

UCLA

Presentations

Title

Data, infrastructure, and stewardship

Permalink

<https://escholarship.org/uc/item/7m59x6mh>

Author

Borgman, Christine L.

Publication Date

2019-11-07

Copyright Information

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

Data, Infrastructure, and Stewardship

Christine L. Borgman

Distinguished Research Professor

Director, Center for Knowledge Infrastructures

<https://knowledgeinfrastructures.gseis.ucla.edu>

University of California, Los Angeles

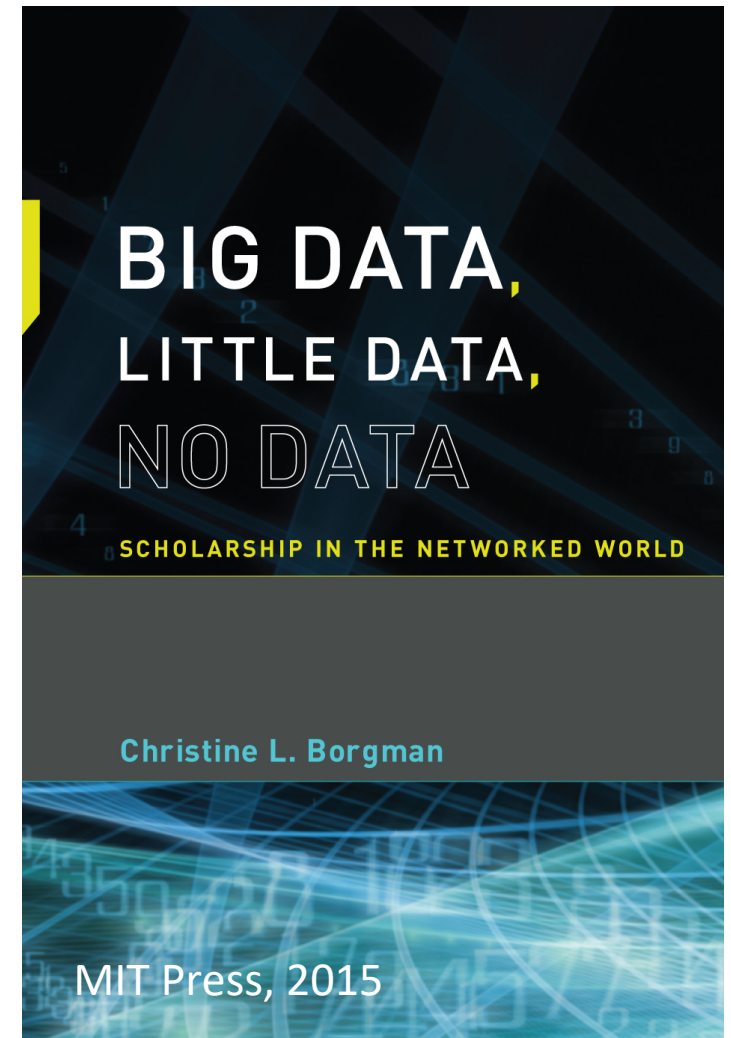
<http://christineborgman.info>

@scitechprof

Space Studies Board, National Academies of Sciences

Keynote Presentation

Beckman Center, Irvine, 7 November 2019





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



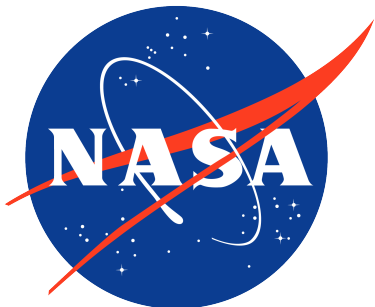
Supported by
wellcometrust



Australian Government
National Health and Medical Research Council



National Science Foundation
WHERE DISCOVERIES BEGIN





Open Data Practices



- Deposit datasets in a data archive
- Link datasets to journal article or publication
- Publish data documentation



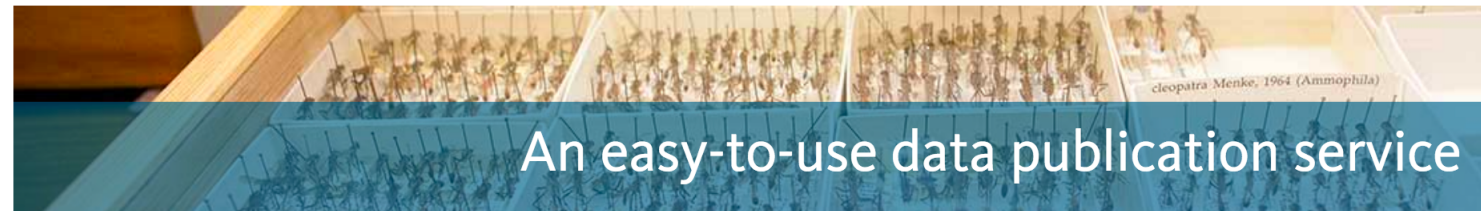
PDS: The Planetary Data System



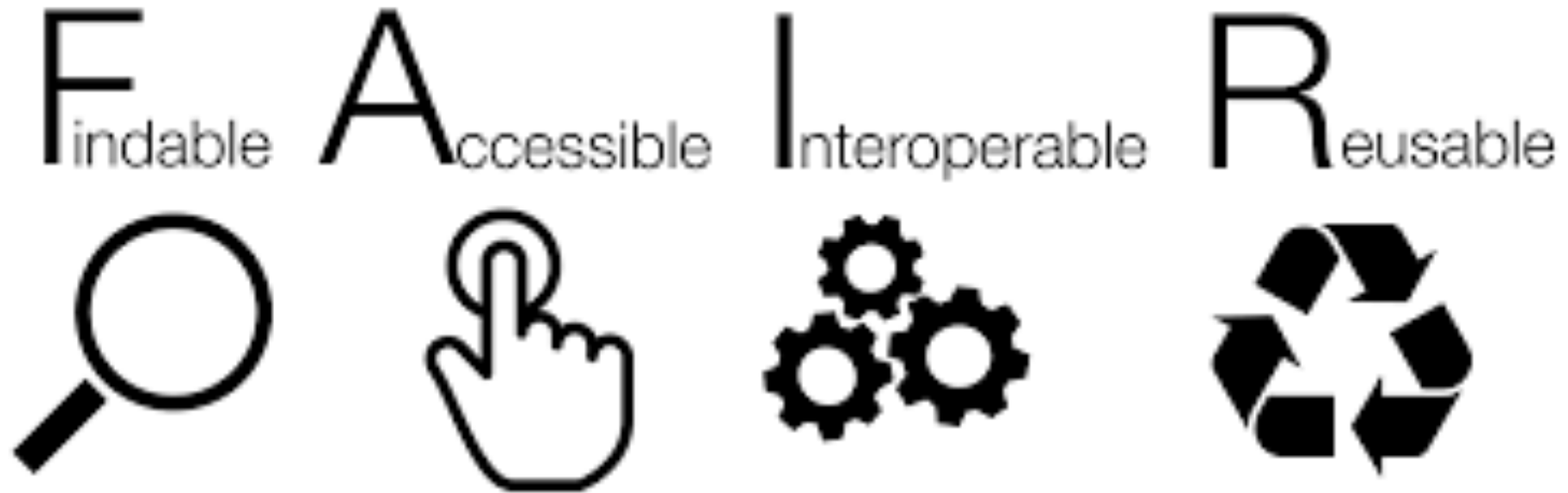
- Research protocols
- Codebooks
- Software
- Algorithms



- Cite data and software



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.

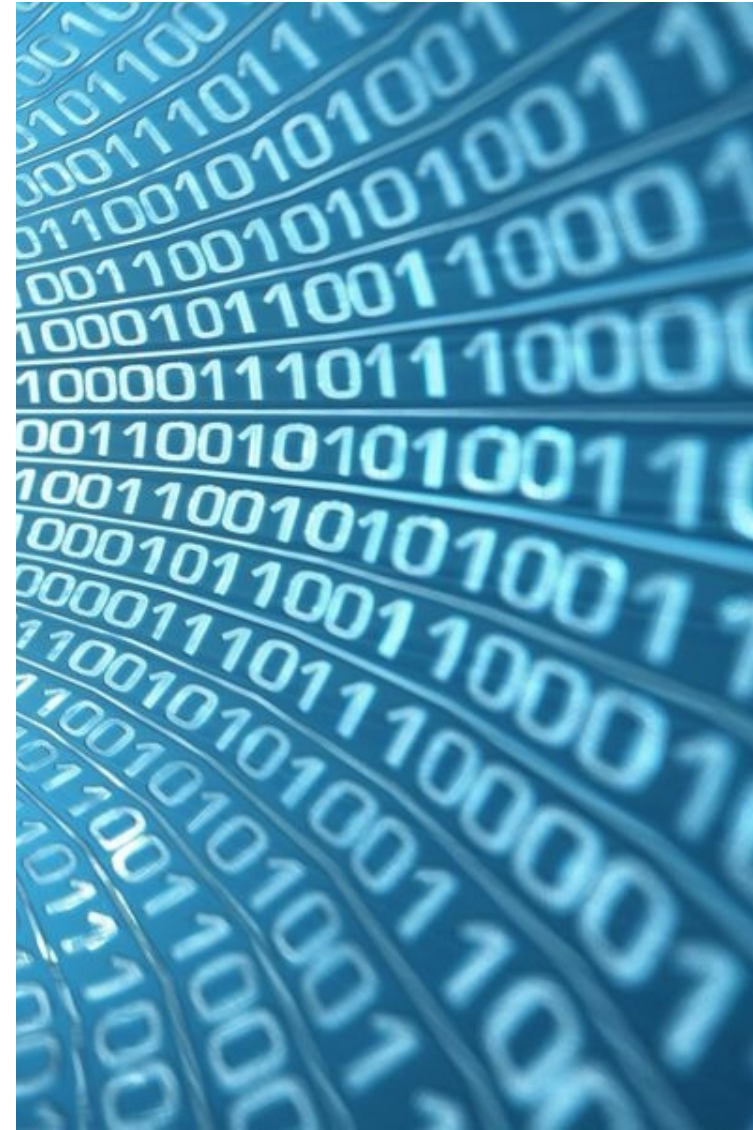


National Aeronautics and Space Administration
NASA Official: Brian Dunbar

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

Data Challenges in Space Studies

- 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?





[Our
Scientists](#)

[Core
Research](#)

[Education
& Outreach](#)

[About
Us](#)

[Donate
Now](#)

Could We Be Alone in the Cosmos?

Data

Cassini-Huygens: Mission to Saturn BY THE NUMBERS

2.5 MILLION
COMMANDS
executed

635 
SCIENCE DATA
collected

6 NAMED MOONS
discovered

162 TARGETED
FLYBYS
of Saturn's moons

27 NATIONS
participated

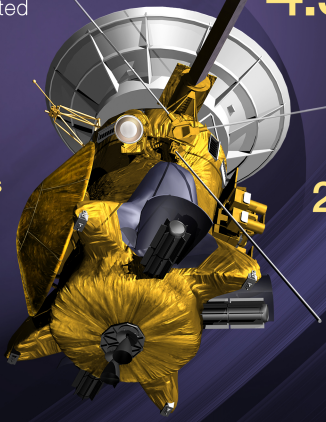
4.9 BILLION
MILES TRAVELED
since launch
(7.9 BILLION KILOMETERS)

3,948
SCIENCE PAPERS
published

294 ORBITS
completed

453,048
images taken

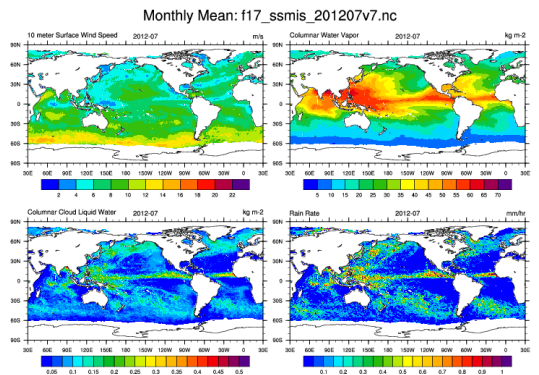
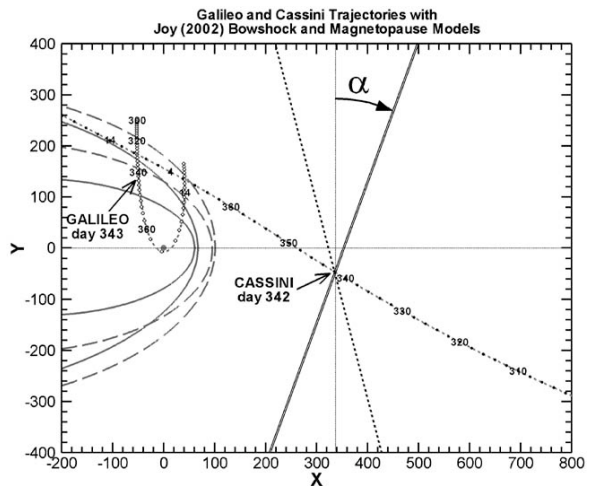
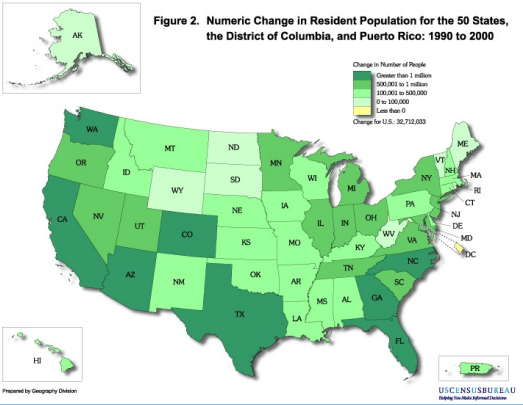
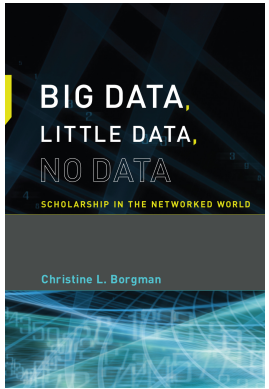
360 ENGINE
burns



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.



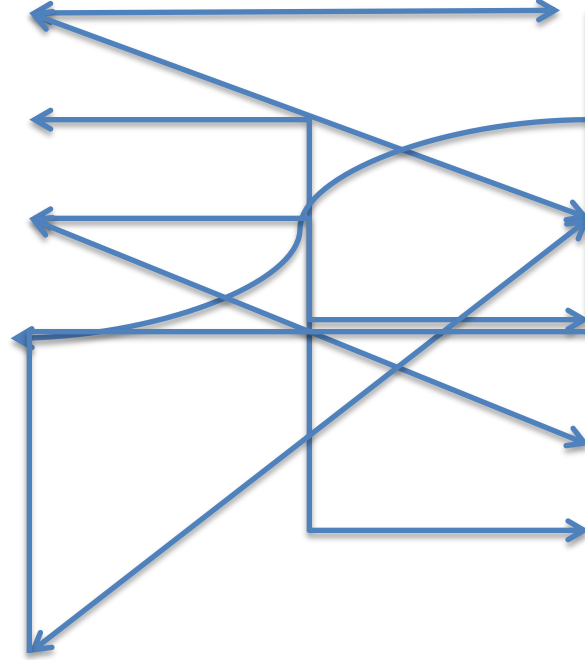
Kivelson, M. G., & Southwood, D. J. (2003). First evidence of IMF control of Jovian magnetospheric boundary locations: Cassini and Galileo magnetic field measurements compared. *Planetary and Space Science*, 51(13), 891–898. [https://doi.org/10.1016/S0032-0633\(03\)00075-8](https://doi.org/10.1016/S0032-0633(03)00075-8)



Publications \leftrightarrow Data: Mapping

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

- Article n

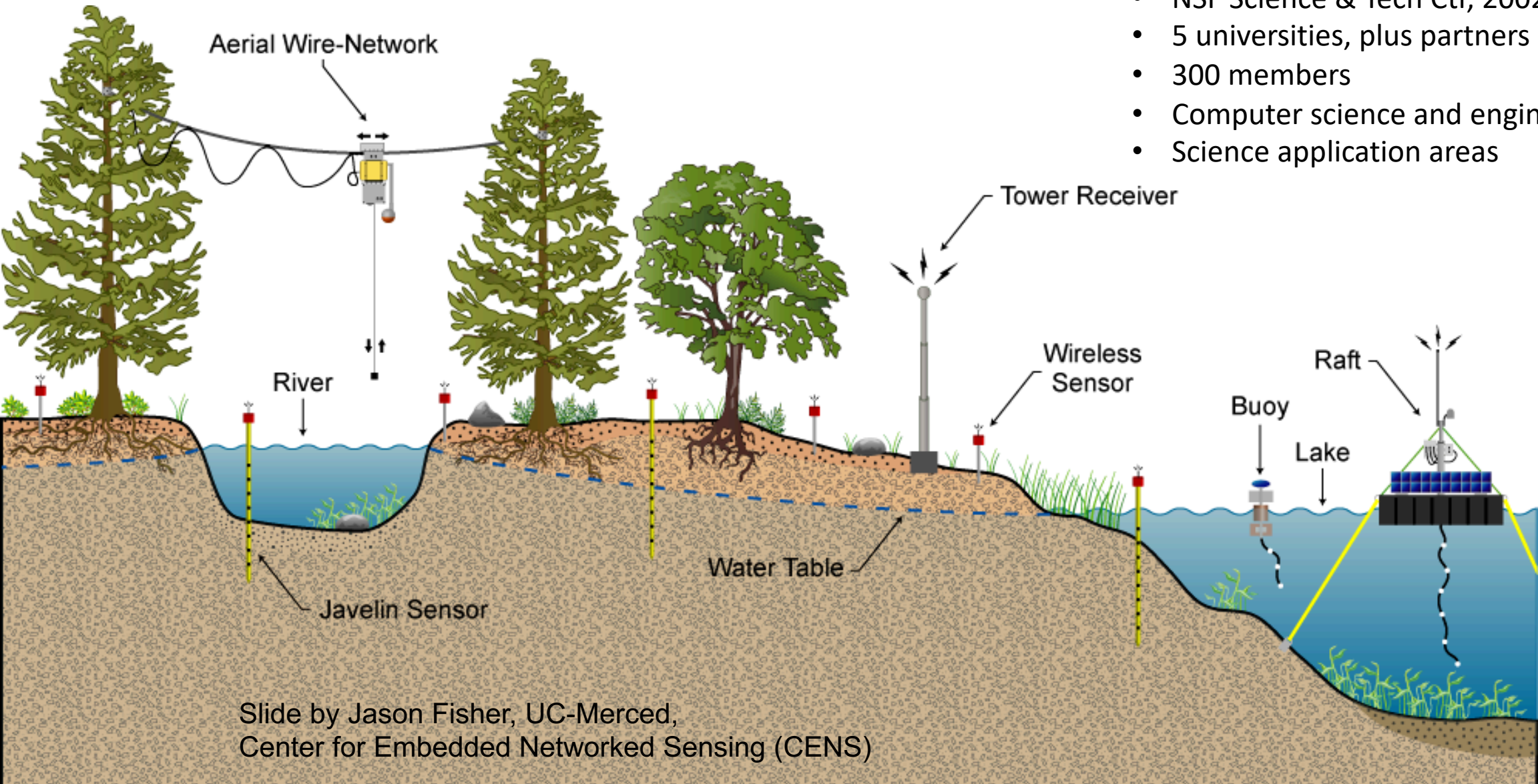


- Dataset time 1
- Dataset time 2
- Observation time 1
- Visualization time 3
- Community collection 1
- Repository 1

Data practices

Center for Embedded Networked Sensing

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

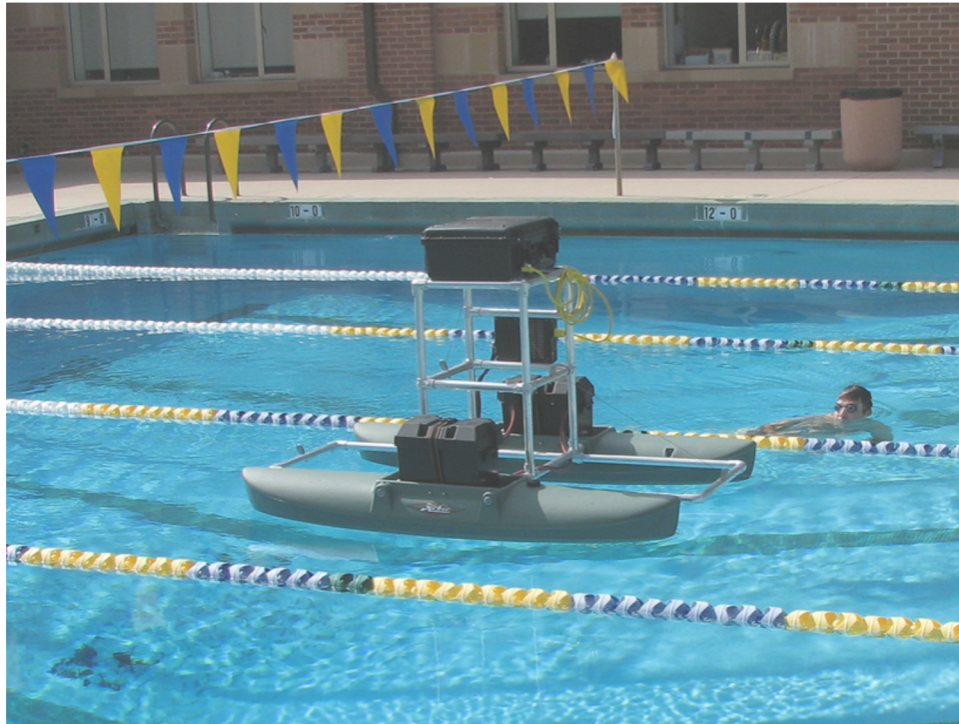


Slide by Jason Fisher, UC-Merced,
Center for Embedded Networked Sensing (CENS)

Science \leftrightarrow Data

Engineering researcher:

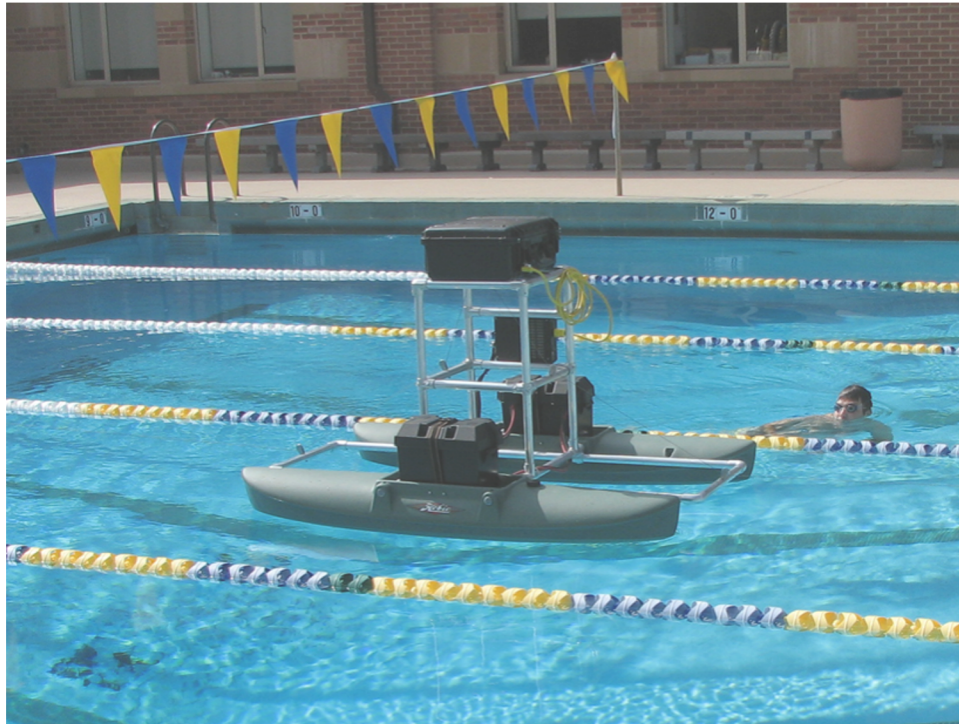
“Temperature is temperature.”



CENS Robotics team

Science \leftrightarrow Data

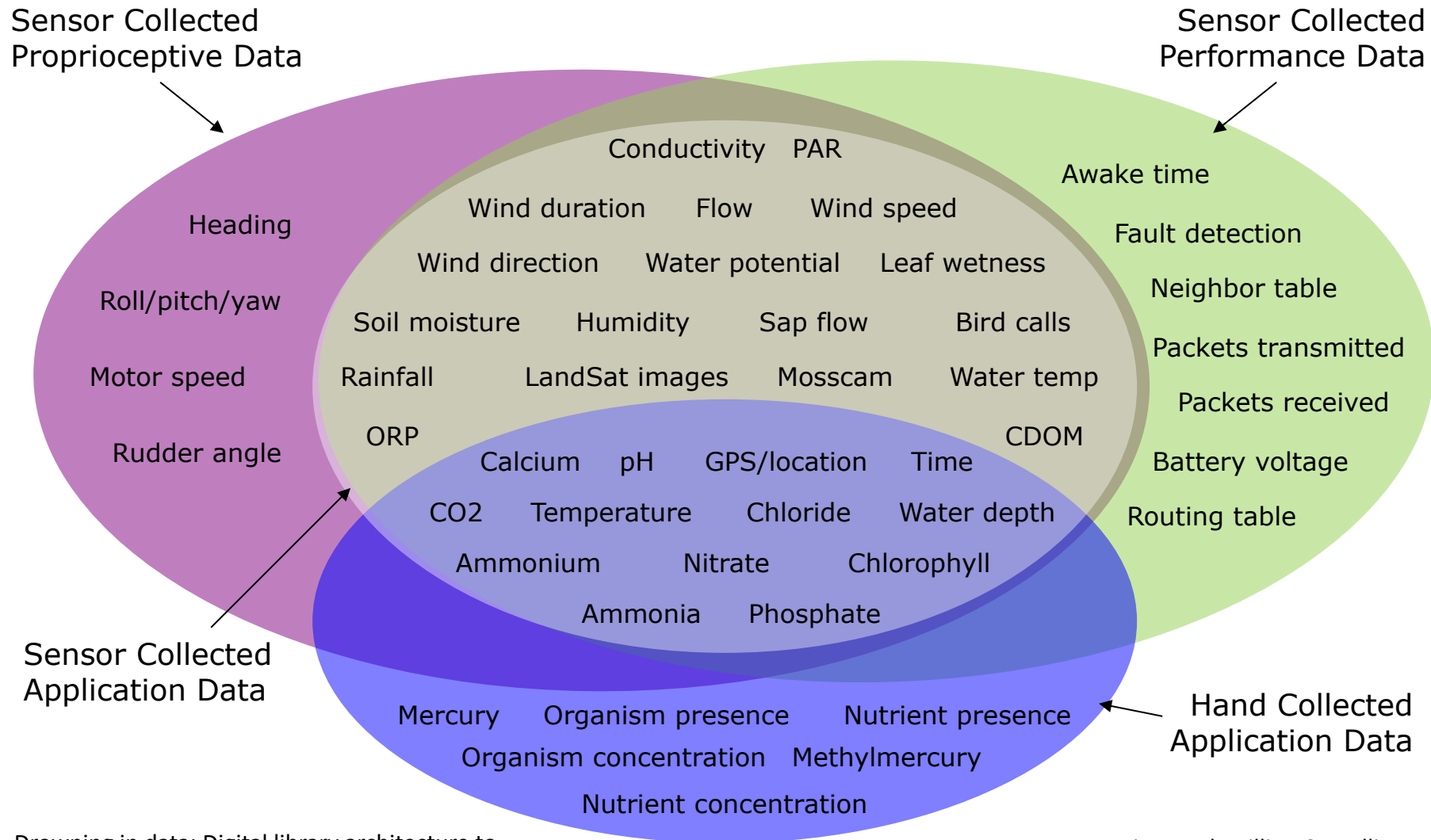
Engineering researcher:
“Temperature is temperature.”



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..”

CENS data variation



Deep Subseafloor Biosphere

- Center for Dark Energy Biosphere Investigations (C-DEBI)
- International Ocean Discovery Program (IODP)
- Microbial communities in the seafloor
- Highly multidisciplinary



Darch, P. T., & Borgman, C. L. (2016). Ship space to database: Emerging infrastructures for studies of the deep subseafloor biosphere. *PeerJ Computer Science*, 2, e97.

<https://doi.org/10.7717/peerj-cs.97>

Center for Dark Energy Biosphere Investigations



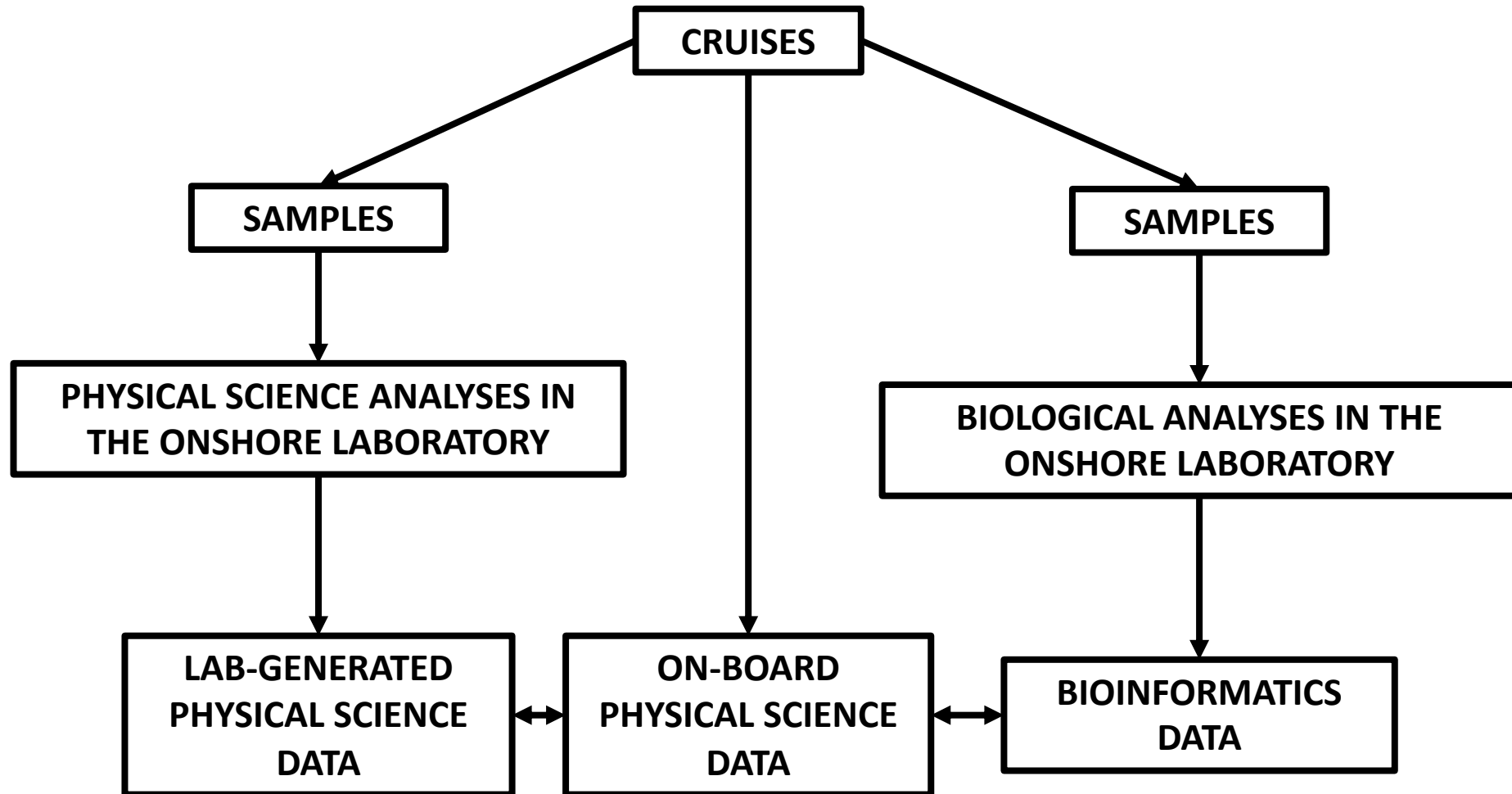
International Ocean Discovery Program
lodp.tamu.org

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

Repository for seafloor cores. Photo: Peter Darch

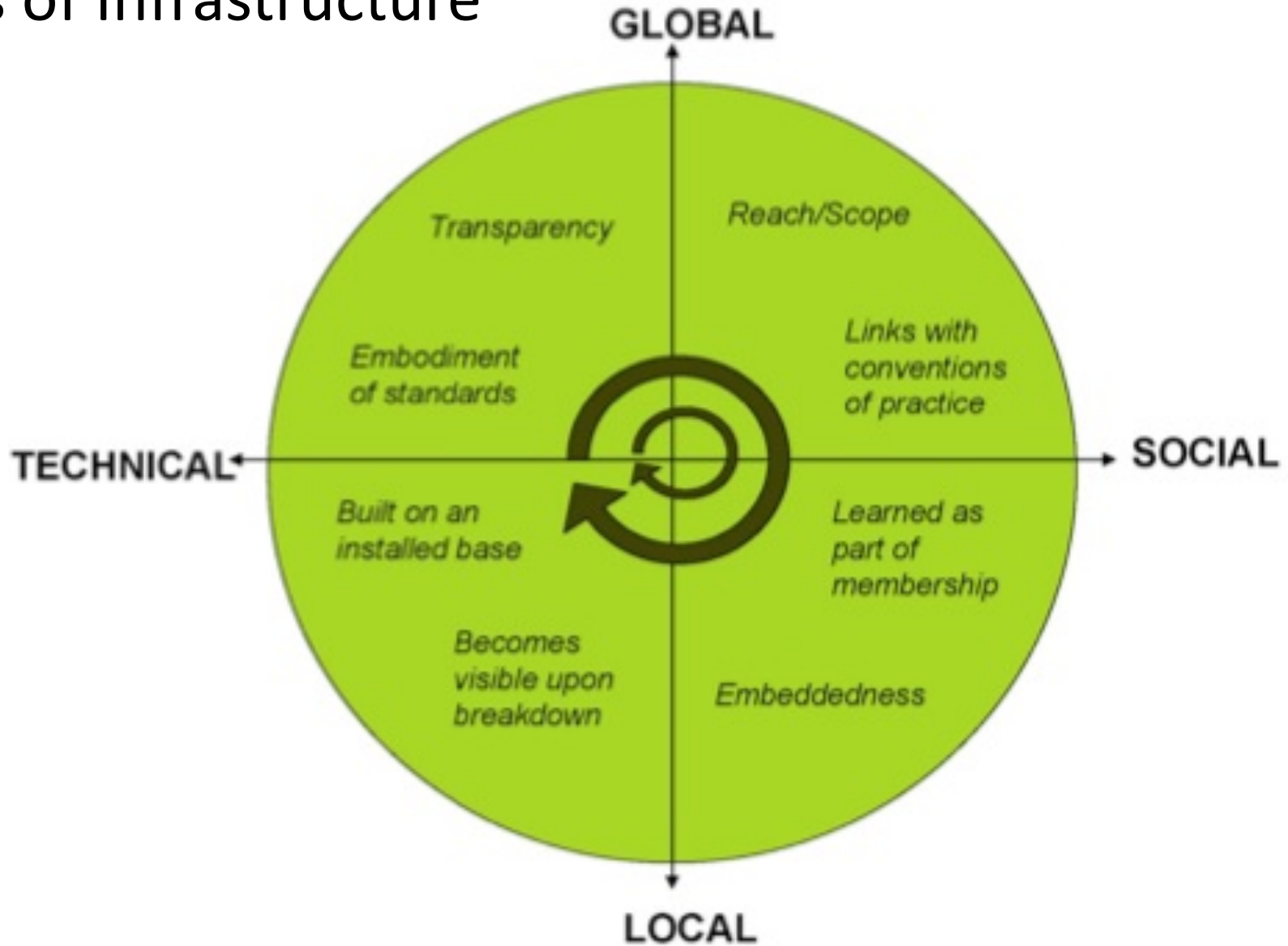


Data Diverge During Scientific Work



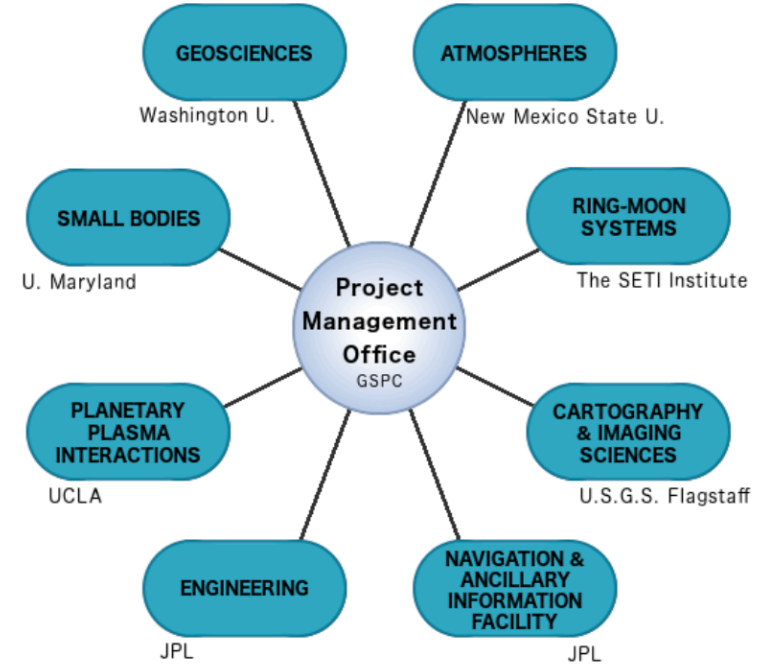
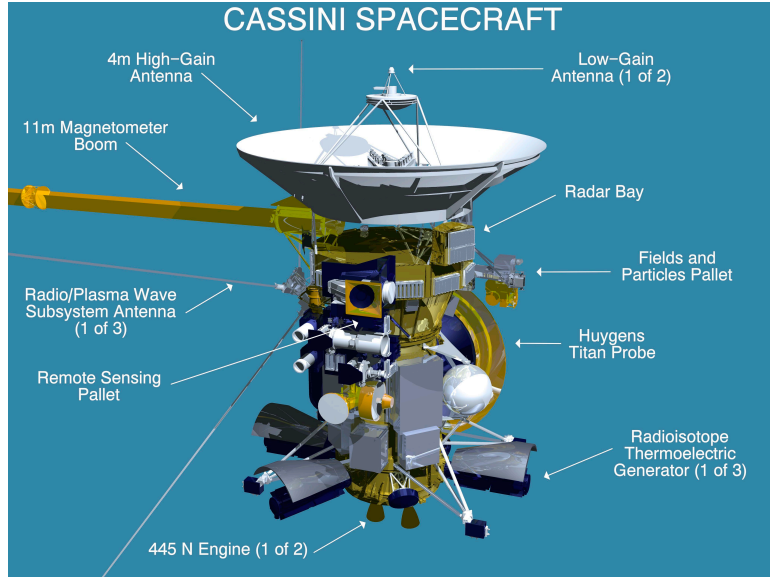
Infrastructure

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>

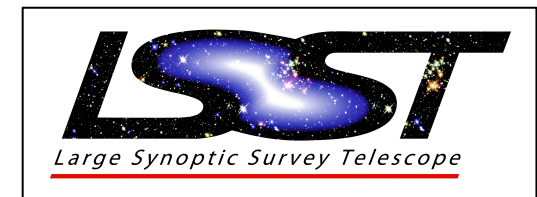
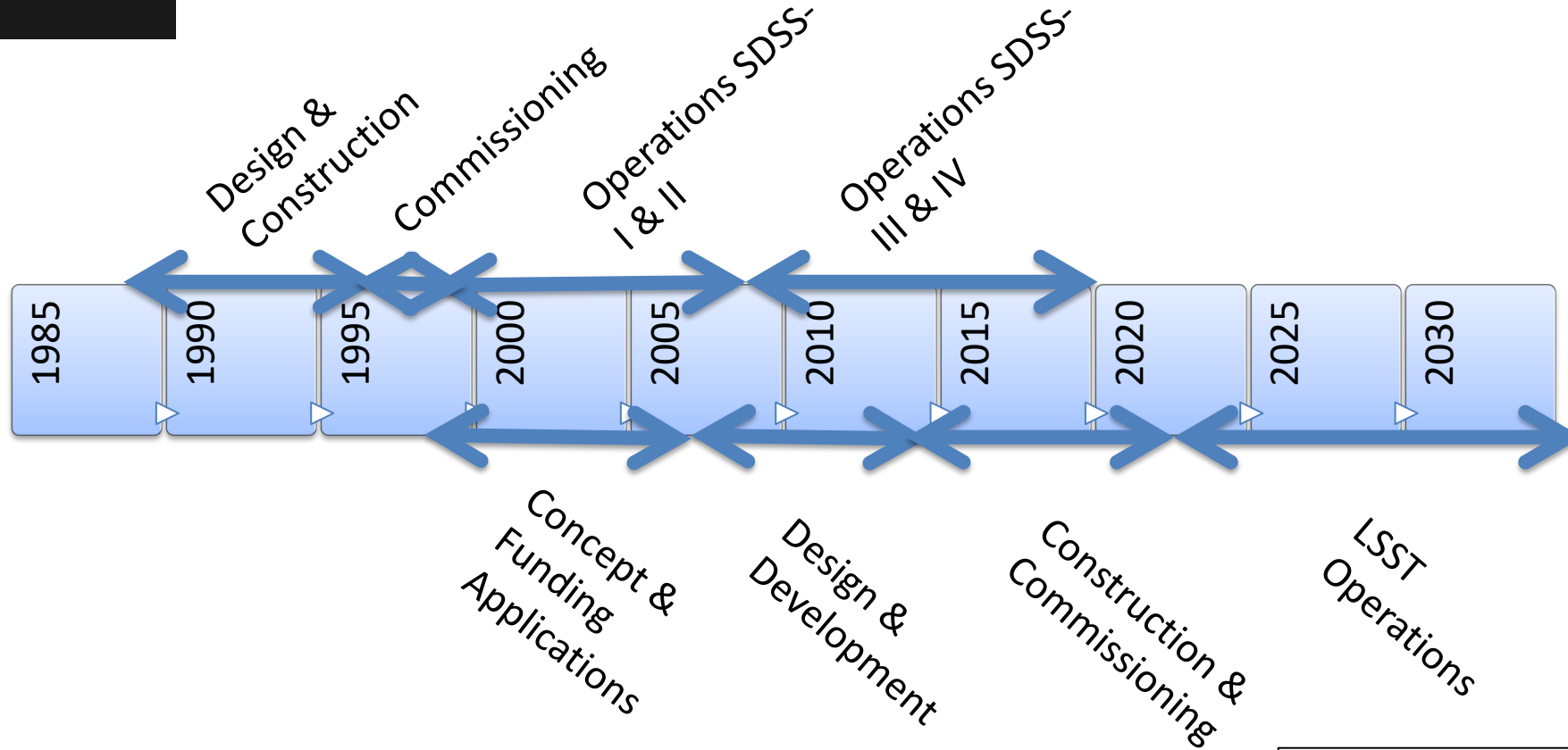
Global and Technical



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



Project Timelines



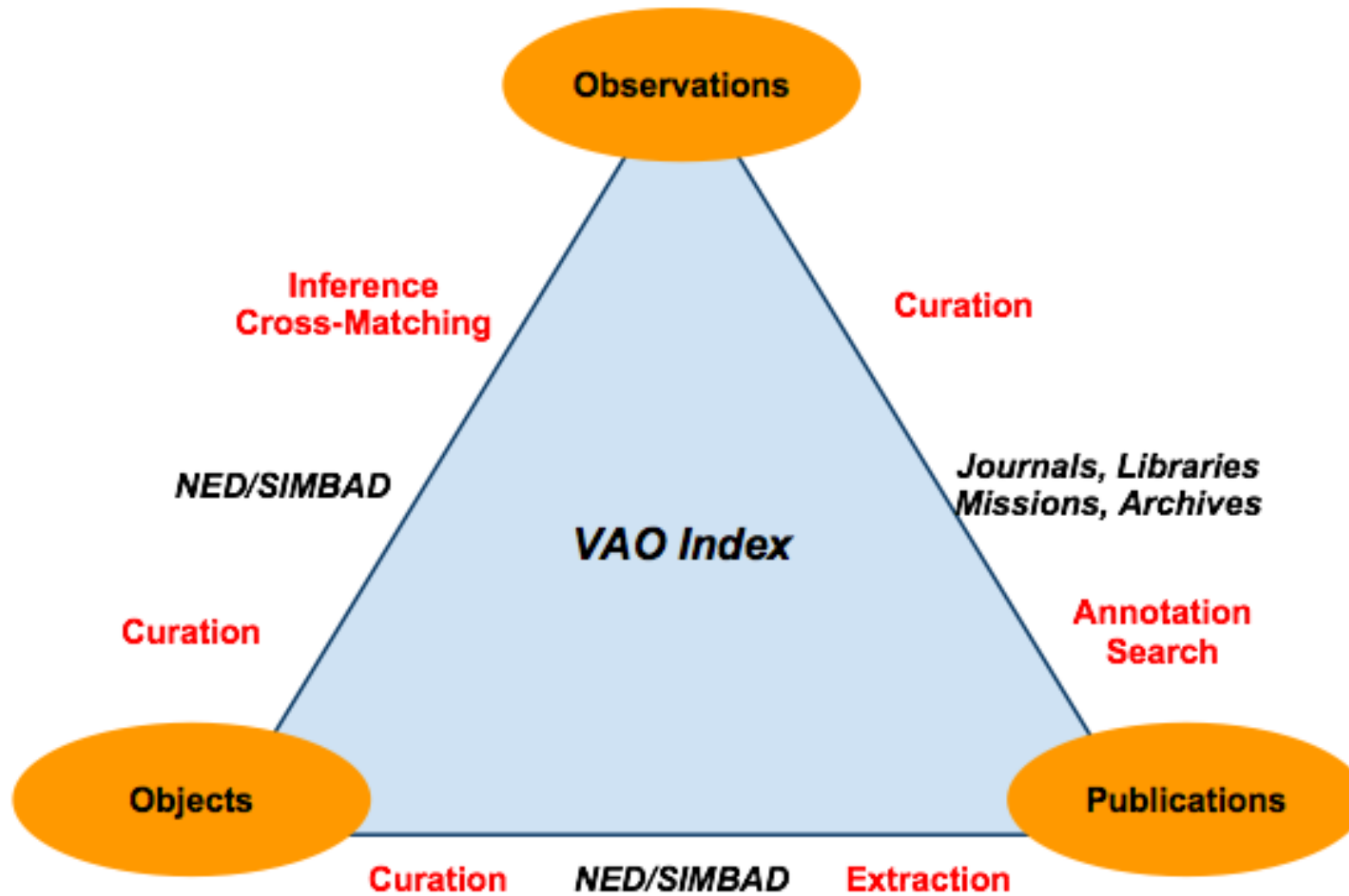
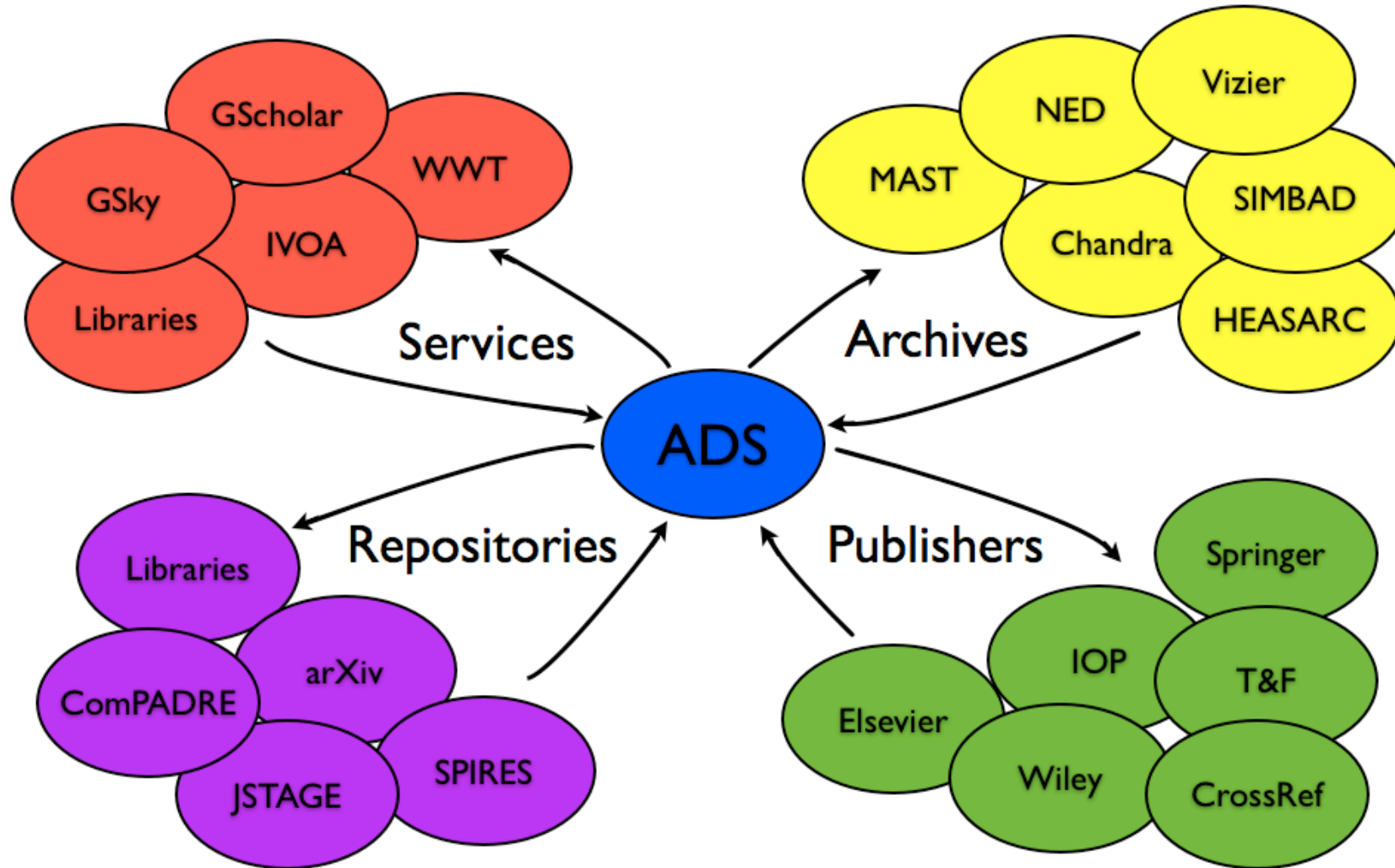


Figure 1. Relationships between Publications, Objects, Observations and the corresponding major actors in the curating process and their activities (in red).

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>

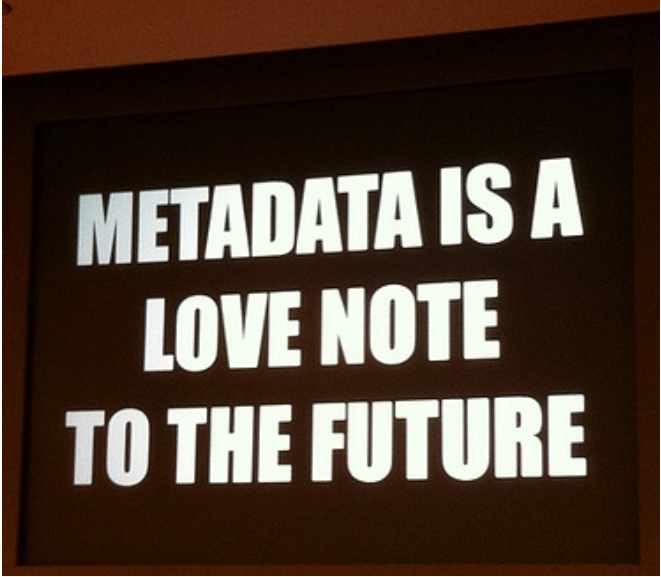
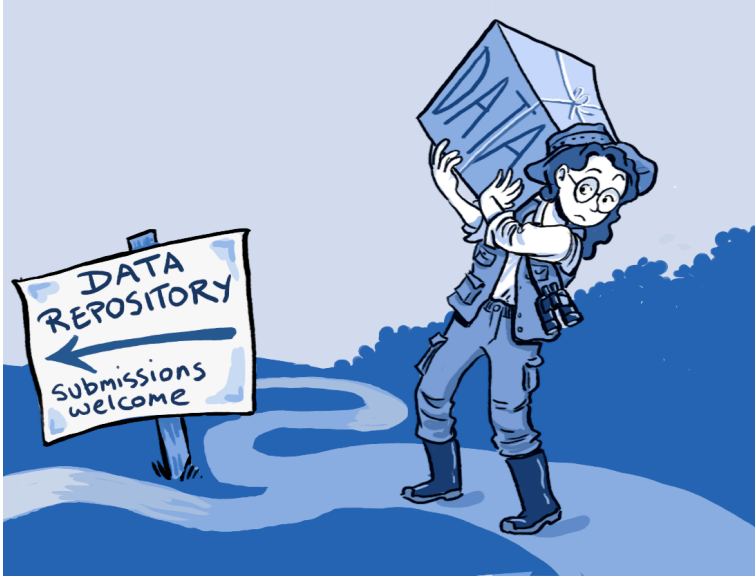
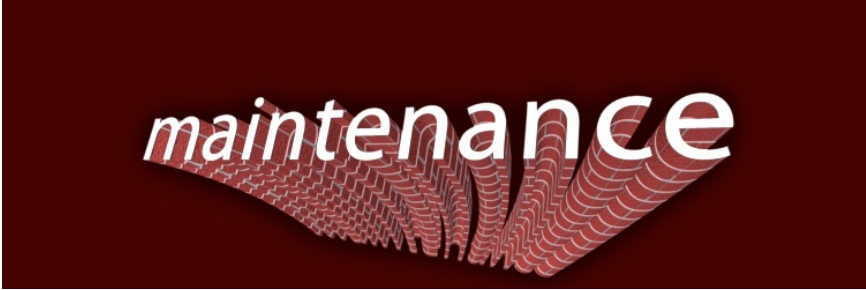


Photo by [@kissane](#); presentation by Jason Scott (@textfiles)



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



CC Sean MacEntee, Flickr

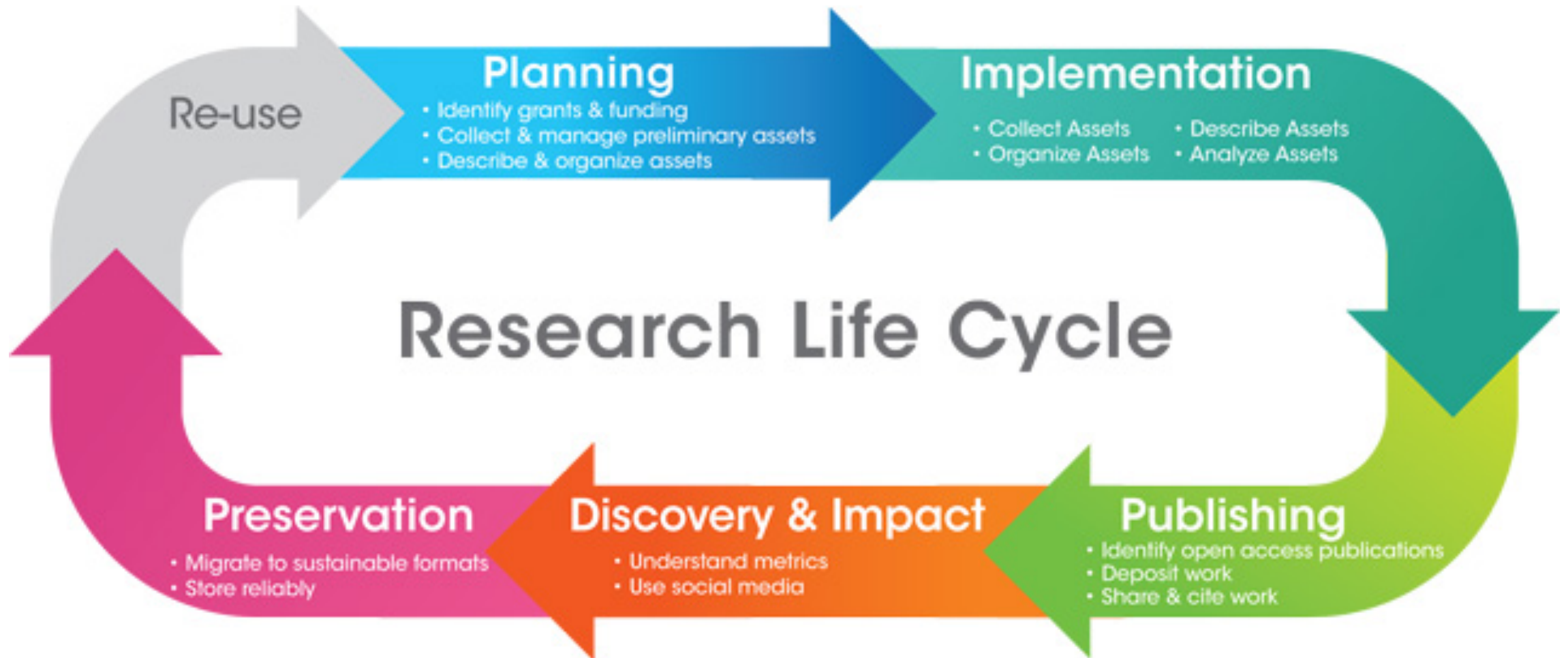
Lack of incentives to share data

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



	Comparative Data Reuse <--> 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

Data creation and reuse: The Ideal



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 ²⁹

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

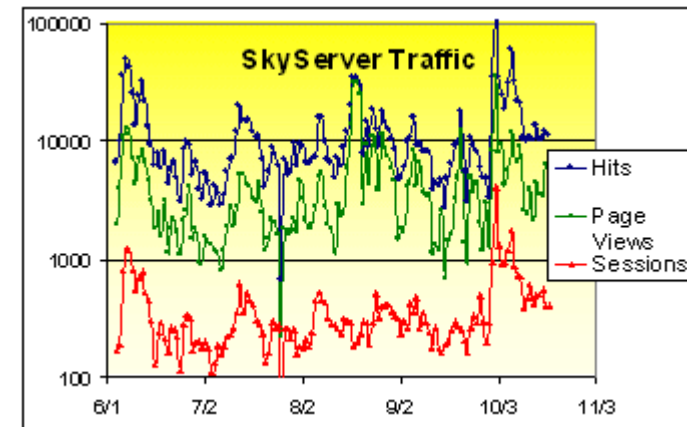
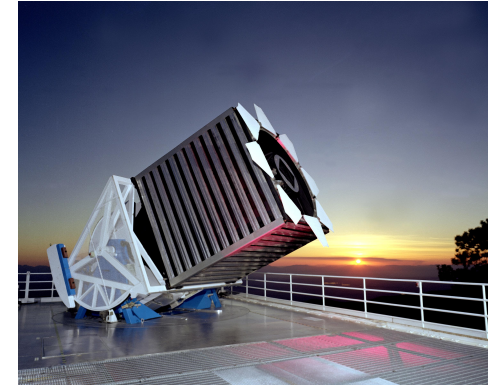
- Investments in data stewardship
 - Mission, instrument
 - Type of research
 - Space-based vs. ground based
 - Large missions vs. observing proposals
 - Shared vs. custom instruments
- Access to data
 - Public archives
 - Local websites
 - Derived data
- Curation investments
 - Open source
 - Proprietary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



Discussion

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



Acknowledgements

UCLA Center for Knowledge Infrastructures



Christine Borgman



Bernie Boscoe



Peter Darch



Milena Golshan



Irene Pasquetto



Michael Scroggins



Cheryl Thompson



Morgan Wofford