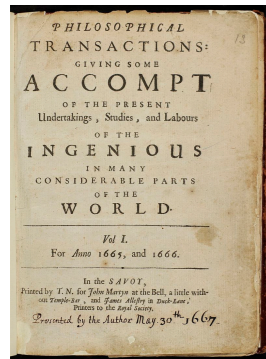


Making Data Accessible

Mercè Crosas
Chief Data Science and Technology Officer
Institute for Quantitative Social Science
Harvard University

Scholarly output doubles every 20 years from 1750s to 2000

Philosophical Transactions of the Royal Society
Nullius in verba



Increase in specialization:
every 100-150 authors =
1 new journal for 500-1000 readers



3 journals

10 journals

400 journals

14, 000 journals
(peer-reviewed)

1665

1700

1800

1900

2000

Now:

- 80,000 total journals
- 33,000 peer-reviewed

Data are part of the article in the form of visuals and tables

30% of articles with visuals; 5% have tables

30% with visuals; use of tables increases

50% with visuals, ~7 each article; integrated with tables

90% with figures and tables; occupy 25% of space

1665

1700

1800

1900

2000

First line graphs and bar charts: Playfair (1786)

First pie chart: Playfair (1801)

First scatterplots: Hershel (1833), Galton (1896)

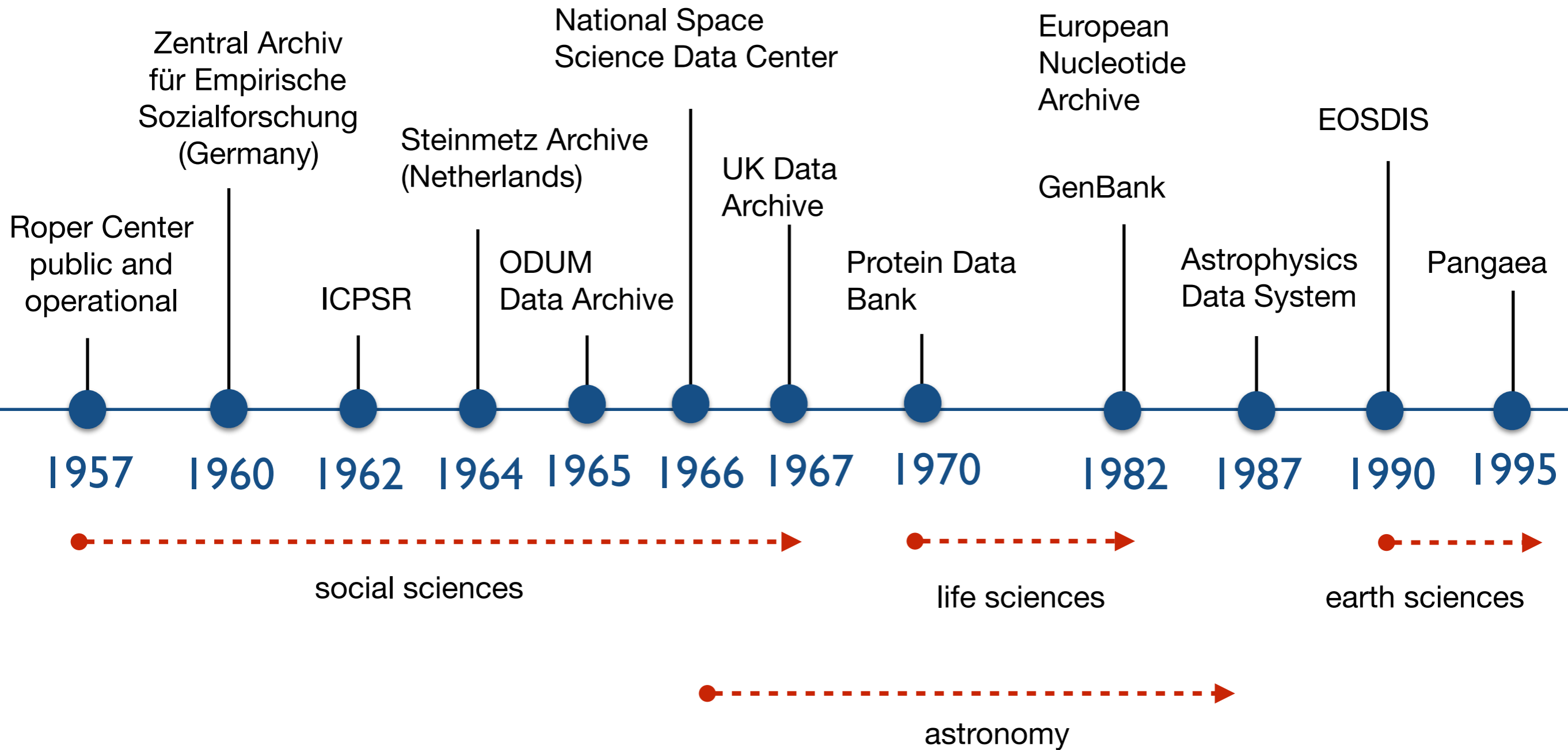
Data visuals evolve from illustrations to scientific arguments

Source: Gross Harmon, Reidy, Communicating Science, 2002

Science communication adapts to the increase in cognitive complexity

- Gross, Herman and Reidy (2002):
 - analyze the style, presentation, and argument of a sample of articles from 17th to 20th century
 - and argue that during that time we developed devices for more efficient communication to compensate the increase of cognitive complexity
- The increase in quantity and complexity of scholarly output is accompanied with an increase in research data
- I argue that in the last decade **data publishing** is born as a new necessary device for more efficient communication of science

Emergence of (domain-specific) digital data archives



Data publishing in the hands of the researcher

Dataverse

Dryad

Figshare

Zenodo

2006

2009

2011

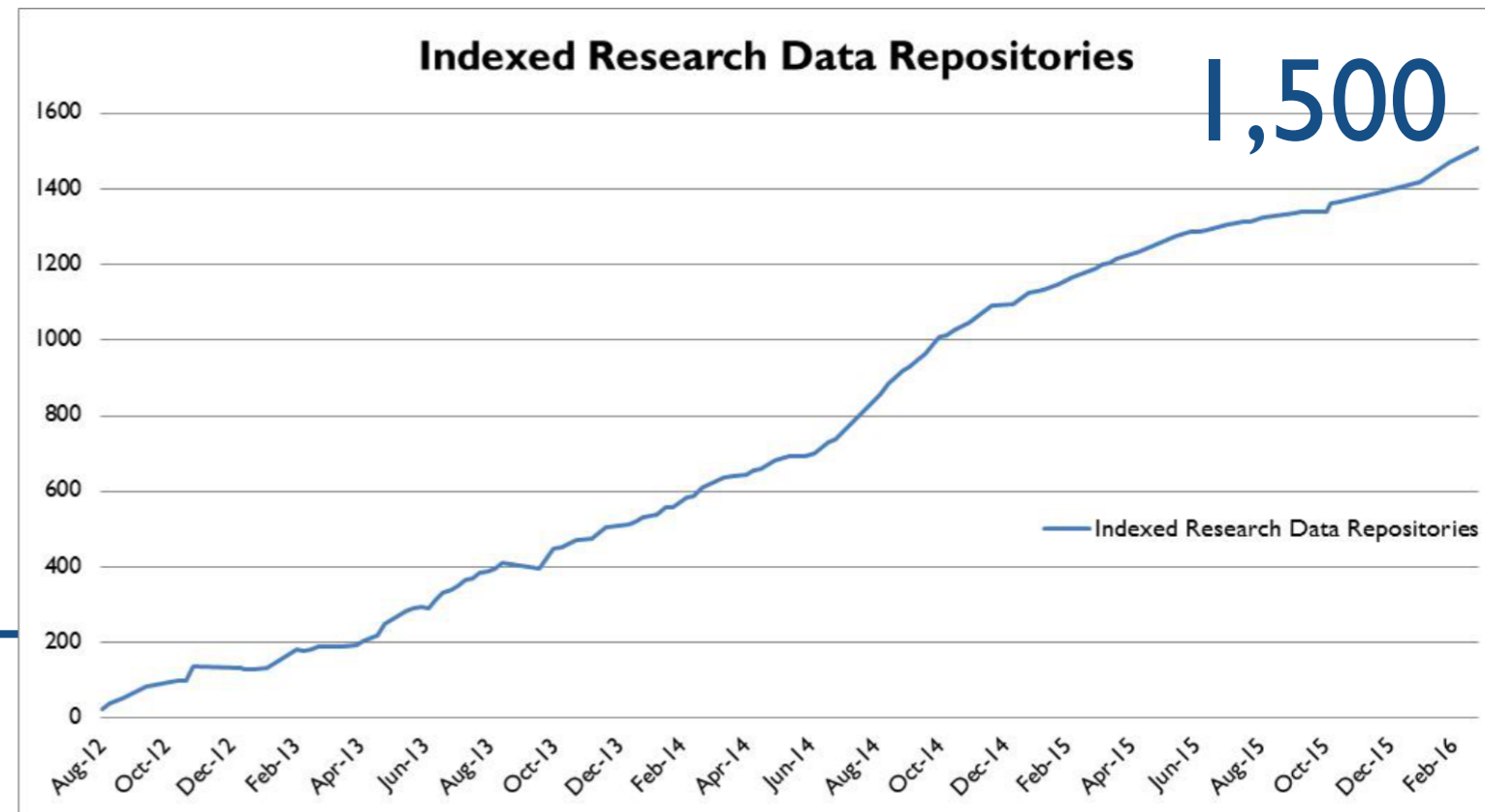
2013

Dublin Core
metadata

DataCite

Data Citation
Principles

Data
Description
Initiative



of (all types of) data repositories from 2012 to 2016

source: r3data.org

*“Research **data publishing** is the release of research data, associated metadata, accompanying documentation, and software code (in cases where the raw data have been processed or manipulated) for re-use and analysis in such a manner that they can be discovered on the Web and referred to in a unique and persistent way.”*

Best Practices for data publishing

- **Data Citation:** to reference and locate with a persistent identifier, and give credit to data authors
- **Metadata:** to discover and reuse
- **Access control rules:** to support publishing workflows and data agreements and licenses, and protect privacy
- **APIs and standards:** to interoperate

The Role of Dataverse

- **Dataverse** is an open-source software platform for building data repositories
- Gives **credit and control** to researchers
- Builds a **community** to:
 - define new standards and best practices
 - foster new research and collaboration in data sharing
- Has brought **data publishing** into the hands of researchers

The Dataverse Software Today

Installed in 20 sites world wide;
Hosts dataverses from > 500 institutions



The Dataverse Community Today

35

GitHub
contributors

Semi-monthly
community calls

300

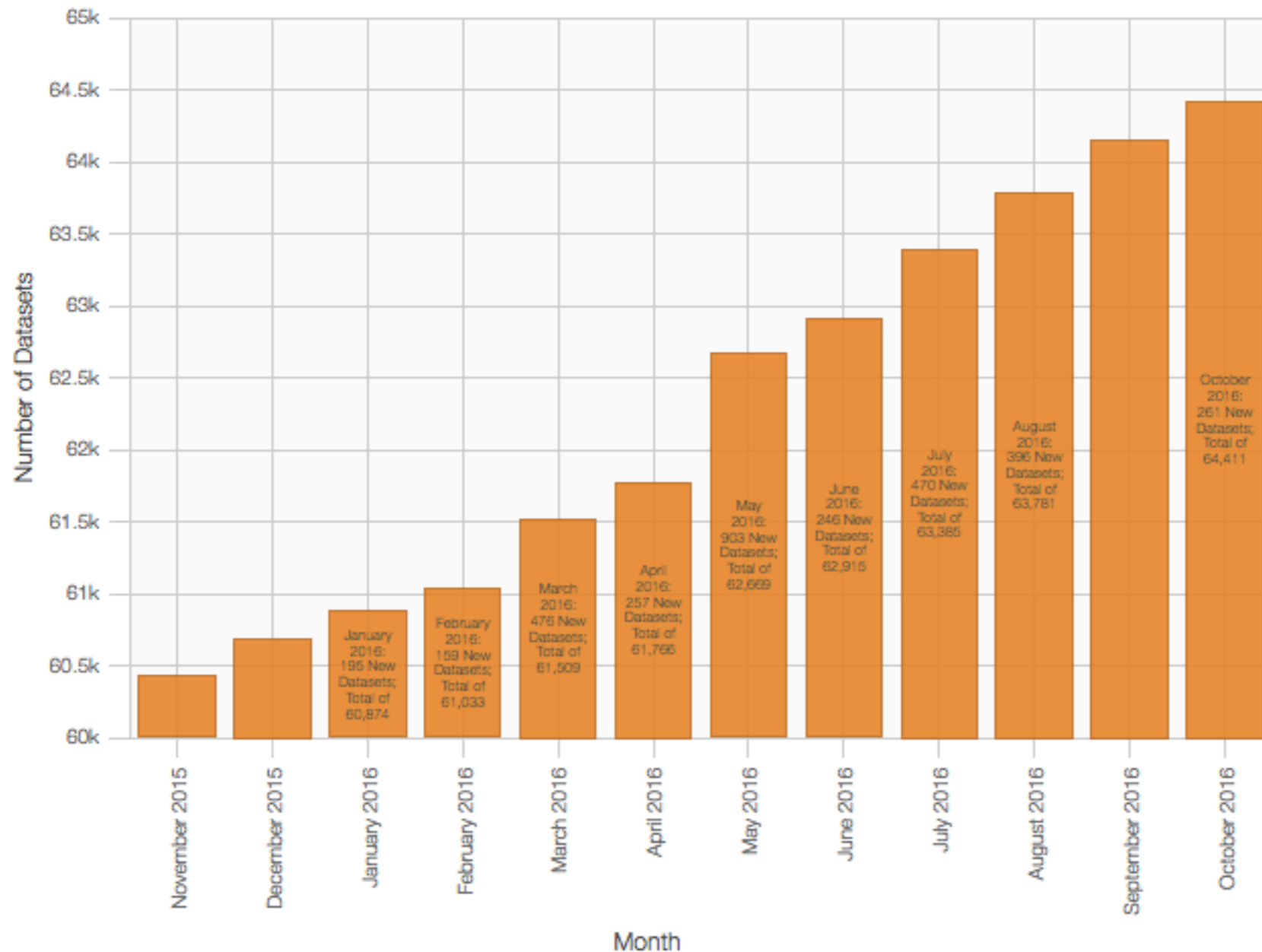
members in the
community list

Annual
Community Meeting,
with near **200**
participants

The Harvard Dataverse Today

Datasets Added Over Time

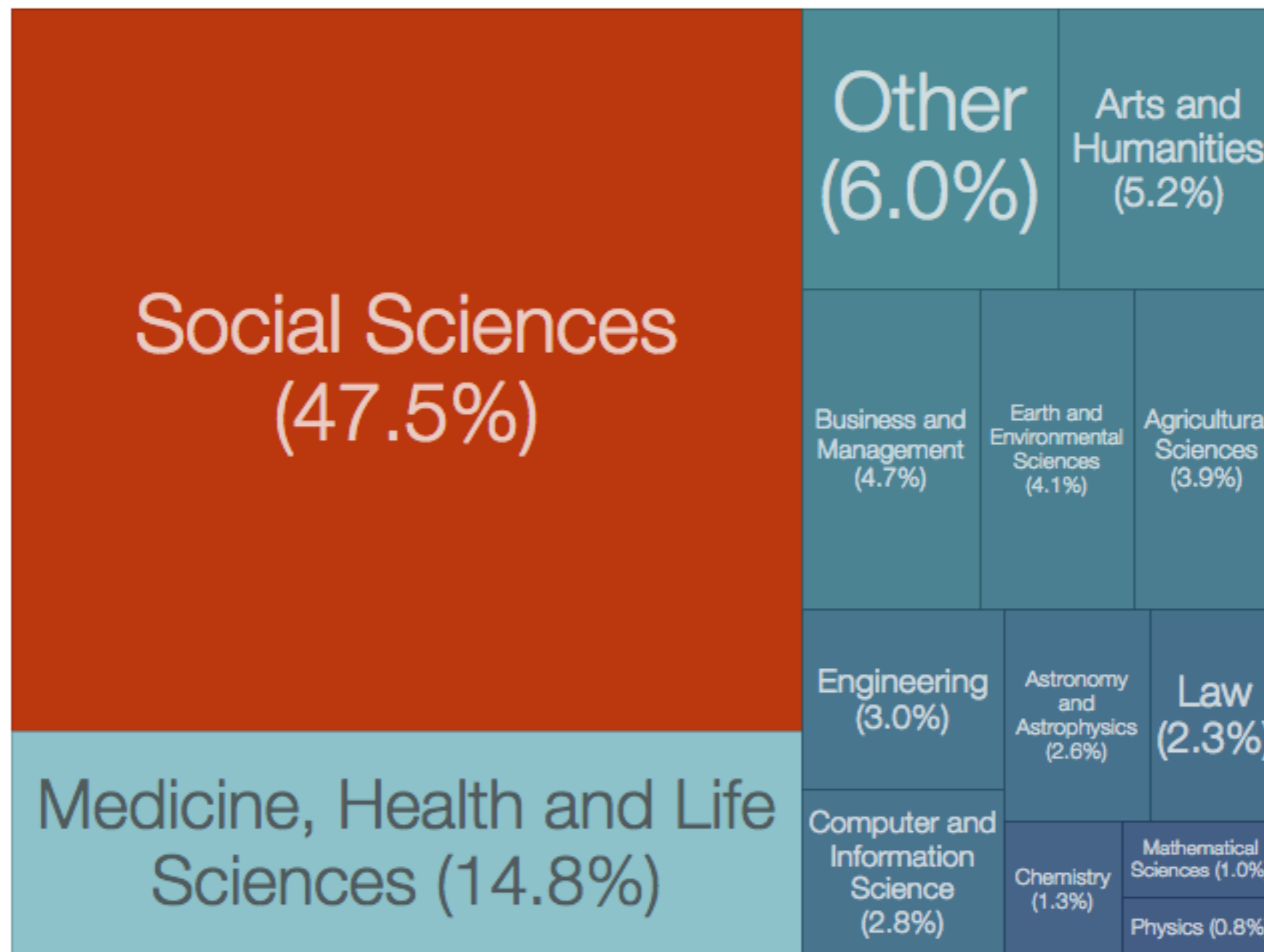
65,000 datasets



~12 new datasets published per day

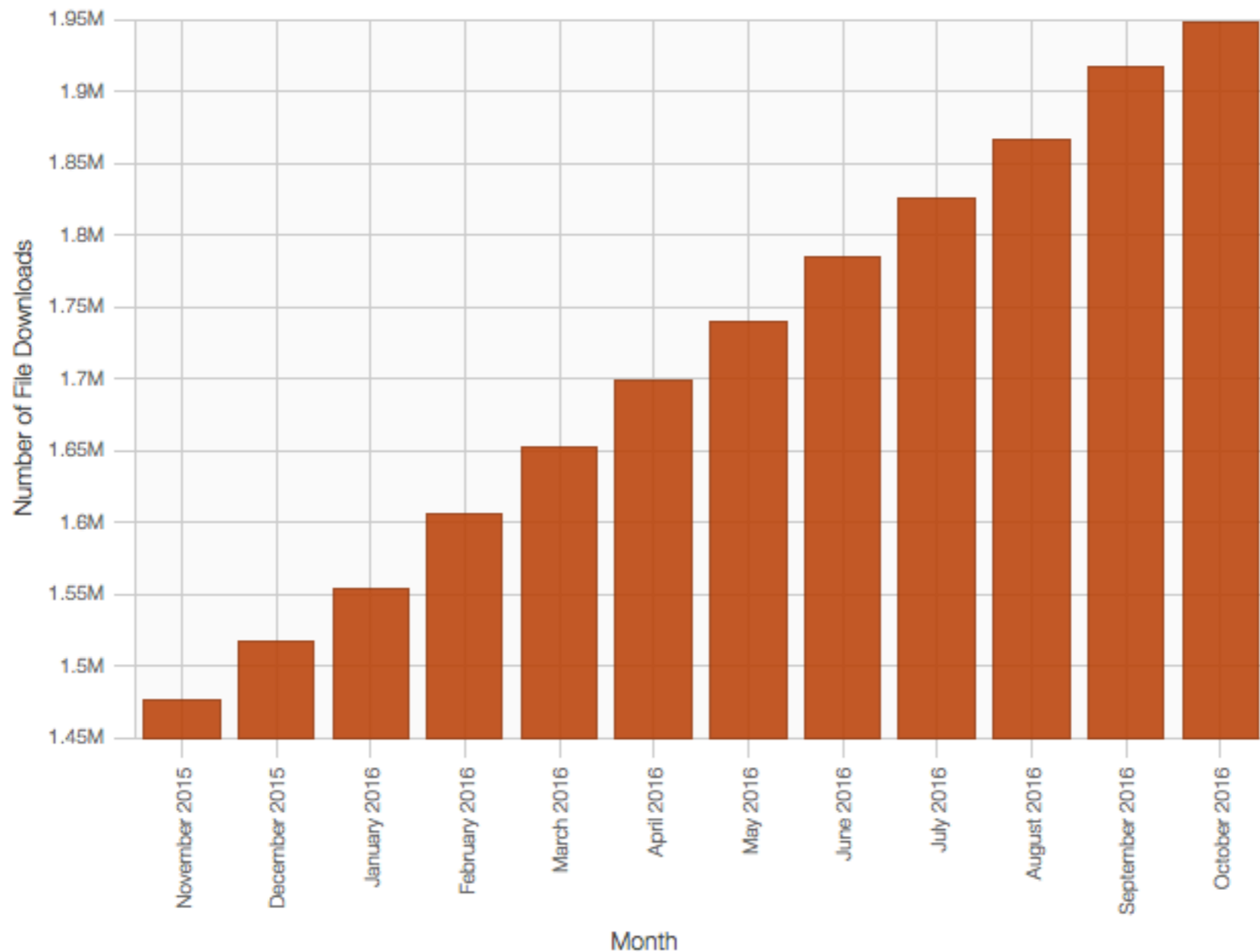
The Harvard Dataverse Today

Datasets by Subject



The Harvard Dataverse Today

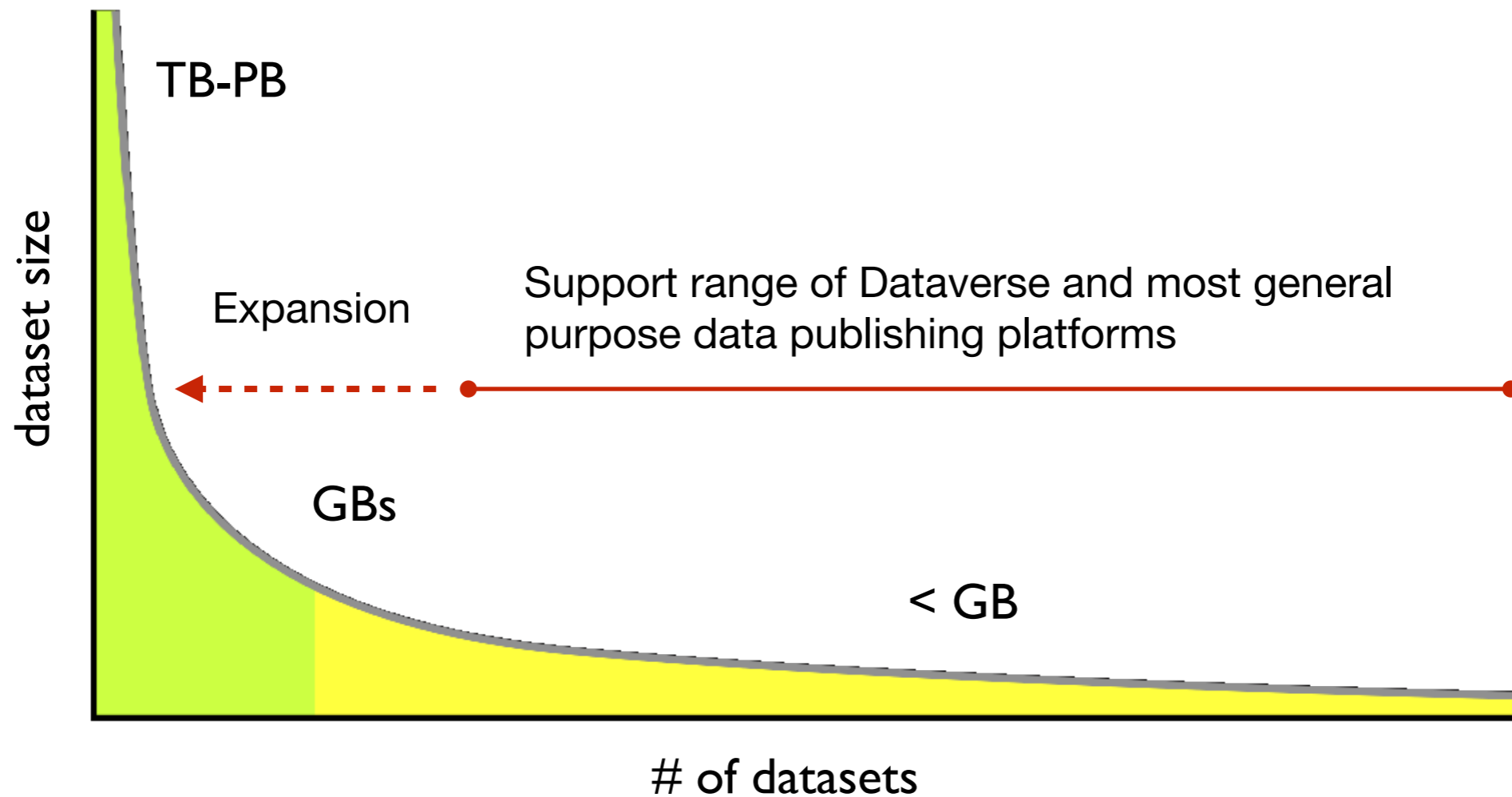
File Downloads Over Time **1.95 Million downloads**



~ 1,500
downloads
per day







Dataverse Today and Tomorrow

custom data
publishing solutions



Distribution of Research Data

Addressing next challenges

- **Provenance:** Seltzer, Crosas, King 
- **Connect journals to data:** King & Crosas  Alfred P. Sloan
FOUNDATION
- **Sensitive data:** Harvard Privacy Tools, led by Vadhan 
- **Large-scale datasets; integration with remote/cloud storage and computing:**
 - Large datasets in biomedicine, Sliz & Crosas  THE LEONA M. AND HARRY B.
HELMSLEY
CHARITABLE TRUST
 - Big Data data in social science, King & Crosas  Alfred P. Sloan
FOUNDATION
 - Cloud Dataverse, with Krieger & Saowarattitada (BU) 

Making Data Accessible use case: Structural Biology Data Grid

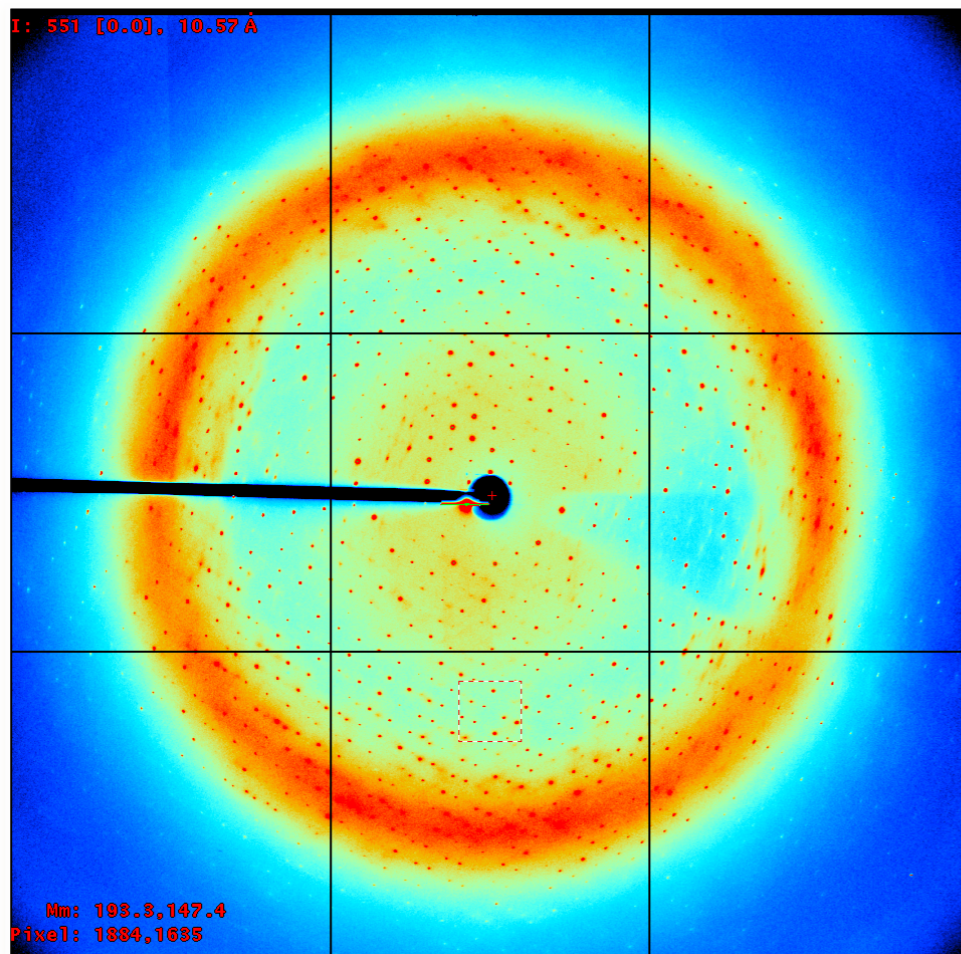


SBGrid
CONSORTIUM



with **Piotr Sliz (Harvard Medical School)**, Pete Meyer, Bill McKinney, Stephanie Socias, Jason Key, Kyle Chard , and IQSS Dataverse team

X-ray diffraction images are the primary data for crystallographic structure determination



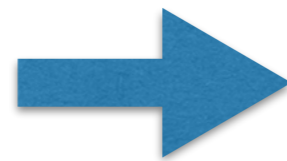
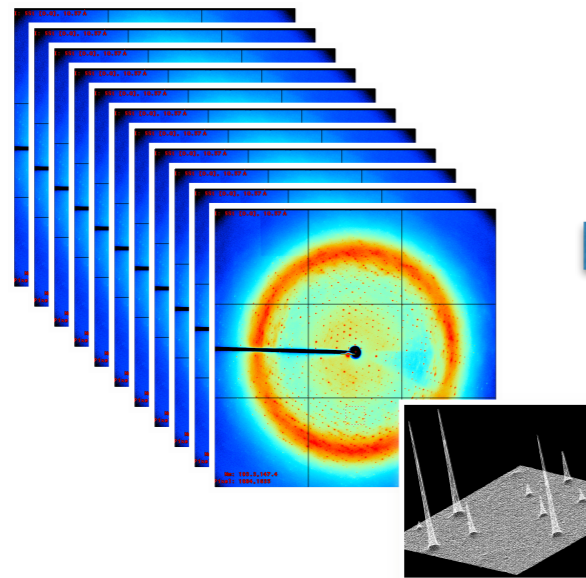
doi:10.15785/SBGRID/19

- X-ray diffraction images collected from frozen protein crystal
- Typically obtained at synchrotron x-ray sources



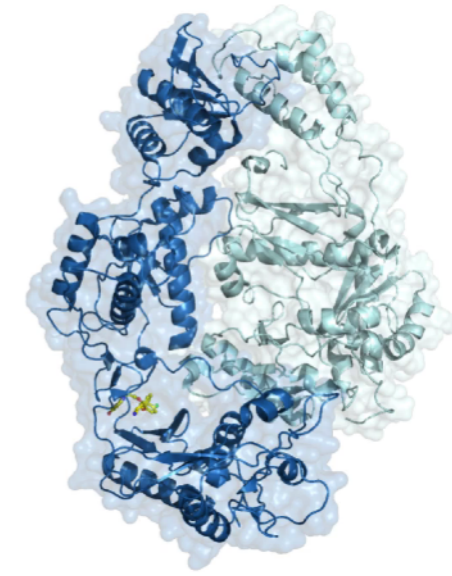
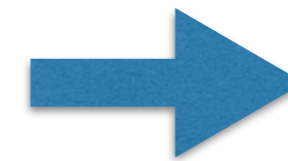
- Primary source for building the protein models published in Protein Data Bank

Why share structural biology primary data?



```

0 0 0 2.257E+02 4.747E+01 4.894E+01 1.082E+02 38.5 0.05310 100 74 -20.27
0 0 -2.647E+03 1.838E+02 894.1 1198.5 38.5 0.05310 100 74 -20.27
0 0 -1.048E+01 1.232E+02 2318.0 1852.1 17.0 0.07354 100 -14 0 0.36
0 0 -2.297E+01 1.222E+01 782.0 1148.5 48.2 0.07375 100 -11 -21.77
0 0 -2.535E+00 1.095E+01 2423.2 1981.5 15.0 0.06413 100 -8 -0.36
0 0 -0.1328E+00 1.532E+01 678.0 1099.7 41.0 0.08442 100 -9 -23.30
0 0 1.168E+02 2.807E+01 2538.0 1938.7 13.0 0.09403 100 46 -2.13
0 0 2.145E+01 2.455E+01 555.0 1849.2 62.0 0.09511 100 25 -24.33
0 0 3.586E+01 2.593E+01 2655.0 1999.3 12.0 0.10547 100 13 -3.65
0 0 -18 -6.824E+00 1.688E+01 2866.0 2095.2 0.0 0.10576 100 -3 -26.39
0 0 -12 2.193E+01 2.958E+01 328.0 951.9 46.0 0.11649 100 9 -27.96
0 0 12 8.822E+00 1.688E+01 2866.0 2095.2 0.0 0.12609 100 -6 -6.41
0 0 -12 -1.698E+01 3.232E+01 288.0 983.0 48.6 0.12728 100 -6 -29.55
0 0 13 -1.154E+01 1.634E+01 981.0 2142.0 6.0 0.12762 100 -12 -6.42
0 0 -13 -3.185E+01 3.328E+01 76.0 856.4 58.3 0.13885 99 -6 -31.17
0 0 -1 1.877E+03 7.188E+01 1835.0 1889.2 35.0 0.02606 99 38 51.38
0 -2 -1 4.818E+02 9.255E-01 1428.4 1558.1 132.0 0.00997 100 6 -3.35
0 2 1 1.773E+03 6.833E+01 1678.2 1451.9 137.0 0.00996 100 29 7.39
0 -2 1 2.187E+03 8.438E+01 1455.1 1481.4 57.0 0.02608 100 32 -52.78
0 -2 2 2.779E+03 1.871E+02 1742.5 1647.0 64.0 0.03818 100 38 38.79
0 -2 -2 -3.174E+00 2.732E+00 1313.1 1508.0 167.7 0.01618 98 -8 -83.53
0 2 2 3.076E+03 1.188E+02 1777.3 1422.0 176.4 0.01682 97 33 181.93
0 2 -2 2.912E+03 1.122E+02 1348.3 1354.9 49.0 0.02619 100 42 -42.69
0 -2 3 6.998E+03 2.691E+02 1858.5 1695.3 38.6 0.04828 100 44 38.75
0 2 -3 6.694E+03 2.577E+02 1283.5 1377.7 12.0 0.02564 100 44 33.45
0 -2 3 6.769E+03 2.611E+02 1887.4 1623.9 0.4 0.02564 84 56 -45.84
0 2 -3 7.134E+03 2.748E+02 1248.2 1386.2 46.4 0.04831 100 53 -37.39
0 -2 4 1.888E+03 7.308E+01 1959.1 1744.5 34.0 0.04056 100 57 29.86
0 -2 -4 1.921E+03 7.454E+01 1895.1 1338.3 28.0 0.03556 99 54 15.82
0 2 4 1.763E+03 6.805E+01 1996.1 1671.4 5.0 0.03537 100 66 -29.62
0 -2 -4 1.882E+03 7.383E+01 1131.7 1256.0 45.2 0.05864 99 71 -34.57
0 2 4 2.285E+03 8.095E+01 2868.0 1794.0 38.7 0.00996 100 71 -28.68
0 -2 -5 2.291E+03 8.099E+01 986.0 1281.9 28.1 0.04573 100 68 4.95
0 2 5 2.238E+03 8.086E+01 2186.3 1718.0 7.8 0.04574 100 75 -22.14
0 -2 6 3.828E+01 1.255E+01 2179.3 1843.7 27.7 0.07144 100 35 17.88
0 -2 -6 4.844E+01 1.803E+01 875.0 1233.0 38.0 0.05683 100 24 -1.67
0 -2 6 5.078E+01 1.877E+01 2217.0 1768.6 7.0 0.05684 100 27 -38.82
0 -2 -6 5.684E+01 1.282E+01 912.2 1157.1 45.2 0.07159 100 27 -32.43
0 2 -6 1.877E+03 4.288E+01 2291.3 1893.3 22.0 0.00138 100 75 -13.88
0 -2 -7 1.854E+03 4.277E+01 764.2 1184.2 33.2 0.06648 100 78 -5.58
0 2 -7 1.823E+03 4.266E+01 2326.0 1817.3 7.0 0.06648 100 87 -17.83
0 -2 -7 1.878E+03 4.454E+01 888.0 1187.2 45.8 0.08217 100 86 -32.38
0 2 -7 8.984E+02 3.288E+01 2444.0 1942.0 22.0 0.06648 100 83 11.28
0 -2 8 6.547E+02 3.873E+01 651.5 1135.1 36.8 0.07691 100 82 -18.52
0 2 8 6.485E+02 3.888E+01 2442.0 1855.0 6.0 0.07691 98 -16.14
0 -2 8 6.817E+02 3.227E+01 688.0 1057.4 46.7 0.09282 100 87 -32.53
0 2 8 9.516E+01 2.363E+01 2519.4 1992.0 28.4 0.18325 99 28 0.66
0 -2 9 8.782E+01 1.943E+01 536.0 1886.6 38.0 0.09742 100 35 -13.87
0 2 9 1.283E+02 2.188E+01 2557.6 1913.9 6.1 0.08725 100 46 -15.88
0 -2 9 1.233E+02 2.160E+01 513.0 1887.8 47.7 0.18354 100 36 -33.83
0 -2 10 1.488E+02 2.793E+01 2635.0 2048.9 18.1 0.11394 95 36 6.29
0 -2 10 9.235E+01 2.586E+01 438.0 1938.1 48.0 0.09897 100 25 -16.83
0 2 10 9.464E+01 2.788E+01 2674.0 1961.7 5.1 0.09774 100 38 -16.84
0 2 -10 7.375E+01 2.771E+01 457.4 958.4 48.0 0.11432 100 14 -33.72
0 -2 11 1.838E+00 3.488E+01 2755.3 2889.4 16.8 0.12471 100 6 4.83
0 -2 -11 -1.886E+01 2.483E+01 82.3 988.3 43.1 0.18807 100 -5 -19.52
0 -2 11 -2.821E+01 4.475E+01 2784.0 2889.8 3.0 0.18826 100 -7 -16.58
0 2 -11 8.443E+01 2.958E+01 339.0 989.3 58.2 0.12512 98 15 -34.57
0 -2 12 3.843E+01 1.235E+01 2877.0 2138.0 13.0 0.13589 100 9 1.86
0 -2 -12 1.492E+00 2.957E+01 181.5 943.5 45.3 0.11917 100 4 -22.84
0 2 -12 -1.864E+01 1.793E+01 2916.0 2055.6 2.0 0.12078 97 -11 -17.19
0 -2 -12 2.645E+01 1.688E+01 218.0 868.7 51.5 0.13681 100 5 -35.54
0 -2 13 8.224E+00 2.138E+01 988.1 2185.5 11.0 0.14656 100 11 -8.25
0 -2 -13 1.565E+01 3.824E+01 57.9 897.0 47.5 0.12978 100 4 -24.43
0 2 13 7.568E+00 1.519E+01 3848.2 2181.5 1.2 0.12958 100 8 -18.86
    
```



1 dataset is 180 - 360 images
3.5 - 7 GB
Primary Data

1 file of 'structure factors':
0.0001 -0.005 GB
Intermediate Data

Atomic Coordinates
0.0001 -0.001 GB
Model

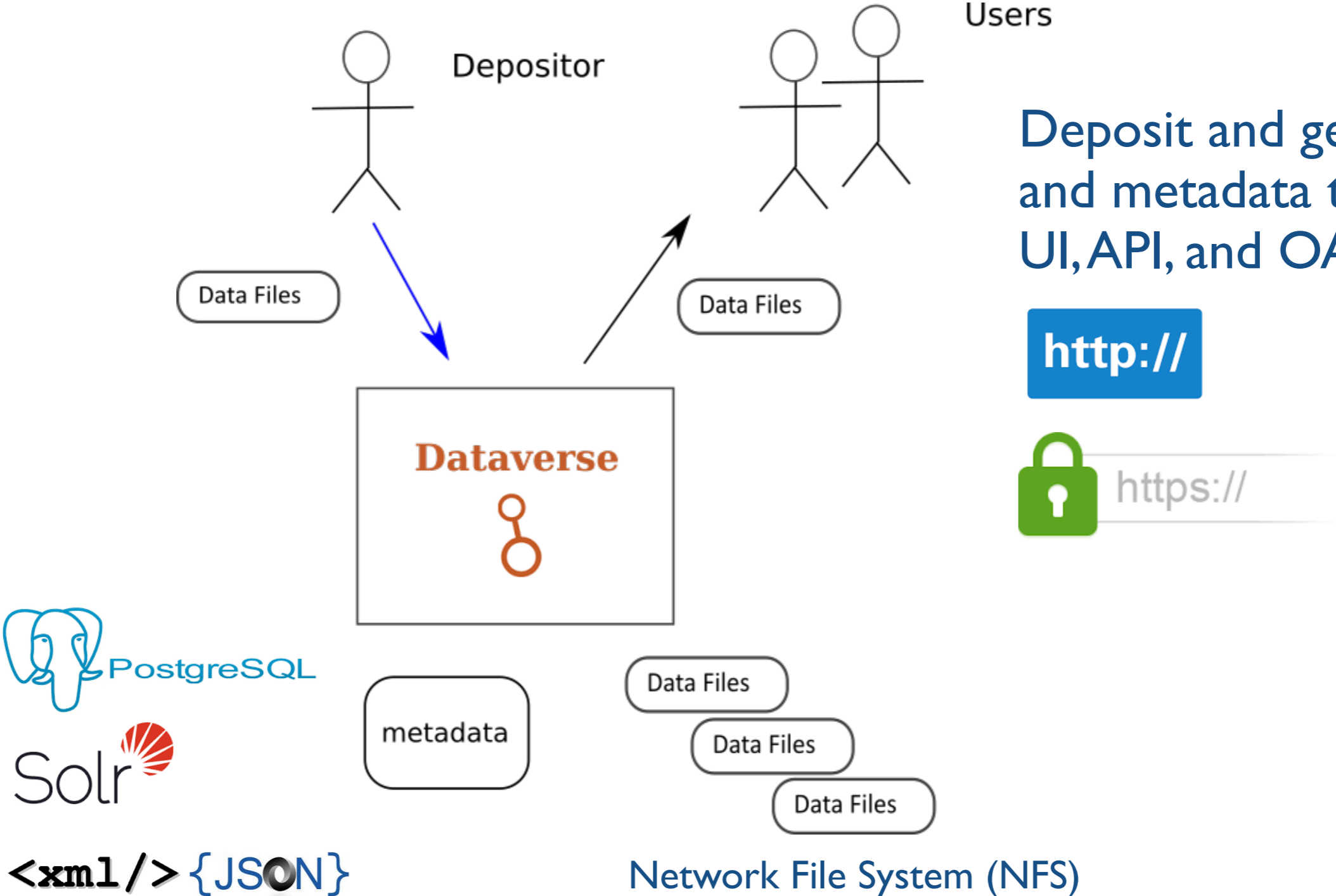


- Replicate published model results
- Reprocess primary data with new, improved analysis

Why Structural Biology Data Grid + Dataverse?


- Support **data publication, data access, preservation, and live analysis** for primary structural biology data
- Apply **best practices** in data sharing: data citation, metadata, standards and protocols
- Use an open-source platform with a growing community as a **sustainable** infrastructure for the repository

Dataverse dataflow now



Deposit and get data and metadata through UI, API, and OAI-PMH

<http://>

 <https://>

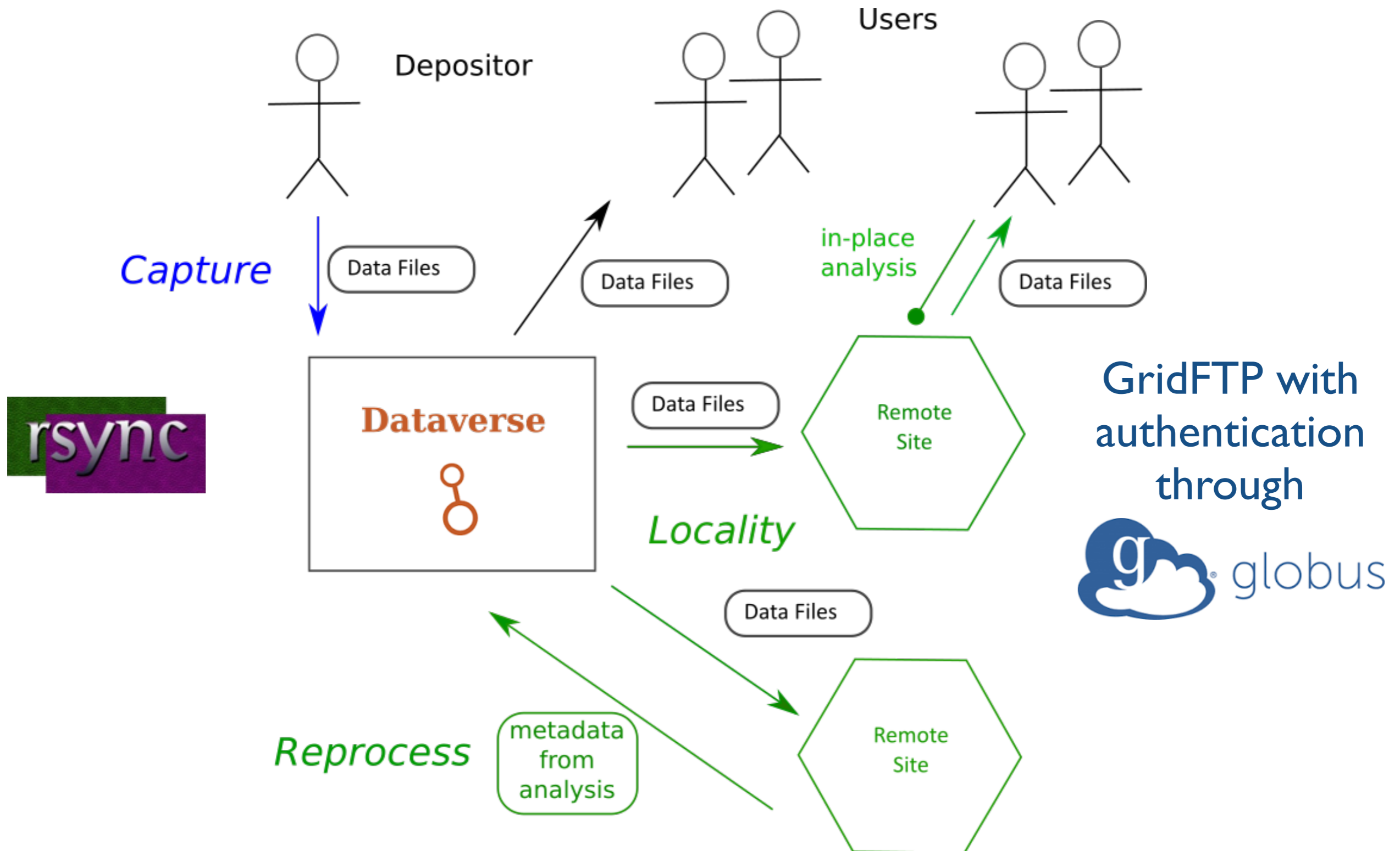
 PostgreSQL

 Solr

<xml /> {JSON}

Network File System (NFS)

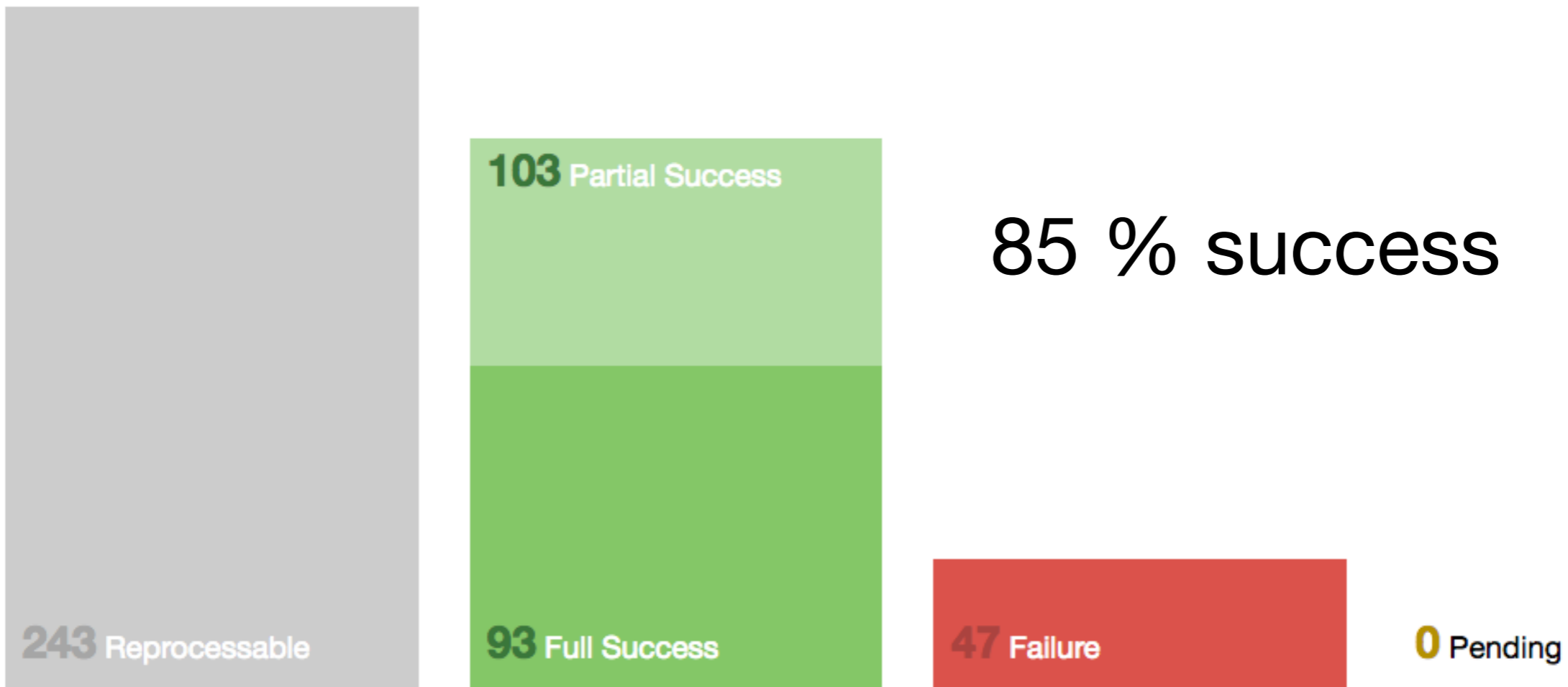
Dataverse expansion



Replication Success Rate

255 Datasets Published

Overall results for datasets with an available reprocessing pipeline shown below. Detailed statistics are reported on dataset landing pages.



Reprocessing Pipeline

X-ray reprocessing summary

Dataset	xia2 -2d	xia2 -3d	xia2 -3dii	xia2 -dials
1	✓	✓	✓	✓
2	✓	✓	✓	✓
3	✗	✓	✓	✓
4	✓	✓	✓	✓
5	✗	✗	✗	✗
6	✗	✗	✗	✗
7	✗	✓	✗	✓
9	✗	⦿	⦿	✗
10	✗	✗	✗	✓
11	✓	✓	✓	✓
12	✓	✓	✓	✓
13	✗	✗	✗	✗
14	✓	✓	✓	✓
15	✗	✓	✓	✓



Sliz Laboratory Dataverse (Harvard Medical School)

Root Dataverse > Sliz Laboratory Dataverse > X-Ray Diffraction data from Lin28/let-7d microRNA complex, source of 3TRZ structure

Metrics

0 Downloads



X-Ray Diffraction data from Lin28/let-7d microRNA complex, source of 3TRZ structure

Sliz, Piotr, 2016, "X-Ray Diffraction data from Lin28/let-7d microRNA complex, source of 3TRZ structure", doi:10.15785/SBGRID/3, Root Dataverse, V1

Cite Data ▾

Learn about Data Citation Standards.

Related Publication

Nam Y, Chen C, Gregory RI, Chou JJ, Sliz P. Molecular Basis for Interaction of let-7 MicroRNAs with Lin28. Cell 2011; 147:1080–1091. doi: 10.1016/j.cell.2011.10.020

3TRZ Coordinates

PDB, MMDB

Biological Sample

Lin28/let-7d microRNA complex

Dataset Type

X-Ray Diffraction

Subject Composition

RNA

Data Creation Date

2010-09-07

Reprocessing Data

```
>_ xia2 -dials /programs/datagrid/3 verbose=True multiprocessing.nproc=1 -ispyb_xml_out apc.xml
>_ xia2 -3dii /programs/datagrid/3 verbose=True multiprocessing.nproc=1 -ispyb_xml_out apc.xml
>_ xia2 -2d /programs/datagrid/3 verbose=True multiprocessing.nproc=1 -ispyb_xml_out apc.xml
Version: 0.4.0.0
Reprocessing failed.
>_ xia2 -3d /programs/datagrid/3 verbose=True multiprocessing.nproc=1 -ispyb_xml_out apc.xml
```

Data Access Alliance:

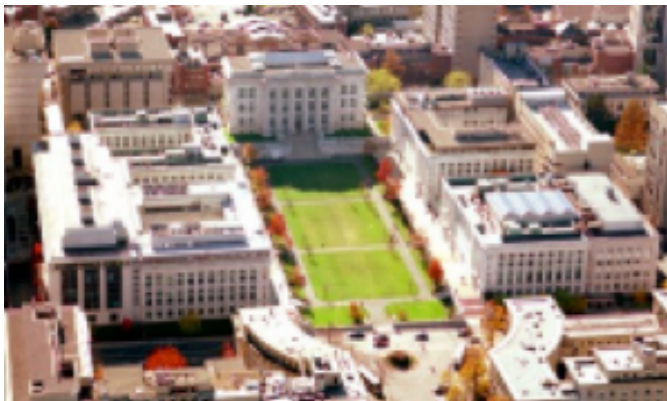
Local Access through a growing list of satellites



Petrel Storage
Argonne Labs



UPPSALA
UNIVERSITET



Harvard Medical
School Orchestra
Cluster



San Diego
Supercomputer
Center



Institut Pasteur de
Montevideo

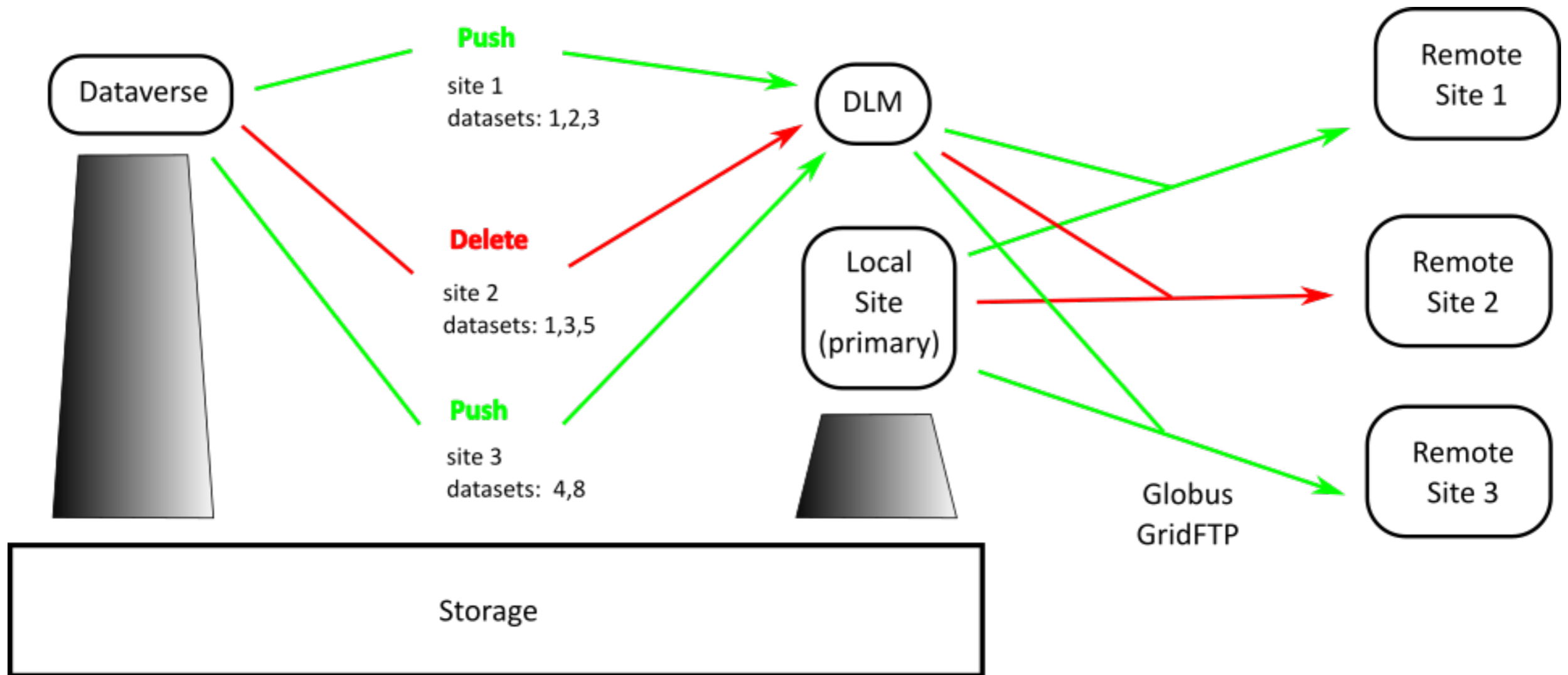


SHANGHAI INSTITUTES FOR BIOLOGICAL SCIENCES

Data Access Alliance Flow

OAI set per remote site

Data transfer through GridFTP  globus

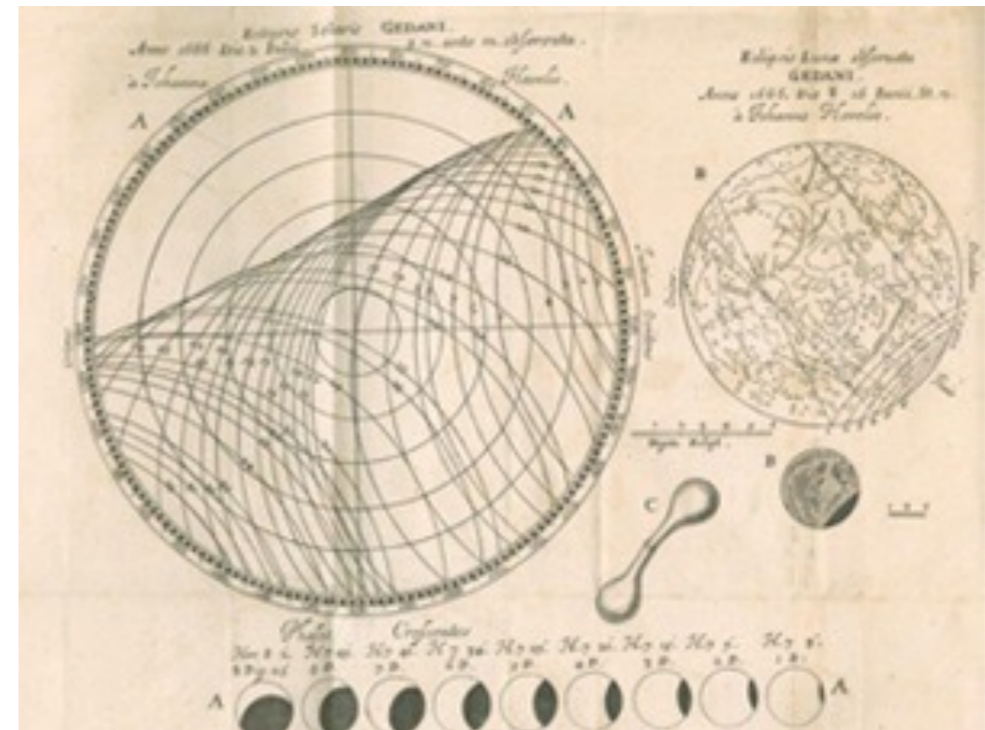


Summary

- In the last decade, *data publishing* emerges as a field in science communication
- Dataverse plays a key role leading data publishing and promoting best practices for making data accessible
- Dataverse is expanding to support data provenance, sensitive data and large scale data, while connecting data with publications and analysis.
- The Structural Biology Data Grid project shows how to:
 - make data accessible in a sustainable way
 - follow best practices in data publishing
 - support larger scale data
 - integrate data with analysis

“The 21st of April, 1665, about eight in the morning, I bored a hole in the body of a fair and large Birch, and put in a Cork with a Quill in the middle; after a Moment or two it [a sap] began to drop, but yet very softly: Some three Hours after I returned, and it had filled a Pint Glass, and then it dropped exceeding fast, viz. every Pulse a Drop: This Liquor is not unpleasant to the Taste, and not thick or troubled; yet it looks as though some few drops of Milk were split in a Bason of Fountain Water.”

(Lister, 1697)



Schematic of solar and lunar eclipses in a 1665 paper by Hevelius

(Philosophical Transactions, Vol 1)