# Modern Cyberinfrastructure: The Ladder to the Shoulders Of Giants

Eli Dart, Science Engagement

Energy Sciences Network (ESnet)

Lawrence Berkeley National Laboratory

PEARC20

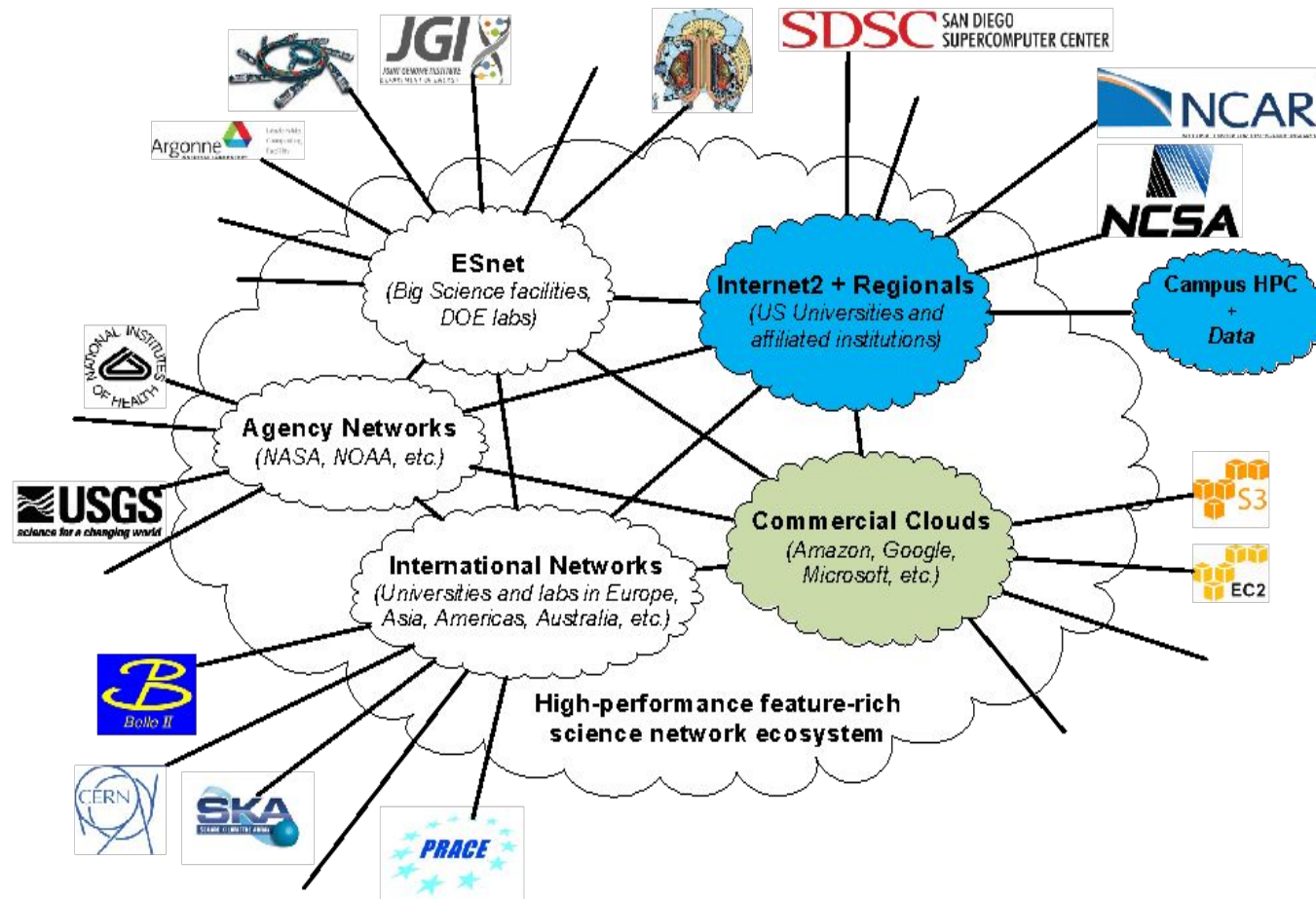Virtual (Coronapocalypse)

July 27, 2020

# 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
  – High speed regional and national networks
  – R&E exchanges
  – Campus networks
  – International connectivity
- Science networks are instrumented for performance
  – perfSONAR
  – Critical for ensuring correct operation
  – Invaluable for timely resolution of problems
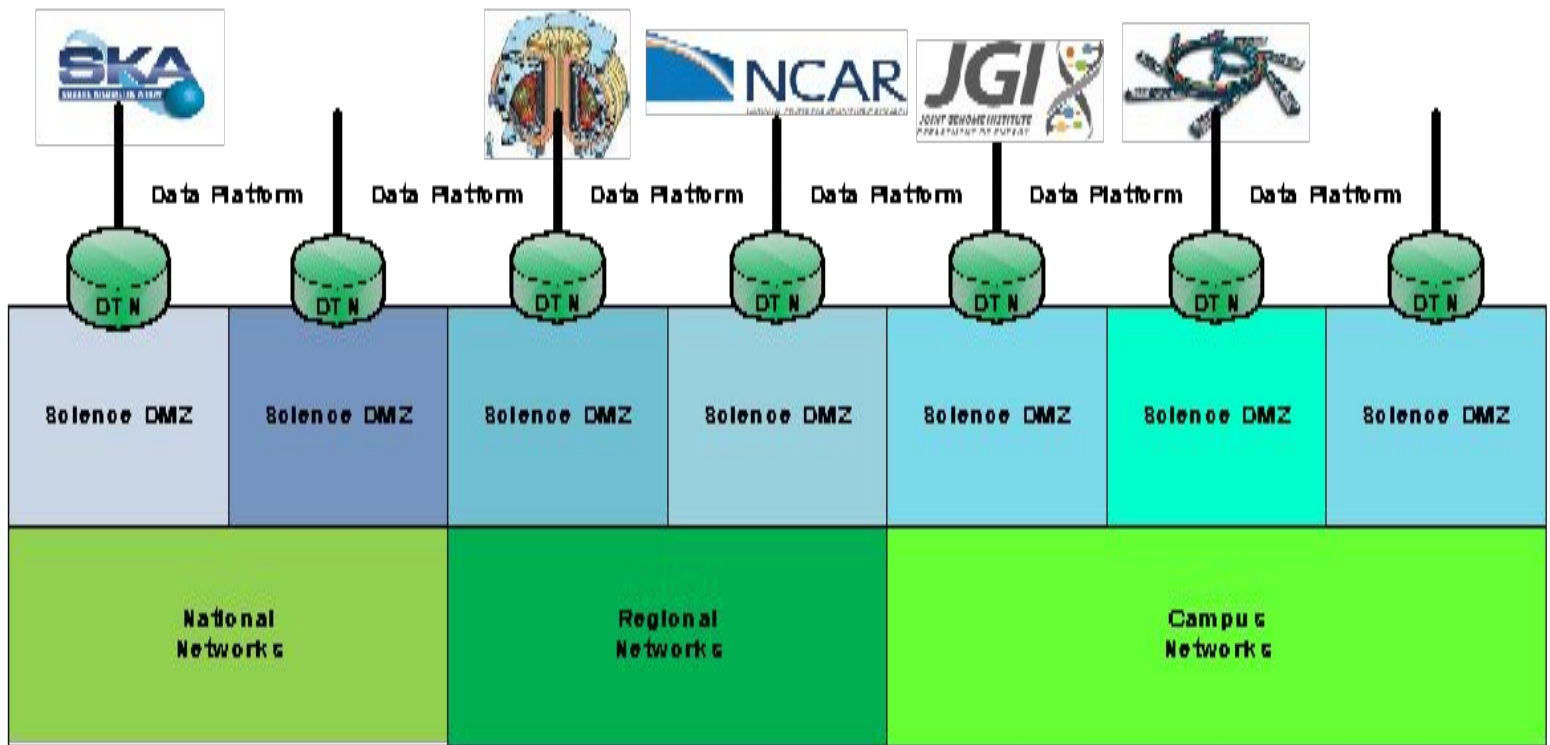
**ESnet**

# Our Community Has Made Great Progress

- The Science DMZ model is widely deployed
  - Campuses, laboratories, experiments
  - HPC facilities
  - Some data portals (more on this later)
- DTNs in the Science DMZs
  - Connect storage to high speed networks
  - HPC filesystems
  - Experiment data acquisition systems
- Data orchestration platforms running on the DTNs
  - This is what the scientist sees
  - Capable platforms allow orchestration rather than clunky user-driven scripting or manual downloads
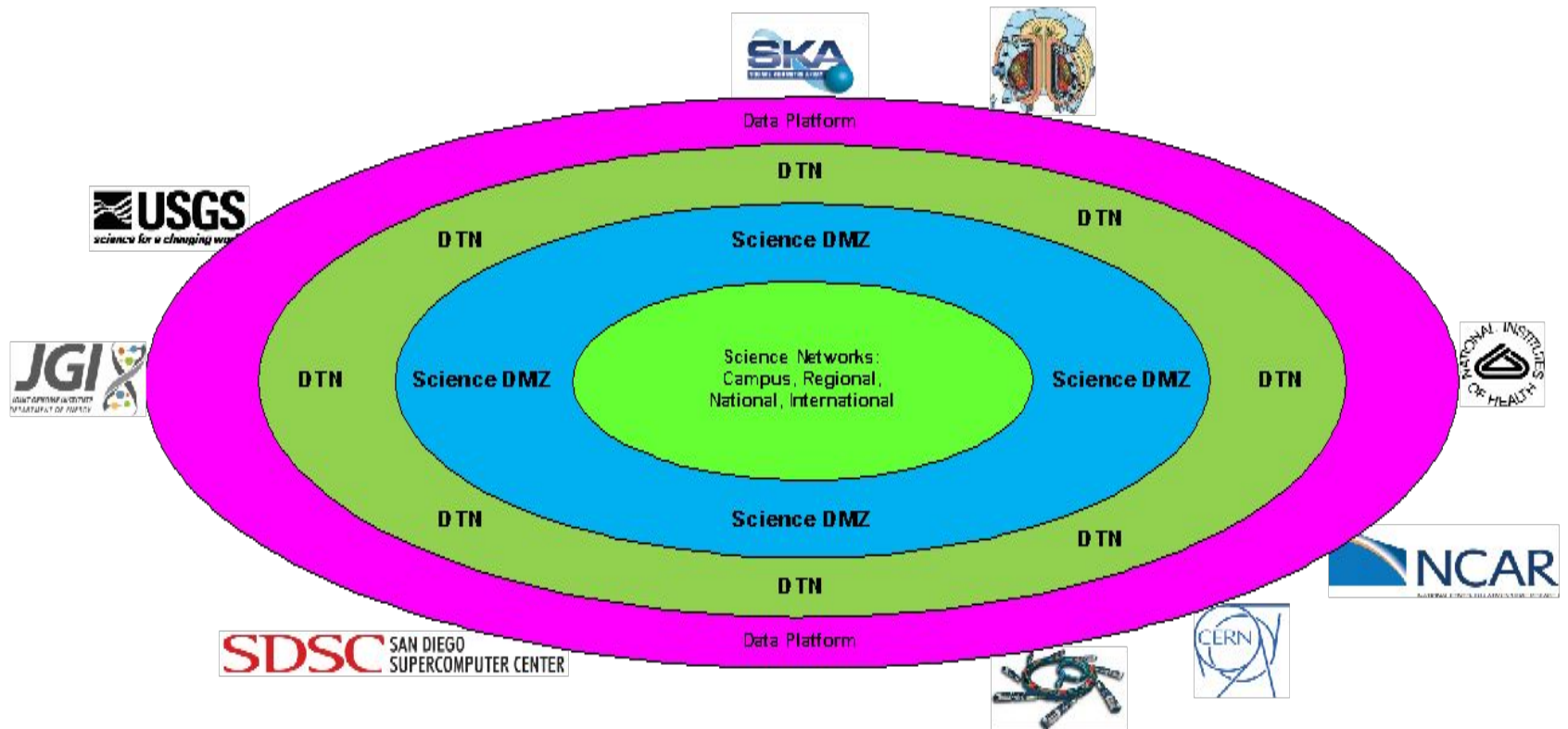
**ESnet**

# Data Ecosystem - Abstract Network Diagram

# Data Ecosystem – Block Visualization

# Data Ecosystem – Concentric View

# What Remains To Be Done?

- We aren't all the way there yet (unfortunately)

- The diagrams show a vision that is not yet fully realized

- Three major tasks remain
  - Deployment of an interoperable platform across Science DMZs
    - This includes test, verification, and performance engineering
    - Partially complete
  - Integrating the major data repositories and portals with the platform
    - 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

# Interoperable Platform Deployment

- This is partially complete

- Necessary features
  - Automation
  - Fault recovery
  - Data integrity
  - Integration with web-based portals

- Several platforms exist
  - Globus
    - Significant deployment in NSF and DOE spaces
    - Basis for examples shown here
  - XRootD (LHC experiments)
  - OSG Stack

- Key point – **the scientist must not be made the integrator**
  - If the scientist is the integrator, they will use HTTP and rsync+SSH forever
  - The old tools don't scale, but the scientists can't build the better platforms themselves
  - **WE** MUST DO THIS

ESnet

# Example Of Platform Power (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

41.0/42.2/43.9
Gbps

43.0/50.0/56.3
Gbps

34.6/47.5/56.8
Gbps

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

**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 DMZ – HPC Center DTN Cluster

# Science Data Portals

- Large repositories of scientific data
  - Climate data
  - Sky surveys (astronomy, cosmology)
  - Many others
  - Data search, browsing, access

- Many scientific data portals were designed 15+ years ago
  - Single-web-server design
  - Data browse/search, data access, user awareness all in a single system
  - All the data goes through the portal server
    - In many cases by design
    - E.g. embargo before publication (enforce access control)
  - Better than old command-line FTP, but outdated by today's standards

ESnet

# Legacy Portal Design



- Legacy portals designed to provide a tiny subset of the data

- Need to integrate legacy portals with modern data ecosystem

  – Better automation, scale, performance

  – Connectivity to HPC

ESnet

# Next-Generation Portal Leverages Science DMZ



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

ESnet

# JGI Data Portal

## Searching for Projects

- Explore what you can do here.
- Search projects/proposals using "Advanced Search" filters.

## Downloading Files

- Download over the web
- Download large number of files with Globus service.
- Download via API using scripting or programming
- Download with "Cart" by collecting projects/portals of your interest.

## Looking for Access

- Looking for data and do not have access to the private portal? Please contact PI
- How to grant access to your proposal/project/genome? Get Instructions.

## JGI Genome Portal

### New Feature: "Bulk Downloads"

collect your favorite projects and download them in **bulk** with our new feature **"CART"**. Ability to download files with Portal or via Globus.

Find out more details

## What's New

**New feature: "Download with Cart"** 🛒

A convenient way to collect projects/genomes/metagenomes of your interest and download all files associated with them in **bulk**.
Read more and provide your comments and suggestions for this feature to our team.

## My Favorites

⭐ My Favorites: New Feature - Based on Your Feedback
This feature allows to save your filtered search results to "My Favorites" and access it later.

## The "Tree of Life"

Please use our powerful search or go to the "Tree of Life" if it is the most convenient way for you to reach your genomes/projects.

ESnet

# NCAR RDA Performance to DOE HPC Facilities

- 1.5TB data set

- 1121 files

NCAR RDA
rda#datashare

DTN

13.9 Gbps    16.6 Gbps    11.9 Gbps

DTN    DTN    DTN

nersc#dtn
NERSC

alcf#dtn_mira
ALCF

olcf#dtn_atlas
OLCF

**ESnet**

# Reasons To Scale Data Portals

- Some reasons are obvious
  - Increase in size of data objects (MB ⯈ GB ⯈ 100s of GB)
  - Number of data objects (many thousands per data set)
- Other reasons are paradigm shifts
  - Modern data analysis on HPC can use a *lot* of data
  - Today's HPC facilities are far more capable than in the past
- Retrofit / rebuilt data portals and data repositories
  - Significant wins from increased data analysis

ESnet

# Science at Scale: Genomics

# Science at Scale: Climate





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 <u>All The Data</u>

- Groups like these need large data sets

- Much of the data in their field is behind legacy portals
  - Significant human effort to retrieve what they 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 <u>perfectly reasonable</u> for a scientist to want <u>all the data</u>
  - Machine learning + HPC
  - But this only works if the scientists can get to the data at scale

ESnet

# We (not the scientists) Have To Do This

- The scientific community cannot do this for themselves

- Individual researchers do not control the resources
  - Computing centers
  - Data repositories
  - Science networks
  - Our community owns these – we have to do the work

- Integration, performance engineering, interoperability

- Science Engagement to teach scientists how to use the better platforms

- This is the path forward

¯\_(ツ)_/¯

ESnet

# Networks Cannot Do This Alone

- We need a whole-community effort
  - Networks
  - HPC facilities
  - Data repositories / Data portals
  - Experimental facilities
  - Science collaborations
  - Science programs
- Networks can help, and must be part of the conversation
  - Heavy lifting is now at the network edge, in collaboration with the network core
  - Need to help them get the architecture right – we know how to do this

ESnet

# Vision – Interoperable Computing And Data

# Cyberinfrastructure Is The Ladder Up

- An integrated, interoperable cyberinfrastructure will allow scientists to make effective use of data, computing, and networks

- This is how we will achieve the advances we need in medicine, energy, climate science, and many other fields

- Large-scale data can only be effectively used if the tools work well together - otherwise the effort is too great



ESnet

Image credit: https://www.flickr.com/photos/amylovesyah/ (CC BY 2.0)

# The Path Forward

ESnet

# The Path Forward

ESnet

# Standing On The Shoulders Of Giants

- Large-scale data sets are the giants of today

- We have all the components we need to give all scientists access to all the data in their fields

- This is not a design problem
  - We have the designs, the technologies, the models
  - We know it works: we have examples

- This is an integration and deployment problem
  - We know what we need to do
  - Let's get to it!

*Isaac Newton*

**ESnet**

Image source: https://commons.wikimedia.org/wiki/File:Sir_Isaac_Newton_by_Sir_Godfrey_Kneller,_Bt.jpg
Godfrey Kneller [Public domain], via Wikimedia Commons

# ESnet
ENERGY SCIENCES NETWORK

# 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/

U.S. DEPARTMENT OF **ENERGY**
Office of Science

BERKELEY LAB

# Extra slides – data download from portal

ESnet

Home | Find Data | Ancillary Services | About/Contact | Data Citation | Web Services | For Staff

# GEOS5 Global Atmosphere Forcing Data
ds313.0 ☆

For assistance, contact Chi-Fan Shih (303-497-1833).

**Description** | **Data Access**

| | |
|---|---|
| **Help with this page:** | RDA dataset description page video tour |
| **Abstract:** | GEOS5 Atmospheric Forcing data, regridded and prepared as meteorological variables to run CESM and WRF simulations. |
| **Temporal Range:** | 2004-01-02 00:00 +0000 to 2017-10-19 21:00 +0000 (Entire dataset)<br>▸Period details by dataset product |
| **Updates:** | Irregularly |
| **Variables:** | Surface Pressure    Upper Level Winds<br>▸Variables by dataset product |
| **Vertical Levels:** | See the detailed metadata for level information |
| **Data Types:** | Grid |
| **Spatial Coverage:** | Longitude Range: Westernmost=180W Easternmost=180E<br>Latitude Range: Southernmost=90S Northernmost=90N<br>▸Detailed coverage information |
| **Data Contributors:** | UCAR/NCAR/ACD | UCAR/NCAR/CGD |
| **How to Cite This Dataset:**<br>**RIS**<br>**BibTeX** | Tilmes, S.. 2016. *GEOS5 Global Atmosphere Forcing Data*. Research Data Archive at the National Center for Atmospheric Research, Computational and Information Systems Laboratory. http://rda.ucar.edu/datasets/ds313.0/. Accessed† dd mmm yyyy.<br>†Please fill in the "Accessed" date with the day, month, and year (e.g. - 5 Aug 2011) you last accessed the data from the RDA.<br>Bibliographic citation shown in [ Federation of Earth Science Information Partners (ESIP) ♦ ] style |
| | Get a customized data citation |
| **Total Volume:** | 449.28 GB |
| **Data Formats:** | netCDF |
| **More Details:** | View more details for this dataset, including dataset citation, data contributors, and other detailed metadata |
| **Data Access:** | Click the **Data Access** tab here or in the navigation bar near the top of the page |
| **Metadata Record:** | Display in [ choose from the list ♦ ] format |

# Portal creates a Globus transfer job for us

# Submit the transfer job, go about our business

# Data Transfer from RDA Portal – Results

# Superfacility: Integrated network of experimental and computational facilities and expertise



A single interconnected "facility" where data is acquired, stored, analyzed and served

- Light Sources
- Sequence
- Telescopes
- Particle Detectors
- Microscopes

- Experimental Facilities
- Computing and Data Facilities
- ESnet
- Science Network
- Expertise — Methods, models, analytics, and software
- More Users
- More Productive Science

# Advanced Light Source Demonstration



GISAXS

HipGISAXS & RMC

ESnet

Liu et al, "Fast printing and in situ morphology …". Adv Mater. 2015