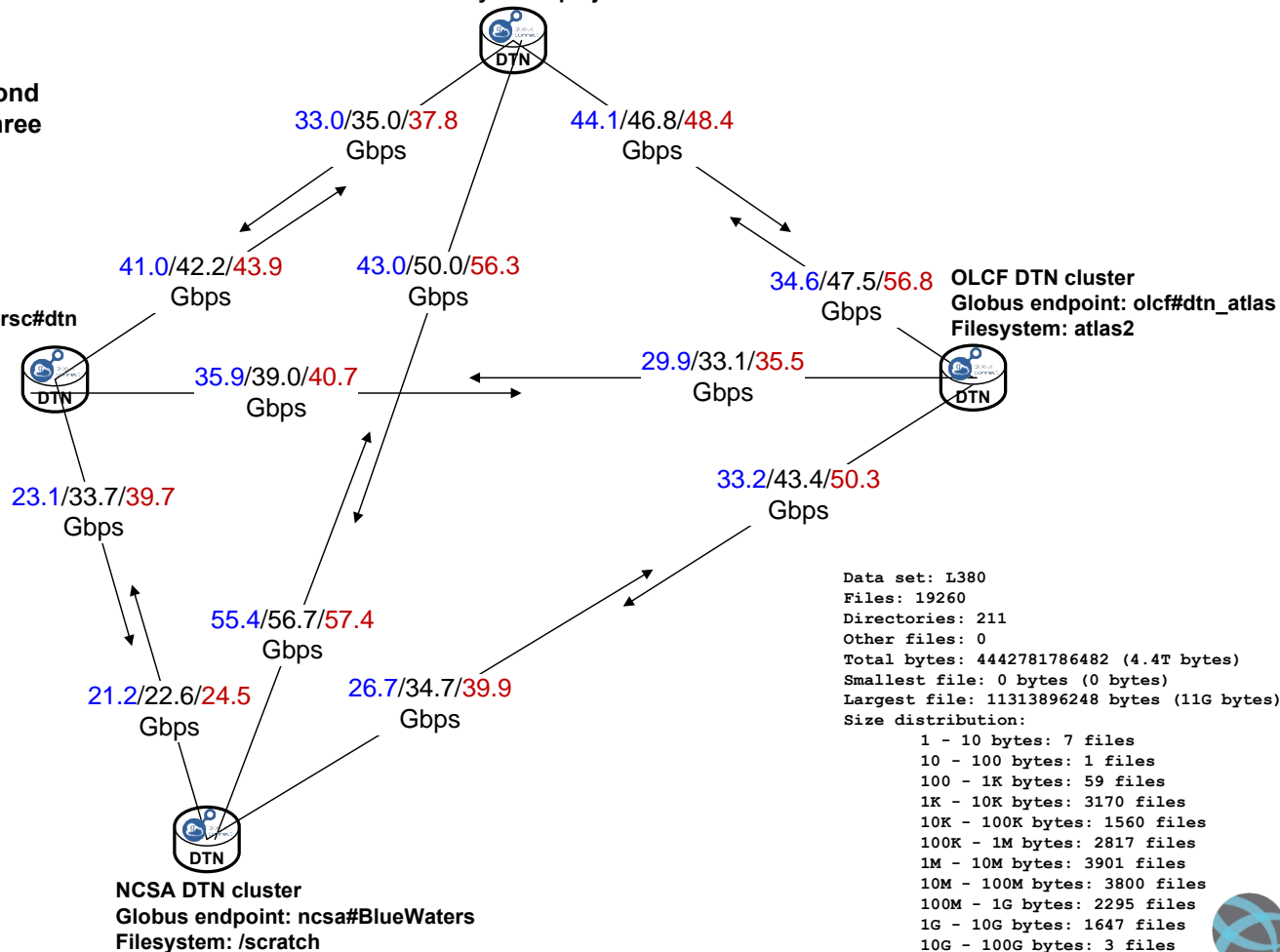
**Petascale DTN Project**

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

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

**ALCF DTN cluster**
**Globus endpoint: alcf#dtn_mira**
**Filesystem: /projects**

33.0/35.0/37.8 Gbps

44.1/46.8/48.4 Gbps

**NERSC DTN cluster**
**Globus endpoint: nersc#dtn**
**Filesystem: /project**

41.0/42.2/43.9 Gbps

43.0/50.0/56.3 Gbps

34.6/47.5/56.8 Gbps

**OLCF DTN cluster**
**Globus endpoint: olcf#dtn_atlas**
**Filesystem: atlas2**

35.9/39.0/40.7 Gbps

29.9/33.1/35.5 Gbps

23.1/33.7/39.7 Gbps

33.2/43.4/50.3 Gbps

55.4/56.7/57.4 Gbps

21.2/22.6/24.5 Gbps

26.7/34.7/39.9 Gbps

```
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
```

**NCSA DTN cluster**
**Globus endpoint: ncsa#BlueWaters**
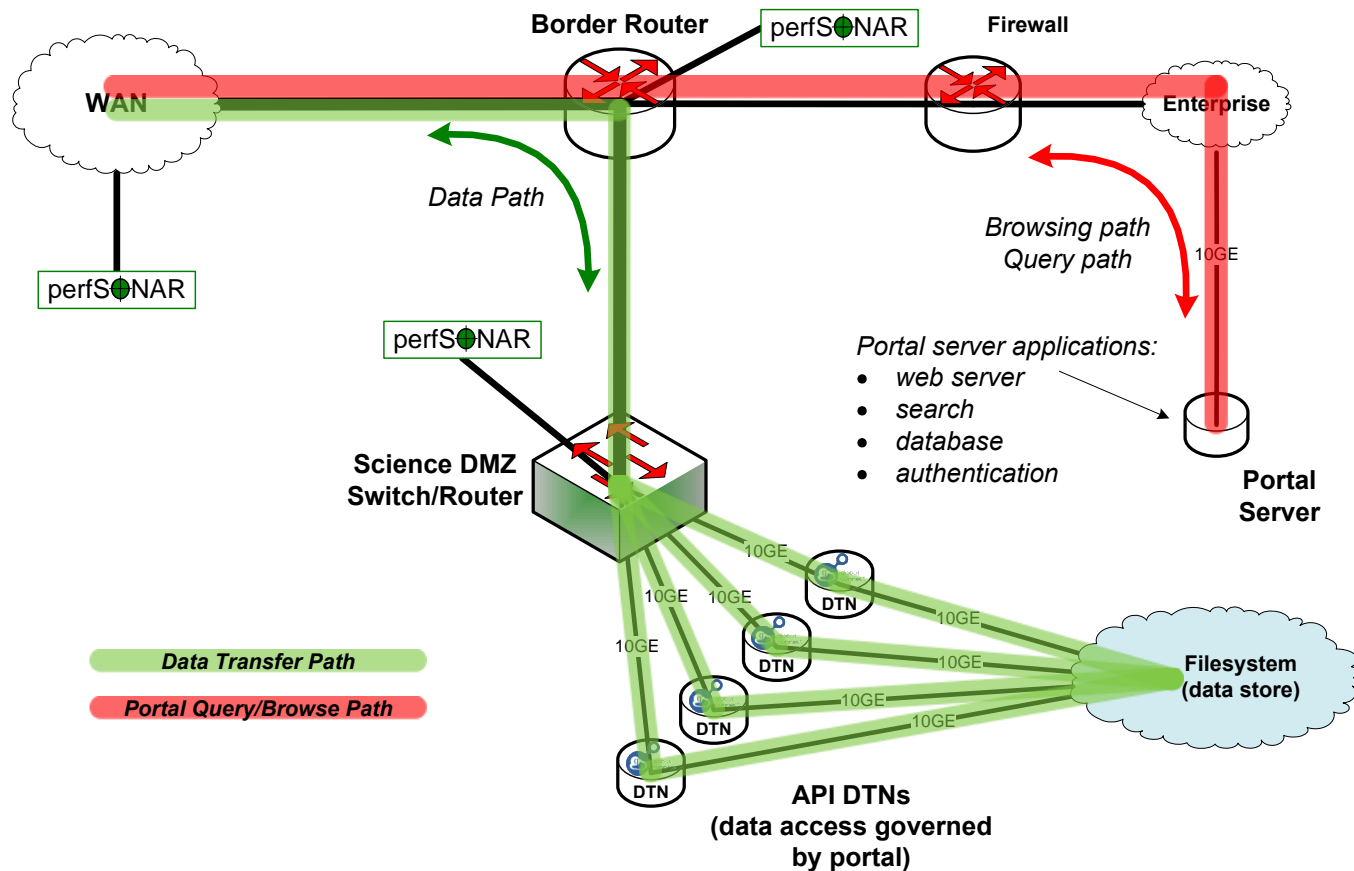**Filesystem: /scratch**

**ESnet**

# Science Data Portals

- Large repositories of scientific data
  - Climate data
  - Sky surveys (astronomy, cosmology)
  - Many others
  - Data search, browsing, access
- Make the data easily accessible on HPC platforms
  - Supercomputers
  - Campus HPC
  - Clouds
- This will feed the rising capabilities of AI/ML and other data analytics

**ESnet**

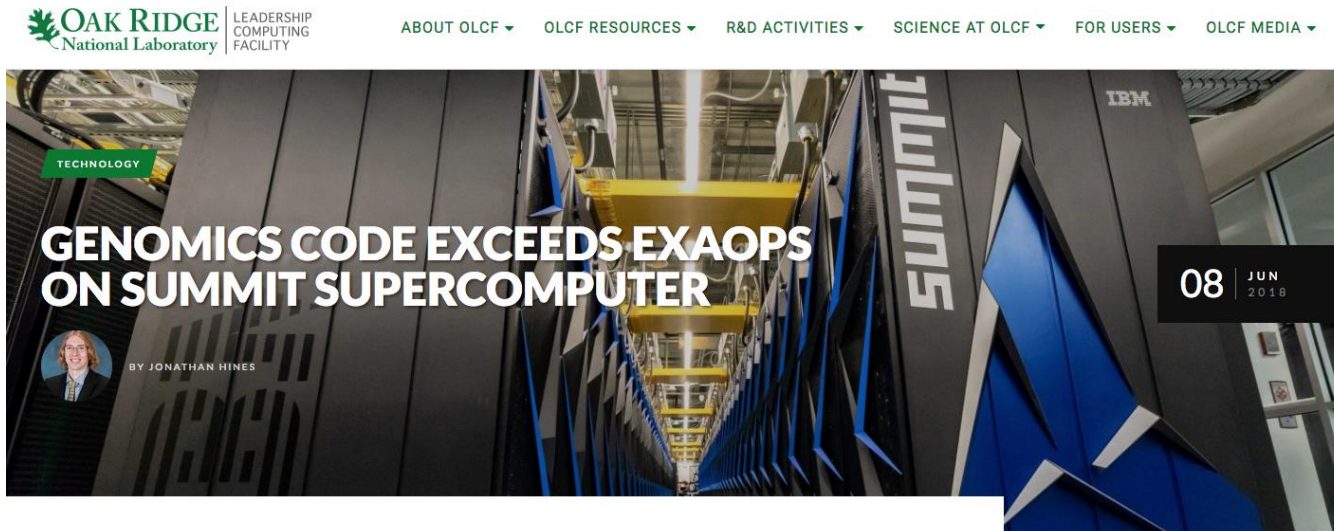# Legacy Portal Design



**Border Router**  perfSONAR  **Firewall**

**WAN**

**Enterprise**

perfSONAR

*Browsing path*
*Query path*
*Data path*

10GE

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

**Portal Server**

10GE

**Filesystem
(data store)**

ESnet

# Modern Research Data Portal Leverages Science DMZ



https://peerj.com/articles/cs-144/

# Science at Scale: Genomics (June 2018)

# Science at Scale: Climate (August 2017)





Figure 3: Sample images of atmospheric rivers correctly classified (true positive) by our deep CNN model. Figure shows total column water vapor (color map) and land sea boundary (solid line).

# They Can Use All The Data

- Groups like these need large data sets
- Much of the data in many fields is behind legacy portals
  - Significant human effort to retrieve what scientists need
  - Legacy systems perform poorly, especially at scale
- Legacy data portals are a product of their time
  - Remember: these were designed to serve small data to small systems
  - We now live in the future from the perspective of those designs
  - Current systems far exceed the capabilities available 15 years ago
  - From the perspective of today's systems, legacy portals are products of a bygone past
- It is now perfectly reasonable for a scientist to want all the data
  - Machine learning + HPC
  - But this only works if the scientists can get to the data at scale

ESnet

# Community DTNs

- Initial tests with some others in the community
- University, Data Portal, and Supercomputer Center endpoints
- We need to do a lot more of this



**NCAR RDA**

**OHSU Portland**

**Talapas University of Oregon**

**Stampede2 TACC**

24Gbps

12.8Gbps    10.3Gbps

9.0Gbps    8.0Gbps

9.3Gbps    7.1Gbps

**Community endpoints to NERSC DTN endpoint**
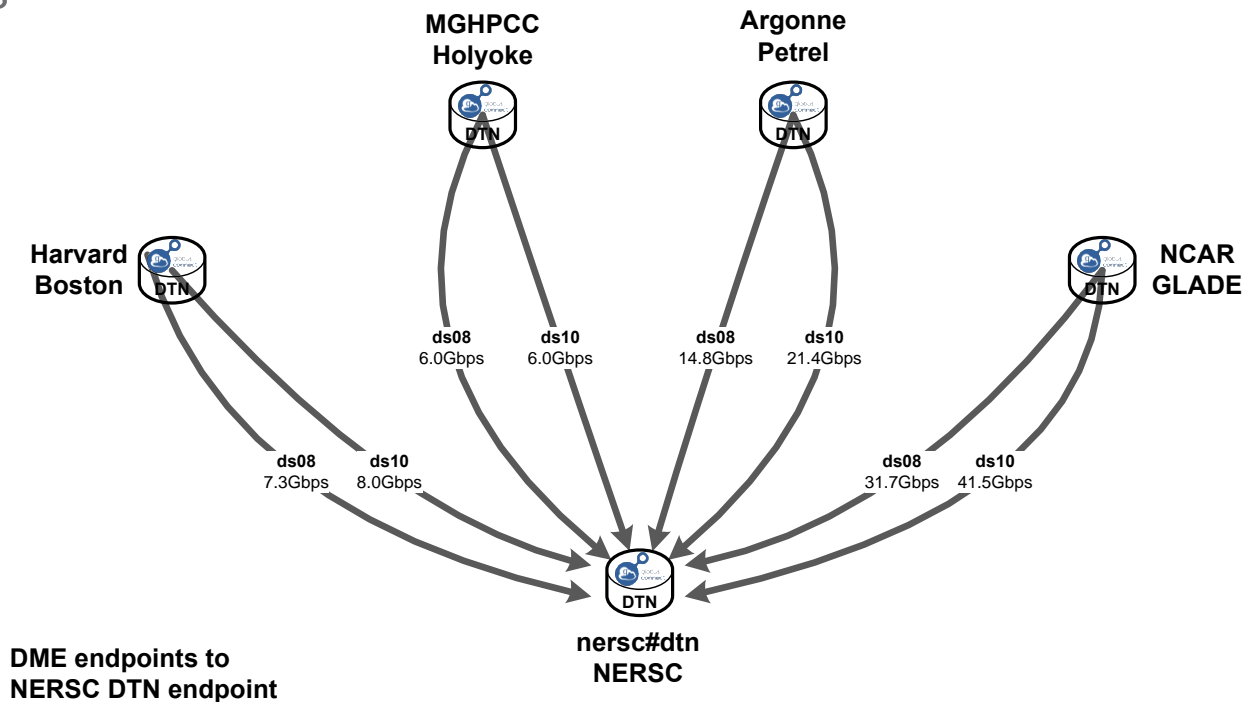
**nersc#dtn NERSC**

# Data Mobility Exhibition

- Expansion of the Petascale DTN effort

- Goal is to dramatically increase the number of endpoints known to interoperate well and at high performance

- Ensure campuses can easily exchange data with national facilities

- This is the Data Mobility Exhibition
  - http://fasterdata.es.net/performance-testing/2019-2020-data-mobility-workshop-and-exhibition/2019-2020-data-mobility-exhibition/
  - Data sets and endpoints are described here:
    - https://www.globusworld.org/tour/data-mobility-exhibition

- Why the term "Exhibition?"
  - The point is to interoperate, not to compete
  - Everyone has different mission, constraints, funding, requirements
  - But everyone has users, and they need this stuff to work well

**ESnet**

# DME Endpoints – Initial Work

- Initial tests with Data Mobility Exhibition endpoints

- Let's expand this – we can work together to achieve something amazing



**MGHPCC Holyoke**

**Argonne Petrel**

**Harvard Boston**

**NCAR GLADE**

| | ds08 | ds10 |
|---|---|---|
| MGHPCC | 6.0Gbps | 6.0Gbps |
| Argonne | 14.8Gbps | 21.4Gbps |

**ds08** 7.3Gbps  **ds10** 8.0Gbps

**ds08** 31.7Gbps  **ds10** 41.5Gbps

**nersc#dtn NERSC**

**DME endpoints to NERSC DTN endpoint**

**ds08 data set: 1TB (~30k files, 1MB to 10GB each)**
**ds10 data set: 1TB (100 files, 10GB each)**

# Vision – Interoperable Computing And Data

# A Call To Action

- We are almost there!

- By ensuring high-performance interoperability, we prove that our systems are well-integrated and ready for use by the scientific community

- In aggregate, these systems will form a data infrastructure across US R&E

- Please participate in the data mobility exhibition
  - [http://fasterdata.es.net/performance-testing/2019-2020-data-mobility-workshop-and-exhibition/2019-2020-data-mobility-exhibition/](http://fasterdata.es.net/performance-testing/2019-2020-data-mobility-workshop-and-exhibition/2019-2020-data-mobility-exhibition/)

- `We will come back in a year with results

**ESnet**

# Thanks!

Eli Dart dart@es.net

Energy Sciences Network (ESnet)

Lawrence Berkeley National Laboratory

engage@es.net

http://my.es.net/

http://www.es.net/

http://fasterdata.es.net/