

UCLA

Presentations

Title

Big Data, Little Data, or No Data? A Social Science Perspective on Data Science [Presentation slides]

Permalink

<https://escholarship.org/uc/item/2911049g>

Author

Borgman, Christine L.

Publication Date

2021-03-19

Supplemental Material

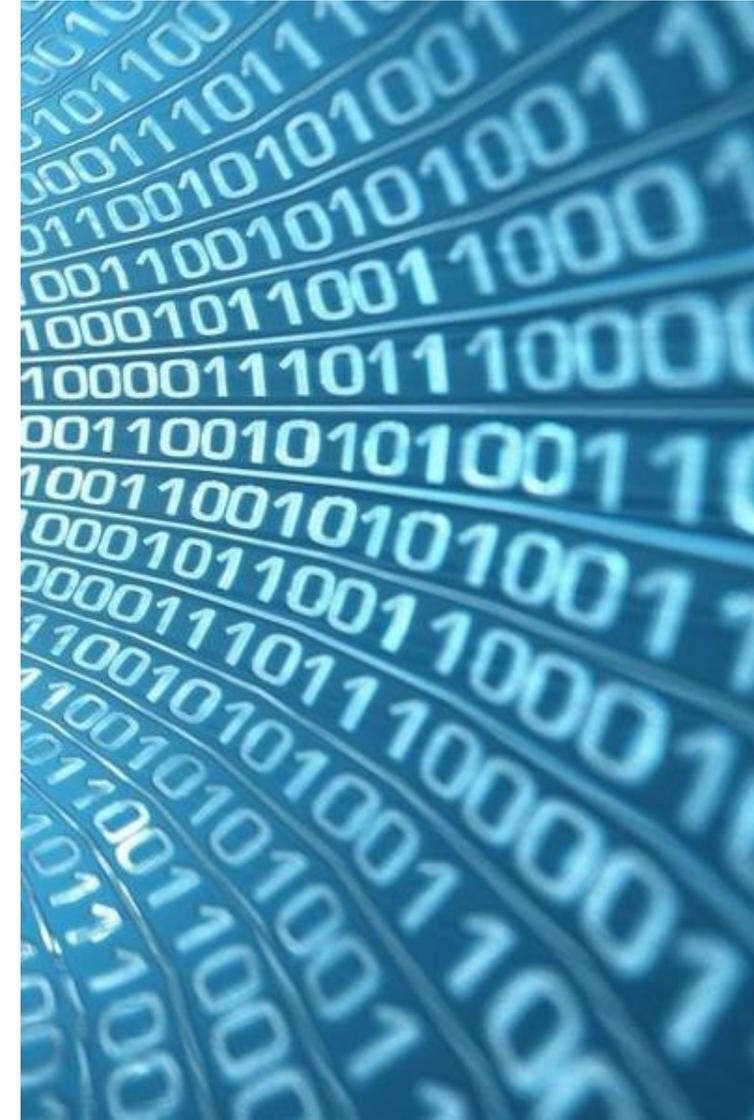
<https://escholarship.org/uc/item/2911049g#supplemental>

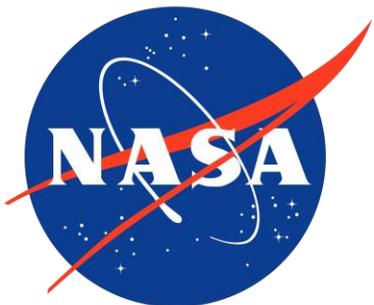
Copyright Information

This work is made available under the terms of a Creative Commons Attribution-NonCommercial License, available at <https://creativecommons.org/licenses/by-nc/4.0/>

Data Challenges in (Data) Science

- How to make data useful and reusable?
- How to decide what data are worth keeping?
- How to balance incentives and benefits?
- How to steward data resources?
- Who pays for infrastructure?

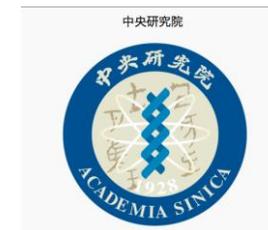




Data sharing policies



- U.S. Federal research policy
- European Research Council
- Research Councils of the UK
- Australian Research Council
- Individual countries, funding agencies, journals, universities



Supported by
wellcometrust



Australian Government
National Health and Medical Research Council



National Science Foundation
WHERE DISCOVERIES BEGIN



National Institutes of Health
Turning Discovery Into Health



U.S. DEPARTMENT OF
ENERGY

Policy RECommendations for Open Access to Research Data in Europe





Open Data Practices



- Link datasets to journal article or publication
- Deposit datasets in a data archive
- Publish data documentation
 - Research protocols
 - Codebooks
 - Software
 - Algorithms
- Cite data and software



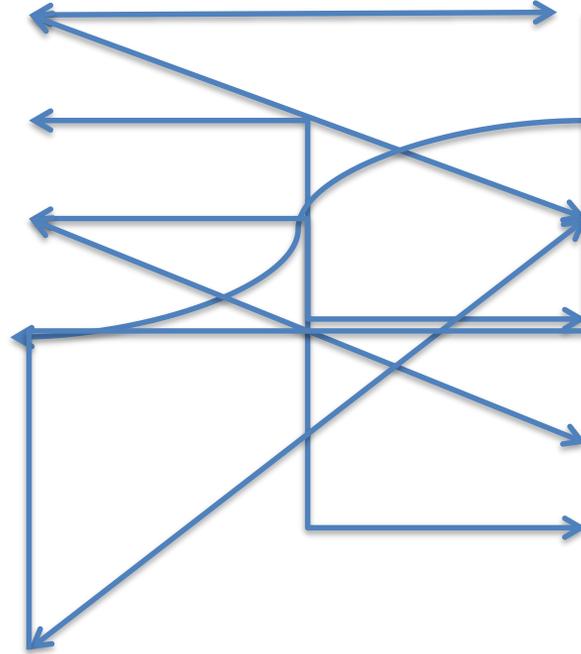
PDS: The Planetary Data System



Publications \leftrightarrow Data: Mapping

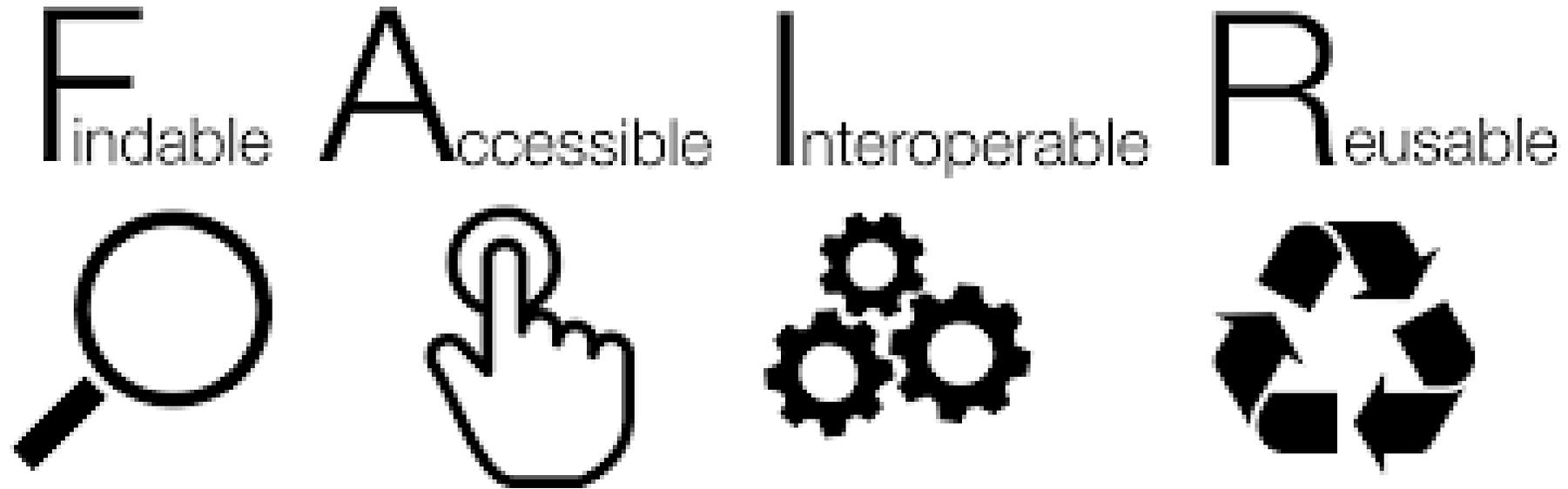
- Article 1
- Article 2
- Article 3
- Article 4

- Article n



- Dataset time 1
- Dataset time 2
- Observation time 1
- Visualization time 3
- Community collection 1
- Repository 1

Data Stewardship: The Ideal



Wilkinson, et al. (2016). The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data*, 3, <http://dx.doi.org/10.1038/sdata.2016.18>

[Topics](#)[Missions](#)[Galleries](#)[NASA TV](#)[Follow NASA](#)[Downloads](#)[About](#)[NASA Audiences](#)

404 The cosmic object you are looking for has disappeared beyond the event horizon.



National Aeronautics and Space Administration
NASA Official: Brian Dunbar

[No Fear Act](#)[FOIA](#)[Privacy](#)[Office of Inspector General](#)[Office of Special Counsel](#)[Agency Financial Reports](#)[Contact NASA](#)

Data

Cassini-Huygens: Mission to Saturn BY THE NUMBERS

2.5 MILLION
COMMANDS
executed

4.9 BILLION
MILES TRAVELED
since launch
(7.9 BILLION KILOMETERS)

635 
SCIENCE DATA
collected

3,948
SCIENCE PAPERS
published

6 NAMED MOONS
discovered

294 ORBITS
completed

162 TARGETED
FLYBYS
of Saturn's moons

453,048
images taken

27 NATIONS
participated

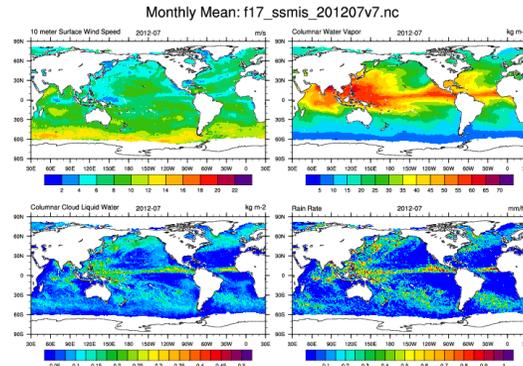
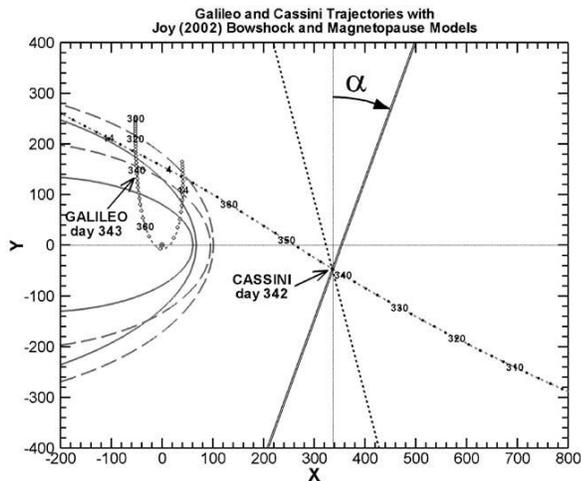
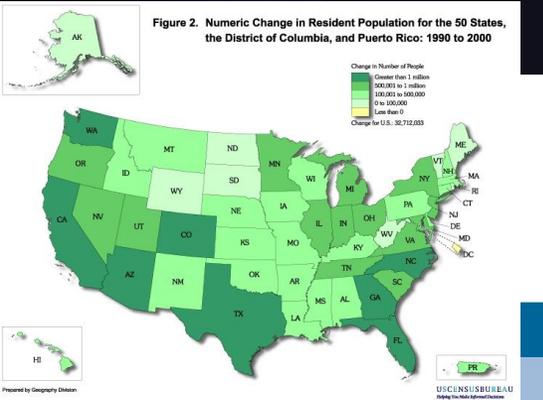
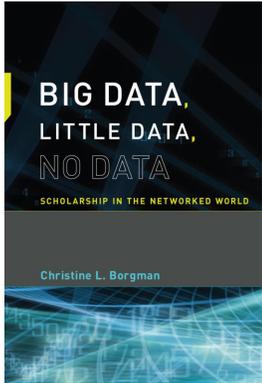
360 ENGINE
burns



NASA Jet Propulsion Laboratory
California Institute of Technology

@CassiniSaturn
saturn.jpl.nasa.gov

Data are representations of observations, objects, or other entities used as evidence of phenomena for the purposes of research or scholarship.

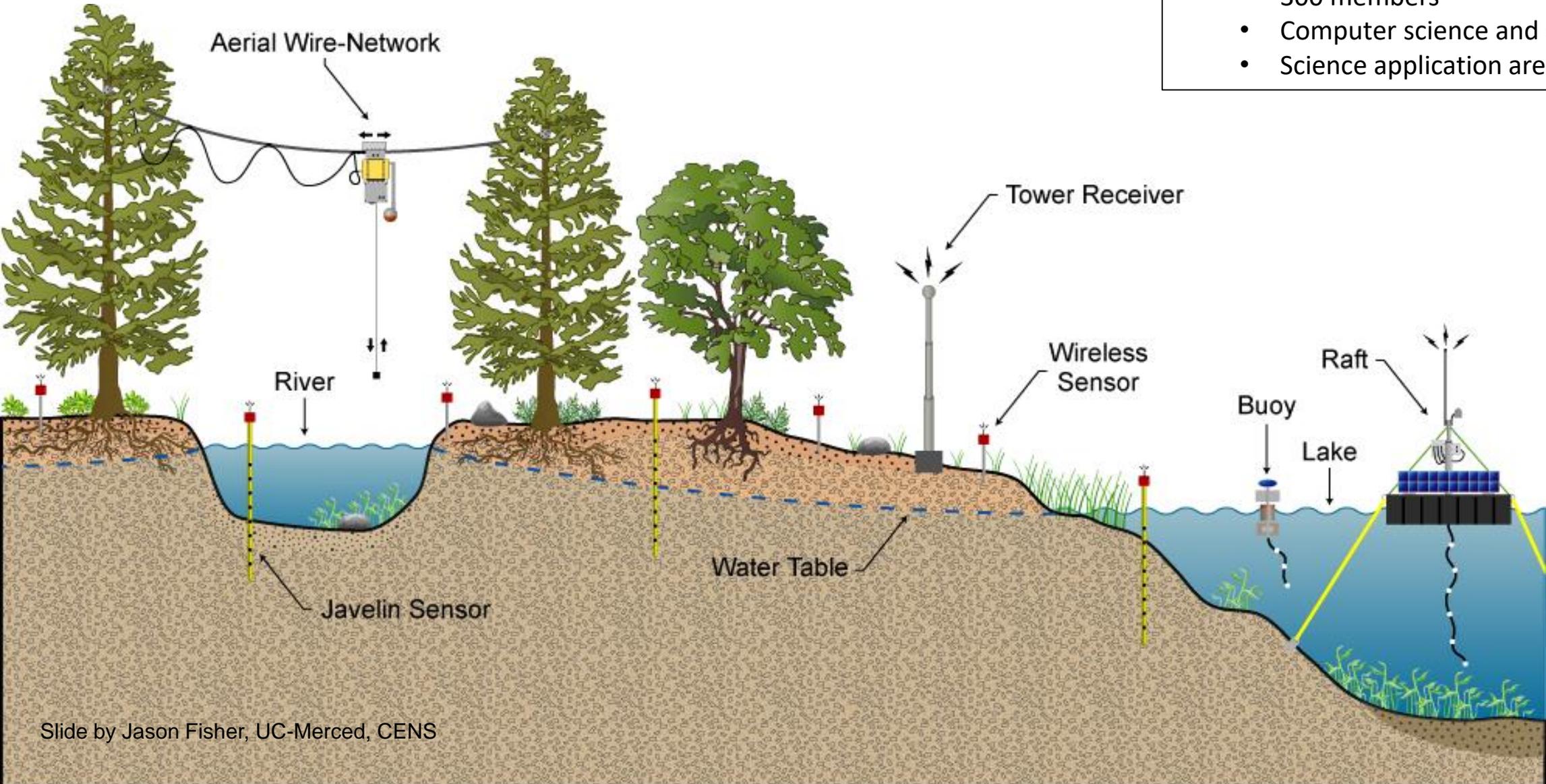


Kivelson, M. G., & Southwood, D. J. (2003). First evidence of IMF control of Jovian magnetospheric boundary locations: Cassini and Galileo magnetic field measurements compared. *Planetary and Space Science*, 51(13), 891–898. [https://doi.org/10.1016/S0032-0633\(03\)00075-8](https://doi.org/10.1016/S0032-0633(03)00075-8)



Center for Embedded Networked Sensing

- NSF Science & Tech Ctr, 2002-2012
- Deborah Estrin, PI
 - 5 universities, plus partners
 - 300 members
 - Computer science and engineering
 - Science application areas



Slide by Jason Fisher, UC-Merced, CENS

Science \leftrightarrow Data

Engineering researcher:

“Temperature is temperature.”



CENS Robotics team

Science \leftrightarrow Data

Engineering researcher:
“Temperature is temperature.”



CENS Robotics team

Biologist: ***“There are hundreds of ways to measure temperature.***
‘The temperature is 98’ is low-value compared to, ‘the temperature of the surface, measured by the infrared thermopile, model number XYZ, is 98.’ That means it is measuring a proxy for a temperature, rather than being in contact with a probe, and it is measuring from a distance. The accuracy is plus or minus .05 of a degree. I [also] want to know that it was taken outside versus inside a controlled environment, how long it had been in place, and the last time it was calibrated, which might tell me whether it has drifted..”

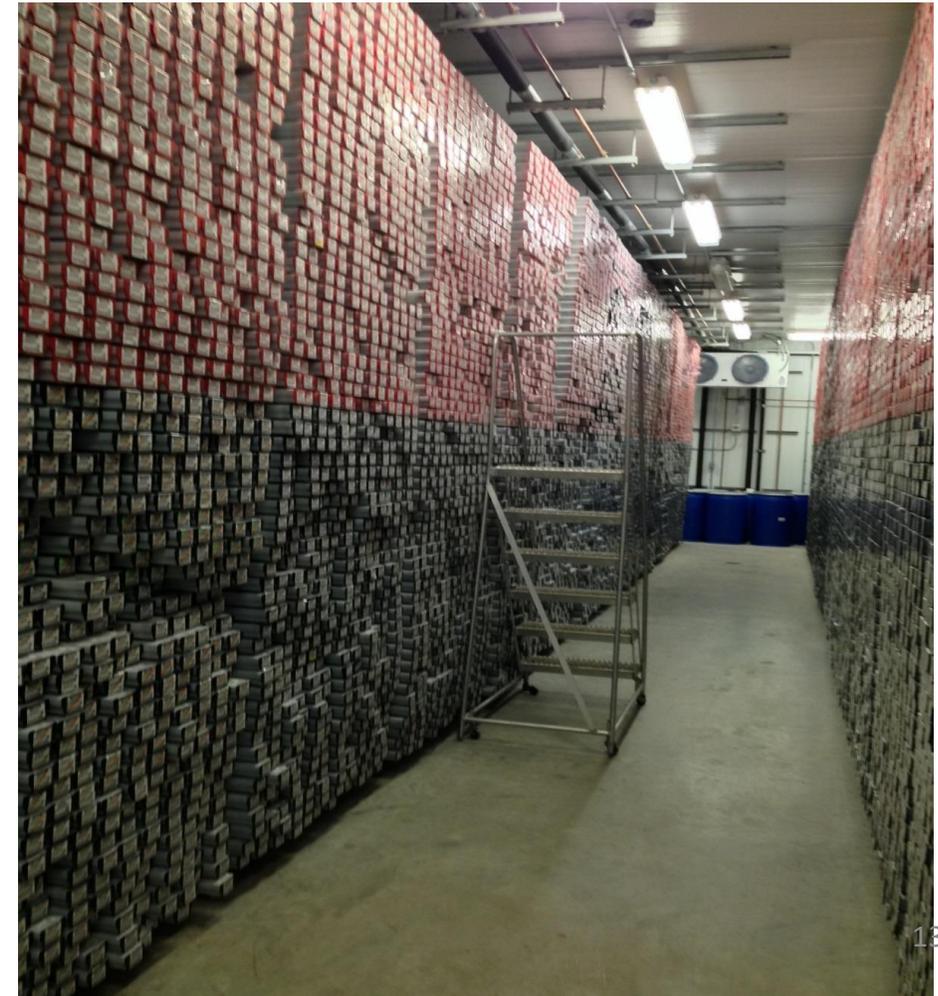
Center for Dark Energy Biosphere Investigations



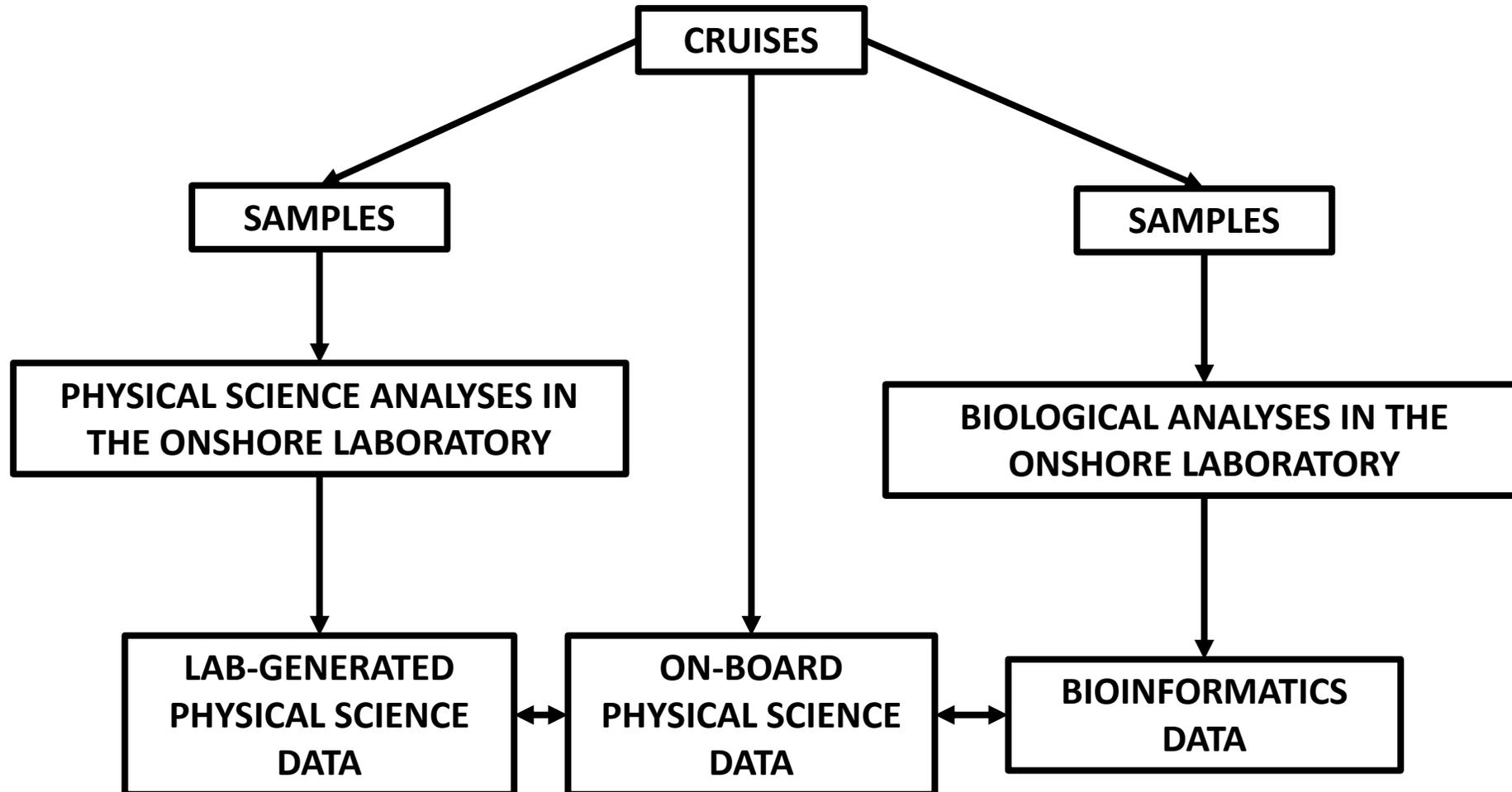
International Ocean Discovery Program
lodp.tamu.org

- NSF Science & Tech Ctr, 2010-2020
- 20 universities, plus partners (35 institutions)
- 90 scientists
- Physical sciences
- Biological sciences

Repository for seafloor cores. Photo: Peter Darch

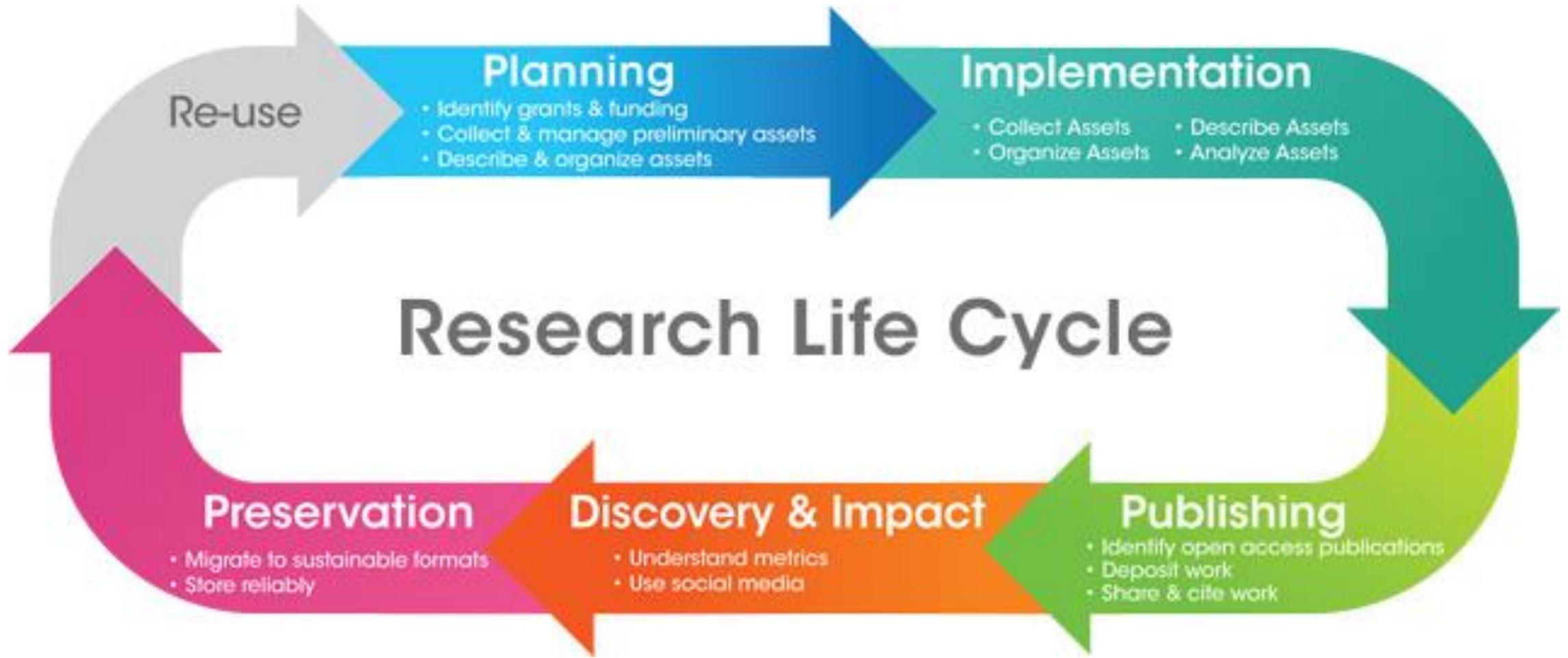


Data Diverge During Scientific Work



Data Practices

Data creation and reuse: The Ideal



Borgman, C. L. (2019). The lives and after lives of data. *Harvard Data Science Review*, 1(1).
<https://doi.org/10.1162/99608f92.9a36bdb6>

Lack of incentives to share data

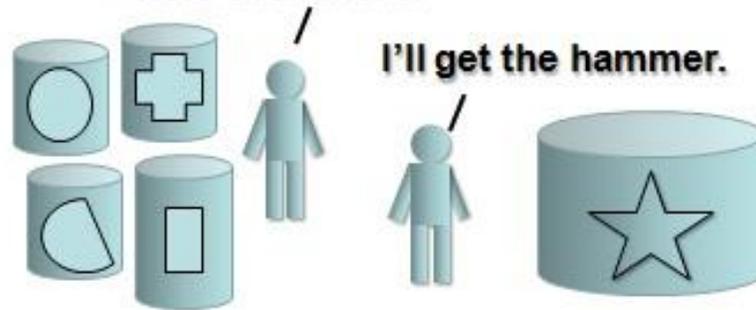
- Labor to document data
- Benefits to unknown others
- Competition
- Control
- Confidentiality
- Lack of expertise and staff
- Lack of sustainability...



Data Stewardship: The Reality



We just need to migrate the data from these systems to fit into that hole over there.



<http://www.datamartist.com/data-migration-part-1-introduction-to-the-data-migration-delema>



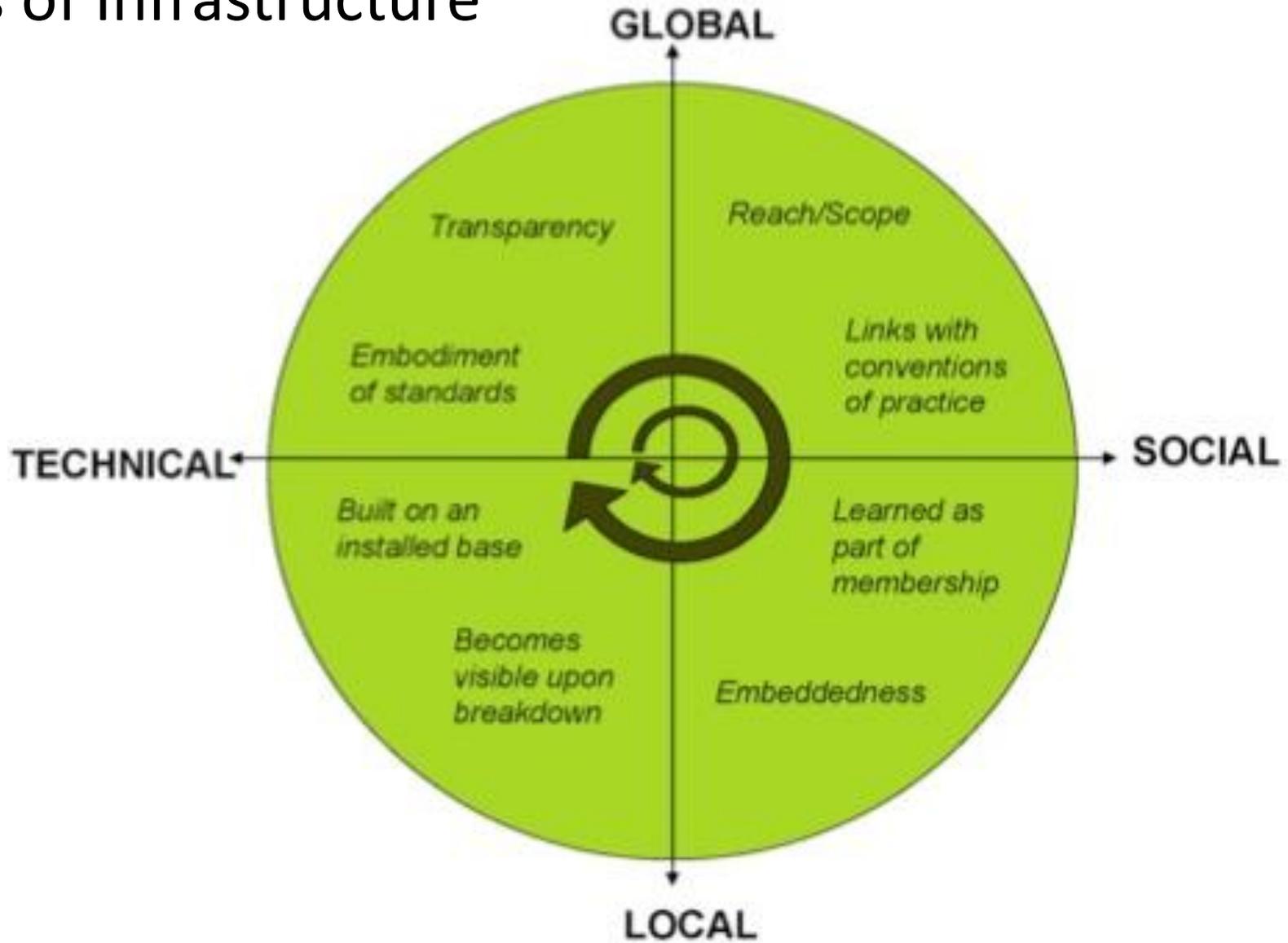
Graduate students



Post-doctoral fellows ¹⁸

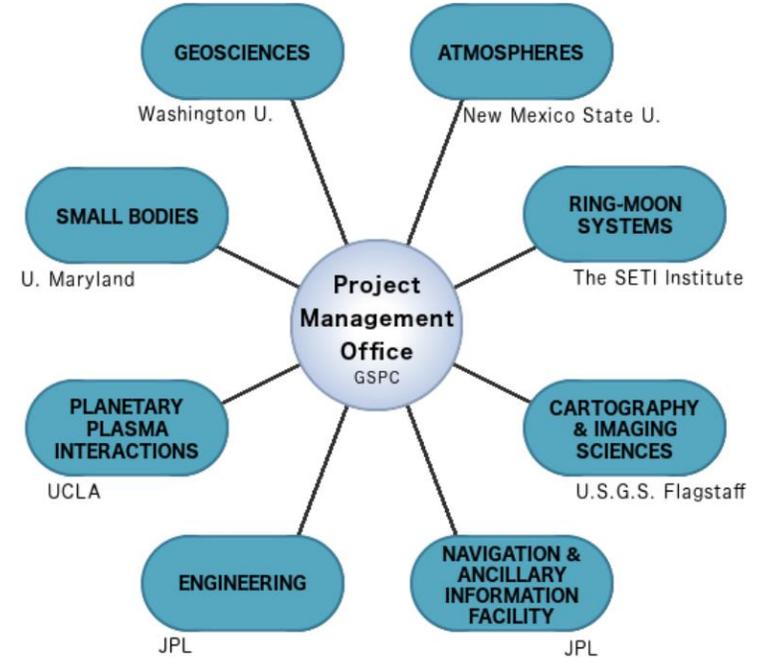
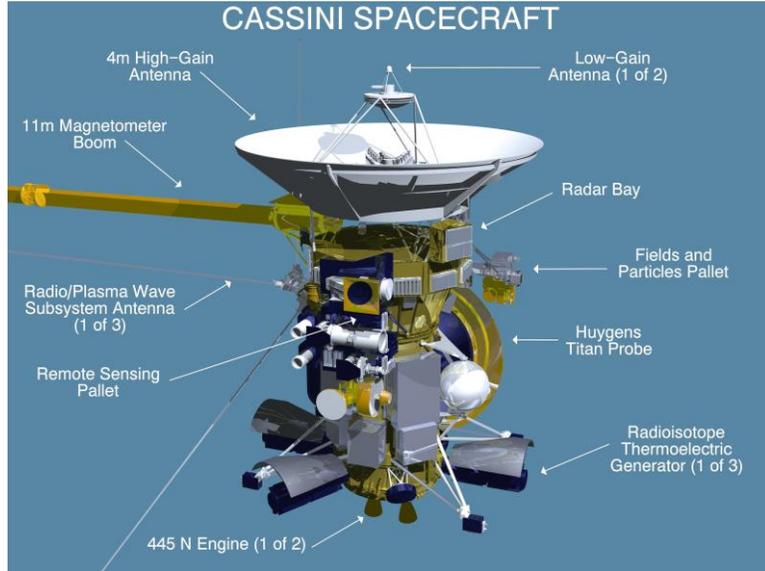
Infrastructure

Dimensions of Infrastructure



Star, S. L. & Ruhleder, K. (1996). Steps toward an ecology of infrastructure: Design and access for large information spaces. *Information Systems Research*, 7(1): 111-134. Figure by Florence Millerand, from: Edwards, P. N., Jackson, S. J., Bowker, G. C. & Knobel, C. P. (2007). *Understanding Infrastructure: Dynamics, Tensions, and Design*. National Science Foundation: University of Michigan. NSF Grant 0630263. <http://hdl.handle.net/2027.42/493530>

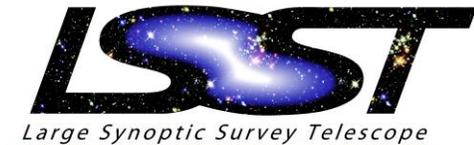
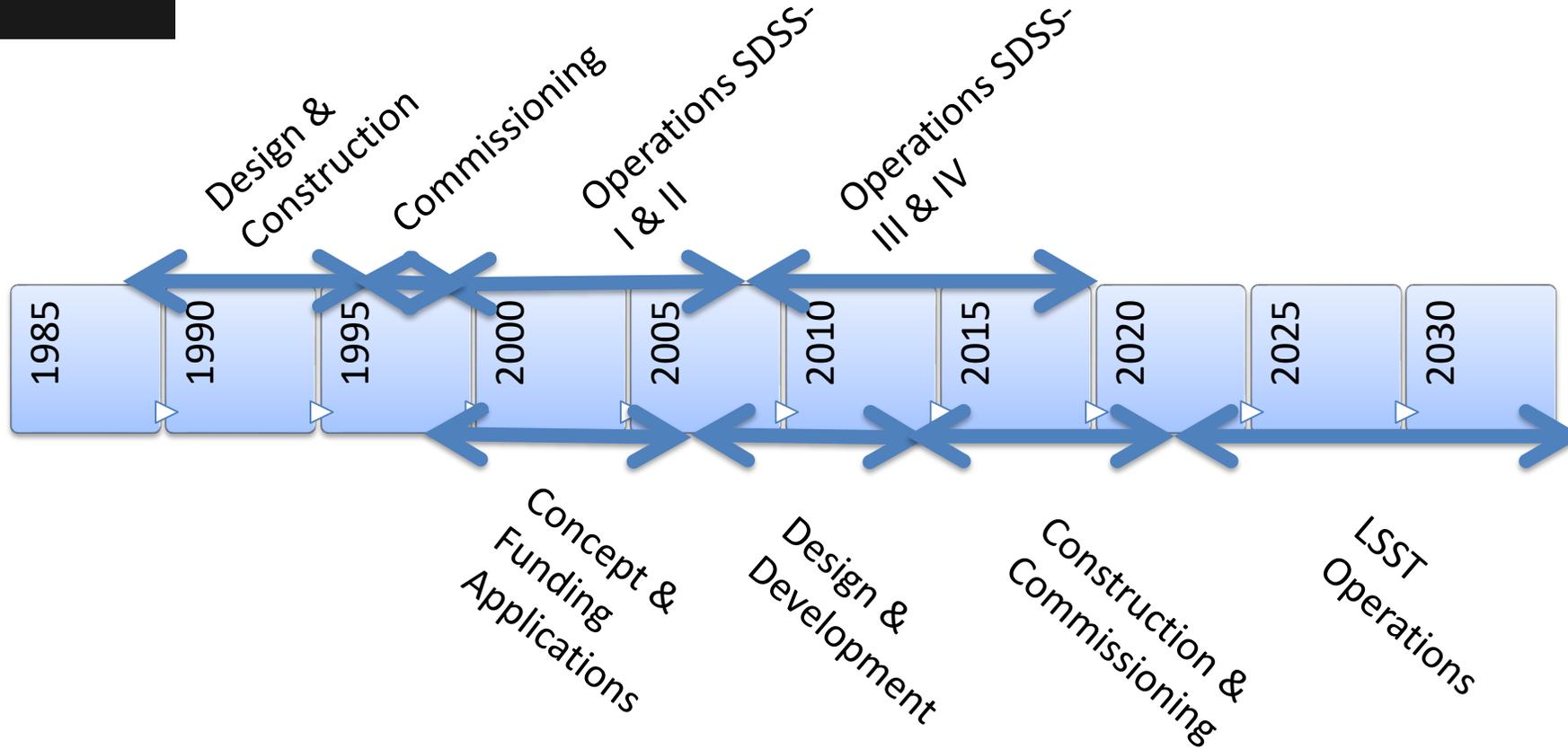
Global and Technical



HOME	DATA SEARCH	TOOLS	DATA STANDARDS	
Home	About PDS	Data Users	Data Proposers	Data Providers



Project Timelines



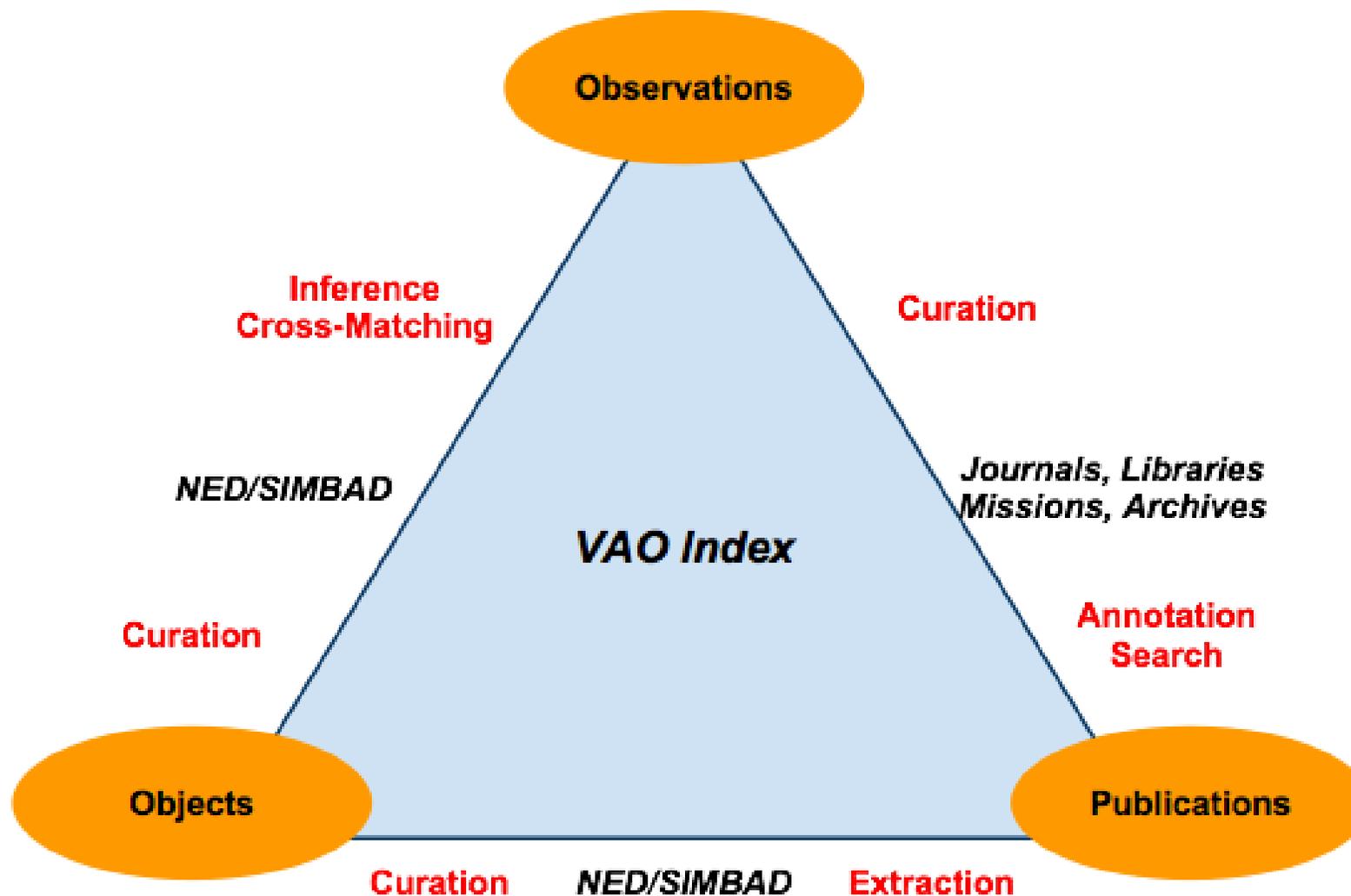
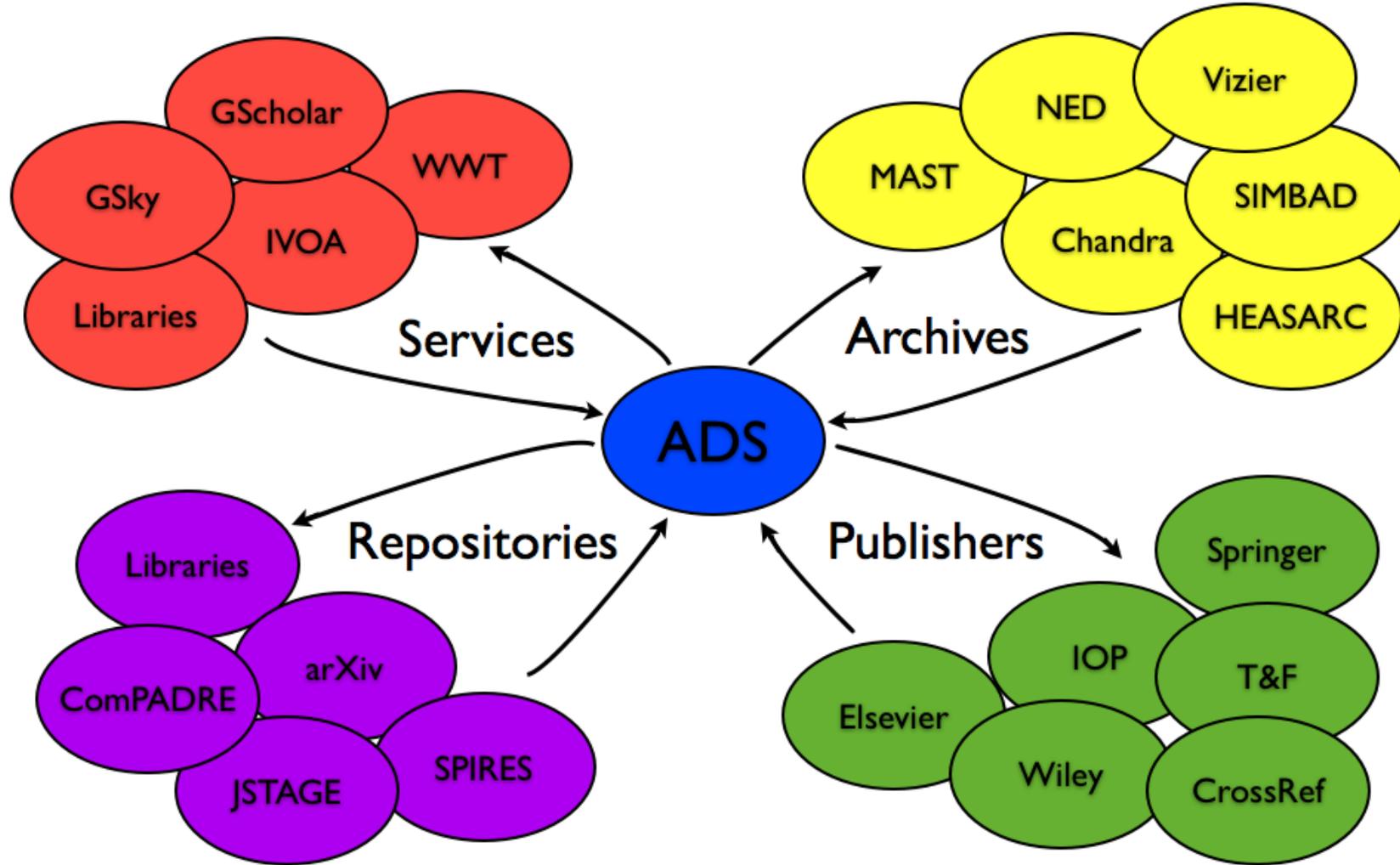


Figure 1. Relationships between Publications, Objects, Observations and the corresponding major actors in the curating process and their activities (in red).

ADS Collaborators



Local and Social

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21st century, requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants

PROGRAMMING & DATABASE

- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing packages, e.g., R
- ☆ Databases: SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and parallel query processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience with xaaS like AWS

DOMAIN KNOWLEDGE & SOFT SKILLS

- ☆ Passionate about the business
- ☆ Curious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative

COMMUNICATION & VISUALIZATION

- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data-driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
- ☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau



<https://github.com/okulbilisim/awesome-datascience>

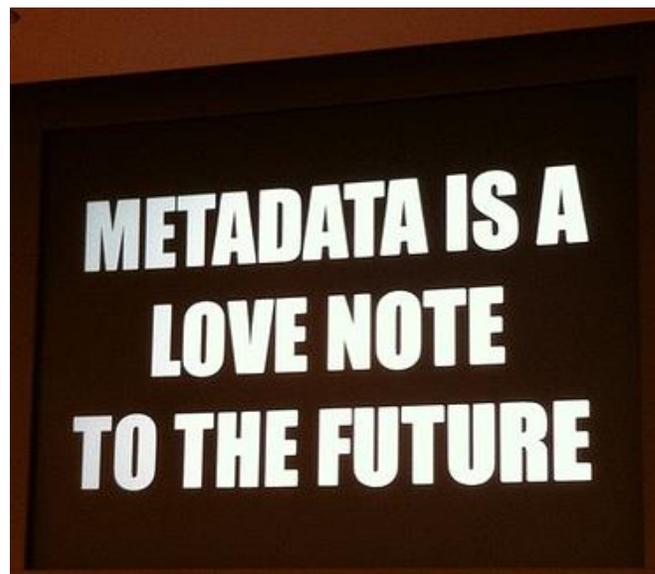


Photo by [@kissane](#); presentation by Jason Scott (@textfiles)



CC Sean MacEntee, Flickr



https://en.wikipedia.org/wiki/Data_sharing

The Data Creators' Advantage

	Comparative Data Reuse \leftrightarrow Integrative Data Reuse	
Goal	'Ground truthing:' calibrate, compare, confirm	Analysis: identify patterns, correlations, causal relationships
Example	Instrument calibration, sequence annotation, review summary-level data	Meta-analyses, novel statistical analyses
Frequency	Frequent, routine practice	Rare, emergent practice
Interpretation	Interactional expertise, 'knowledge that'	Contributory expertise, 'knowledge how,' tacit knowledge

Infrastructure: Durability



- Collaboration and openness
- International coordination
- Long-term value of data
- Agreed standards
 - Units of measurement
 - Coordinate systems
 - Data structures
- Shared resources
 - Missions, instruments
 - Data archives
 - Tools and technologies

Infrastructure: Fragility

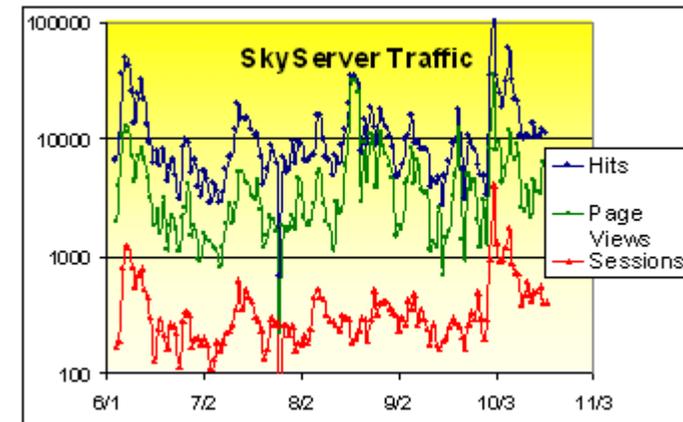
- Investments in data stewardship
 - Mission, instrument
 - Type of research
 - Space-based vs. ground based
 - Large missions vs. observing proposals
 - Shared vs. custom instruments
- Access to data
 - Public archives
 - Local websites
 - Derived data
- Curation investments
 - Open source software
 - Proprietary tools
 - Local pipelines, tools, scripts



Summary

Scientific Data and Infrastructure

- Infrastructures are fragile
- Visible infrastructure
 - Instruments
 - Institutions
- Invisible infrastructure
 - Data, metadata, provenance...
 - Information work
- Interdisciplinary science
 - Global science
 - Local practices



Telescope for the Sloan Digital Sky Survey, Apache Point, New Mexico

LSST All Hands Meeting, August 2014, Arizona State University. Arrow to Peter Darch

Data, Infrastructure, and Stewardship

- Whose data?
 - Global, comparative, fungible
 - Local, integrative, specific
- Whose infrastructure?
 - Funders, universities, companies
 - Individual investigators
- Whose stewardship?
 - Maintain collections, models, instruments, technology, code...
 - Invest in people, skills, collaborations





Alberto Pepe, David Fearon, Katie Shilton, Jillian Wallis, Christine Borgman, Matthew Mayernik (2009)



Christine Borgman



Peter Darch



Ashley Sands



Irene Pasquetto



Milena Golshan



Bernie Boscoe



Cheryl Thompson



Morgan Wofford



Michael Scroggins



Sharon Traweek