# UCLA
## Presentations

**Title**
If Data Sharing is the Answer, What is the Question?

**Permalink**
https://escholarship.org/uc/item/9dx15801

**Author**
Borgman, Christine L.

**Publication Date**
2016-09-13

# If Data Sharing is the Answer, What is the Question?

## Christine L. Borgman

Distinguished Professor and Presidential Chair in Information Studies

University of California, Los Angeles

http://christineborgman.info

https://knowledgeinfrastructures.gseis.ucla.edu

@scitechprof

Closing Keynote

SciDataCon, Denver, CO

September 13, 2016



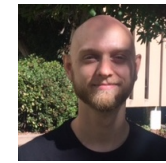Christine Borgman    Peter Darch    Ashley Sands    Irene Pasquetto

Bernie Randles    Milena Golshan    Pietro Santachiara

**UCLA** Center for Knowledge Infrastructures

# Data sharing policies

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

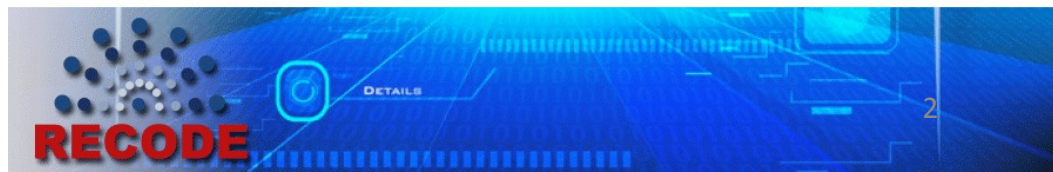Policy RECommendations for Open Access to Research Data in Europe

2

# Why Share Research Data?

- To reproduce research
- To make public assets available to the public
- To leverage investments in research
- To advance research and innovation

MIT Press, 2015
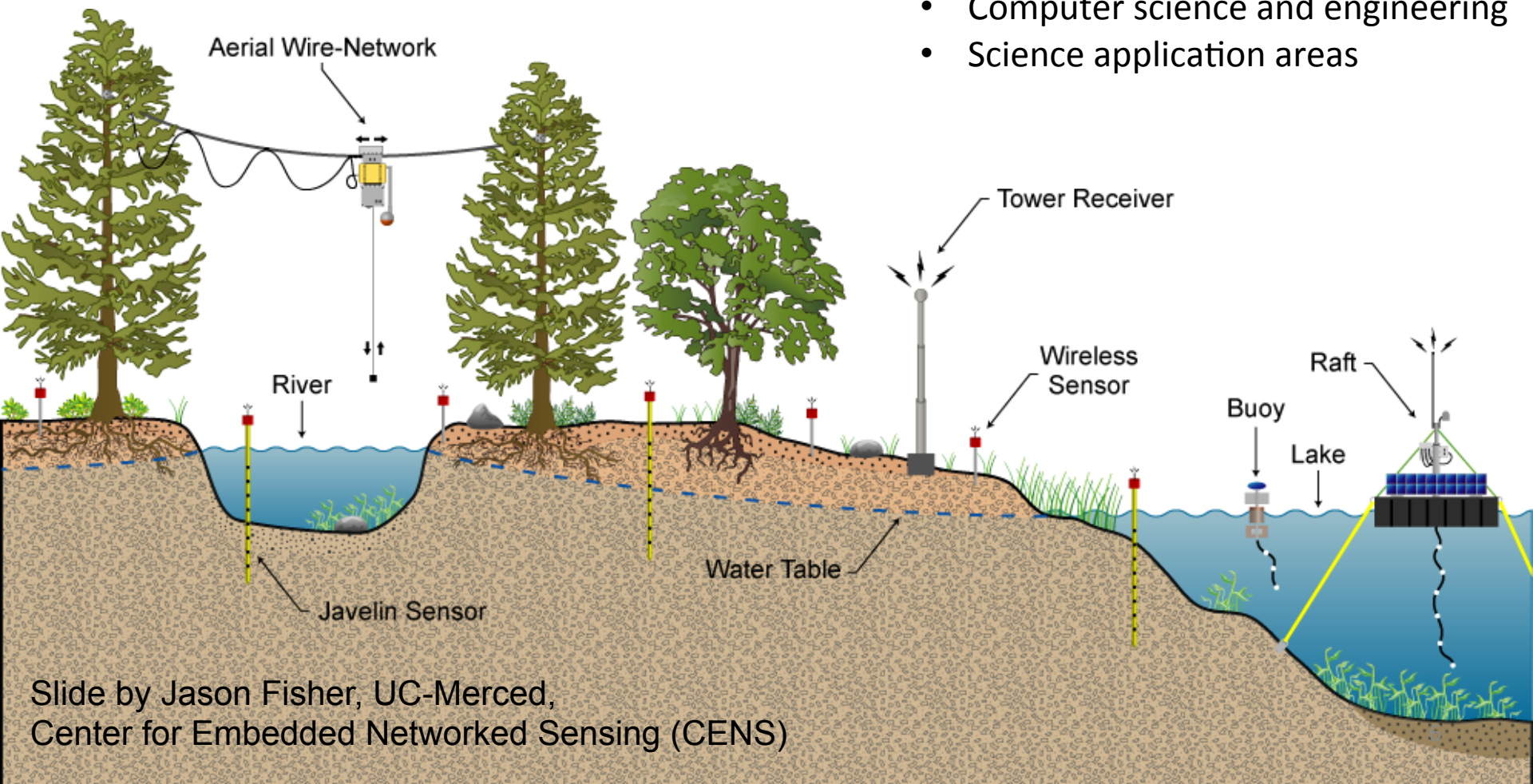
# Lack of incentives to share data



- Rewards for publication
- Effort to document data
- Competition, priority
- Control, ownership

http://www.buildingsrus.co.uk/.../ target1.htm
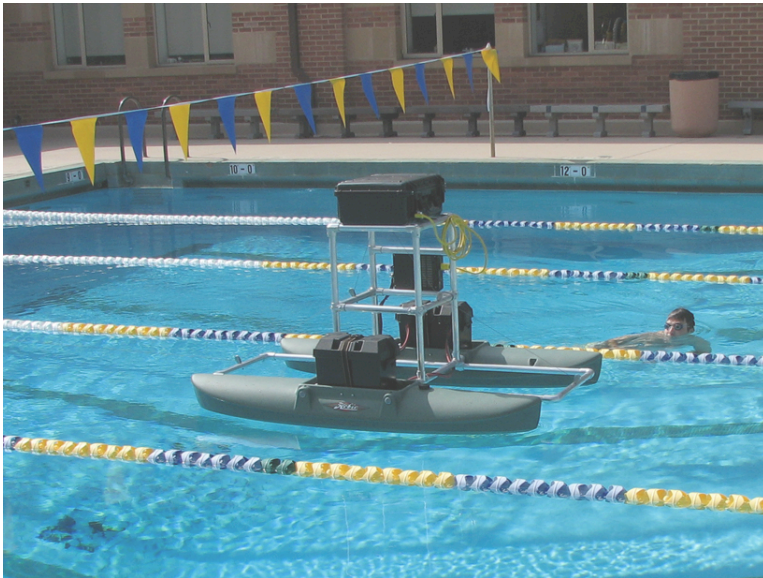
# Data

# Center for Embedded Networked Sensing

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



Aerial Wire-Network

Tower Receiver

River

Wireless Sensor

Raft

Buoy

Lake

Javelin Sensor

Water Table

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

# Documenting Data for Interpretation

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.."*

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

C.L. Borgman (2015). *Big Data, Little Data, No Data: Scholarship in the Networked World*. MIT Press

8

# If Data Sharing is the Answer, What is the Question?

- Research Design, 2015-2018

- Methods

- Questions

- Findings

- Conclusions

- DANS study

- Recommendations

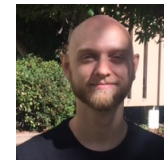Christine Borgman  Peter Darch  Ashley Sands  Irene Pasquetto

Bernie Randles  Milena Golshan  Pietro Santachiara

**UCLA** Center for Knowledge Infrastructures

# Research Design

- Goals
  - Explicate data, sharing, reuse, openness, infrastructure across scientific domains
  - Identify new models of scientific practice

- Dimensions
  - Mixtures of domain expertise
  - Factors of scale
  - Centralization of data collection and analysis

**UCLA** Center for Knowledge Infrastructures

# Qualitative Methods

- Document analysis
  - Public and private documents and artifacts
  - Official and unofficial versions of scientific practice
- Ethnography
  - Observing activities on site and online
  - Embedded for days or months at a time
- Interviews
  - Questions based on our research themes
  - Compare multiple sites over time

**UCLA** Center for Knowledge Infrastructures

# Current Research Sites

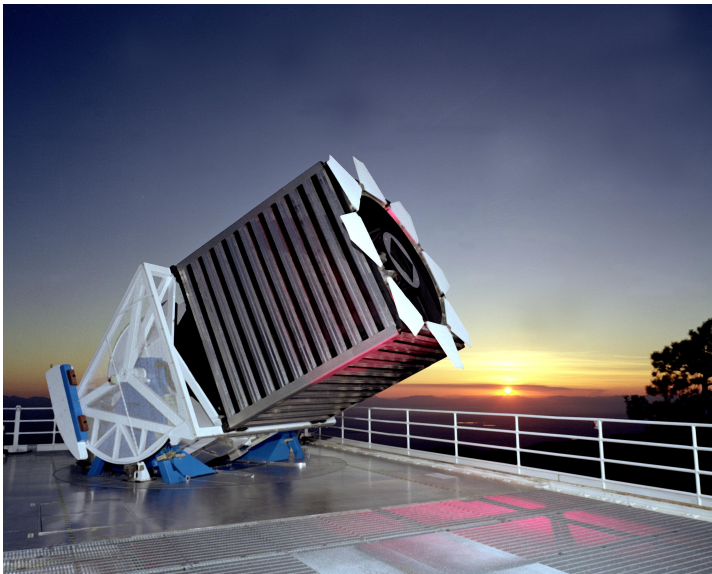| Domain | Focus | Topic |
|---|---|---|
| Astronomy sky surveys | Place: sky and universe | Survey of night sky |
| Deep subseafloor biosphere | Place: under ocean floor | Microbial life and environment |
| Craniofacial research | Problem: Craniofacial birth defects in humans | Genomics of four model organisms |
| Computational science | Problem: Data analysis at scale | Computing platform for astronomy, physics, turbulence, soil science, genomics… |

# Research Question 1

How do the *mixtures of domain expertise* influence the collection, use, and reuse of data – and vice versa?

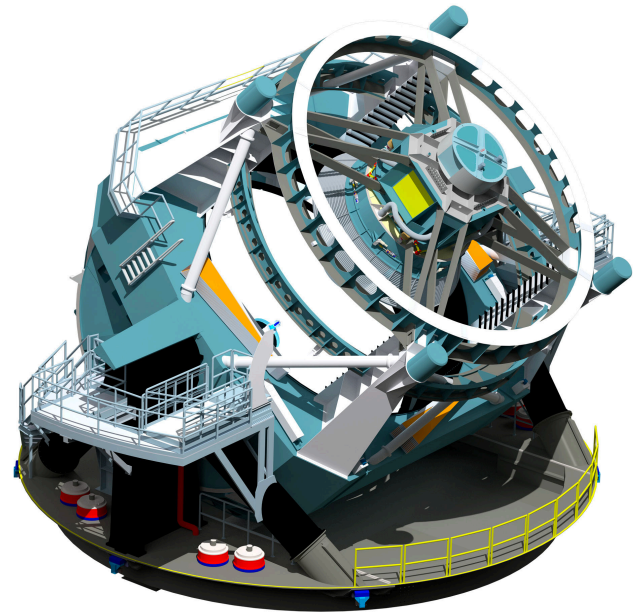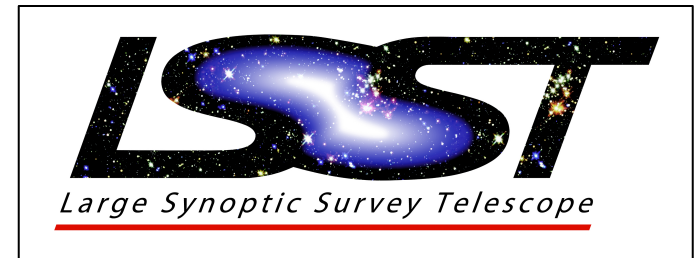| Domain |
| --- |
| Astronomy sky surveys |
| Deep subseafloor biosphere |
| Craniofacial research |
| Computational science |

**UCLA** Center for Knowledge Infrastructures

# Sloan Digital Sky Survey (SDSS-I/II)



- Survey from 2000-2008
- 160+ TB data total
- Tens of millions of dollars
- Open data
- Proprietary software



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

# Large Synoptic Survey Telescope (LSST)

- Survey from 2022-2032

- 15 TB data per night

- 1+ Billion dollars

- Data open to partners

- Open source software





https://news.slac.stanford.edu/sites/default/files/images/image/lsst_h_0.jpg

LSST telescope, Chile

# Mixtures: Astronomy sky surveys

- Domains
  - Astronomy
  - Computer science
- Project characteristics
  - Mature discipline
  - Abundant data
  - Trusted archives
  - Shared tools, methods
  - Established infrastructure for data access and use

# Center for Dark Energy Biosphere Investigations





Repository for seafloor cores. Photo: Peter Darch

International Ocean Discovery Program
lodp.tamu.org

- NSF Science & Tech Ctr, 2010-2020
- 35 institutions
- 90 scientists
- Biological sciences
- Physical sciences

17

18

# Mixtures: Deep subseafloor biosphere

- Domains
  - Biological sciences
  - Physical sciences
  - 50+ self-identified specialties
- Project characteristics
  - Emergent scientific problem area
  - Scarce data
  - Disparate, exploratory methods
  - Building capacity for data collection
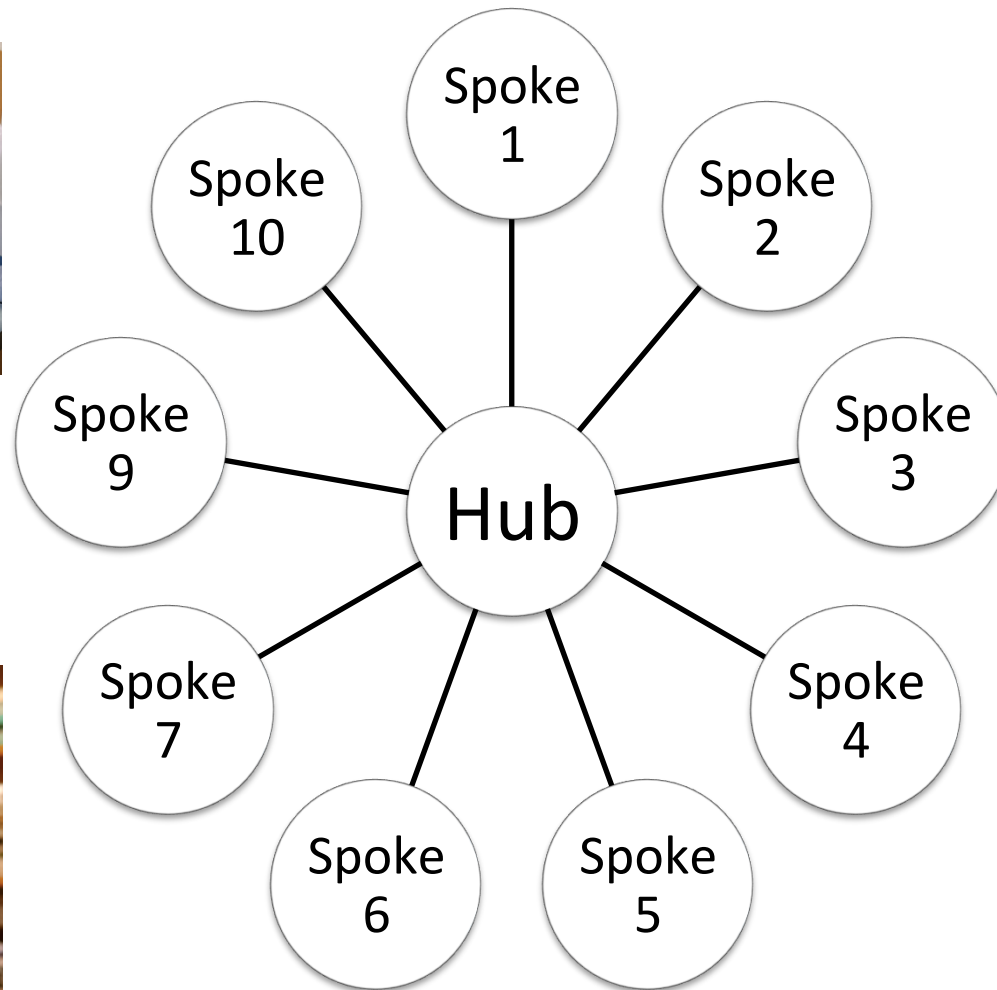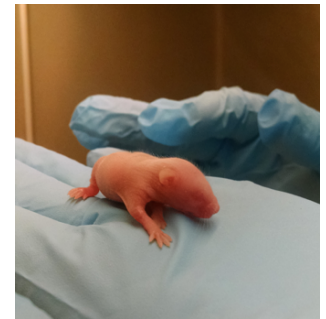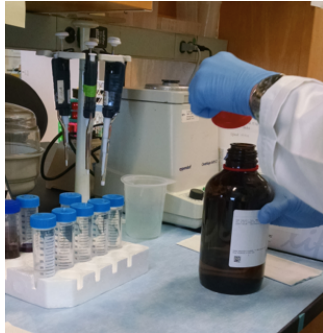  - Sharing established infrastructures
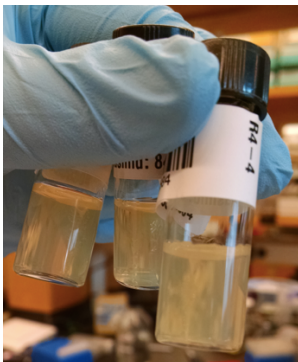
# FaceBase Consortium

- National Institute for Dental and Craniofacial Research
- Genetics, imaging data: craniofacial development
- 11 projects: clinical, biology, bioinformatics
- 4 model organisms: human, primates, mice, zebrafish
- Make data available on hub **www.facebase.org**

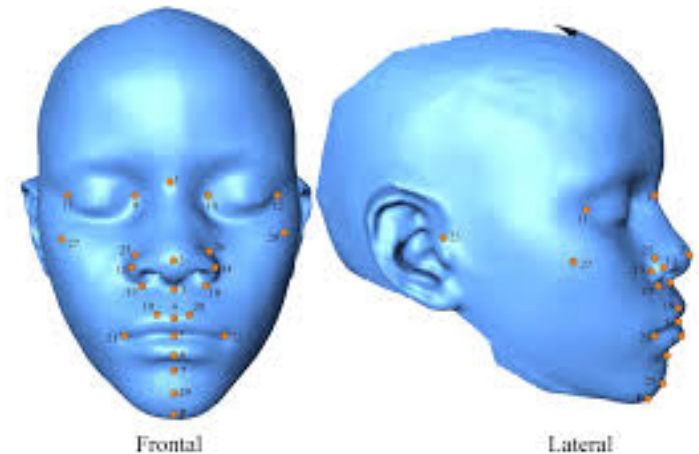

Frontal          Lateral

# FaceBase Spokes and Hub

1 coordinating center
10 spokes

# Mixtures: Craniofacial deformities

- Domains
  - Genomics, bioinformatics
  - Molecular, developmental biology
  - Dentistry, plastic surgery
- Project characteristics
  - Urgent medical problem
  - Species-specific data
    - Humans
    - Primates
    - Mice
    - Zebrafish
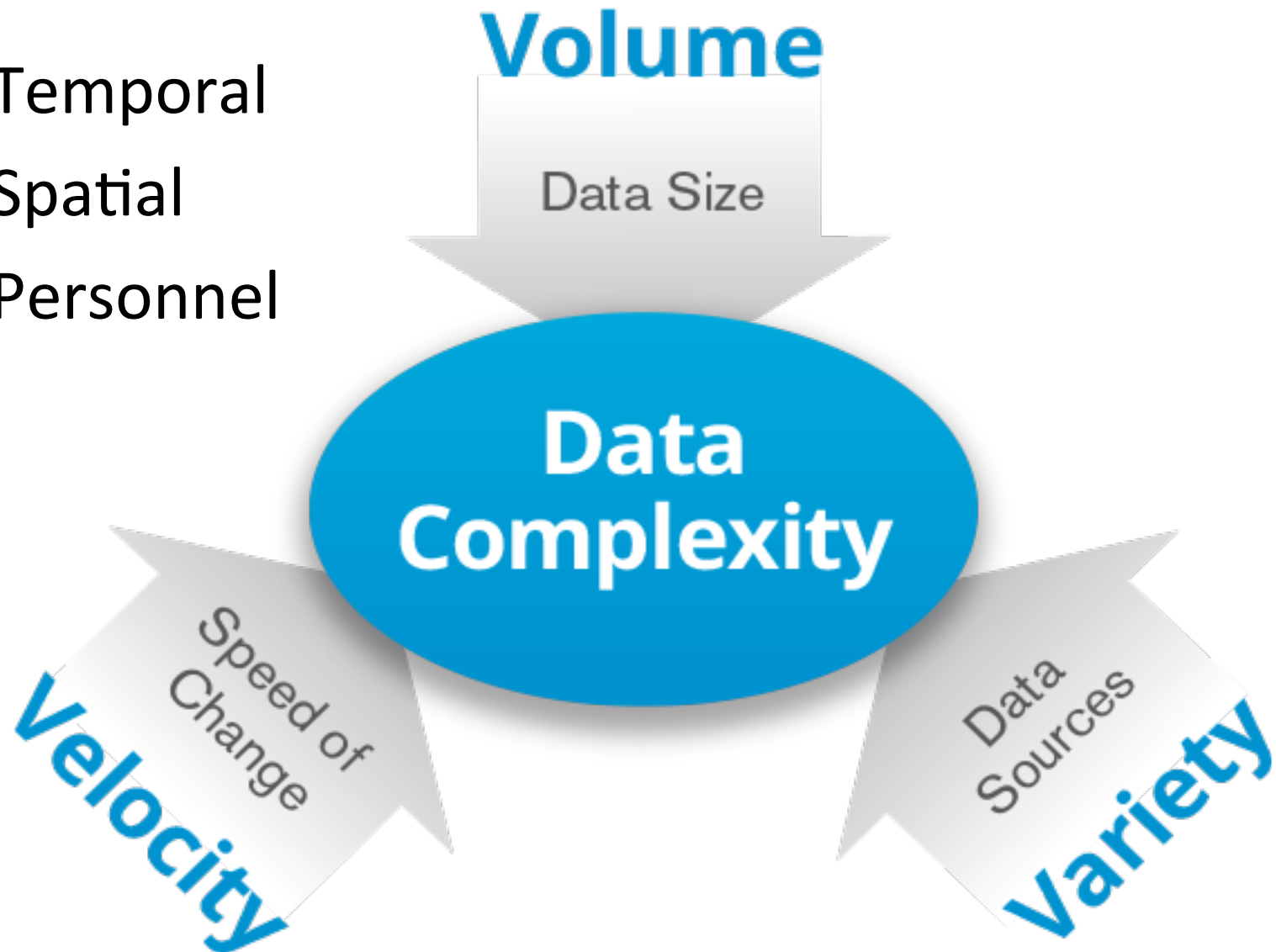  - Competing tools, methods
  - Multiple established infrastructures



Frontal    Lateral



FaceBase

# Research Question 2

What *factors of scale* influence research practices, and how?

| Domain |
| --- |
| Astronomy sky surveys |
| Deep subseafloor biosphere |
| Craniofacial research |
| Computational science |

**UCLA** Center for Knowledge Infrastructures

# *Scale factors*

- Temporal
- Spatial
- Personnel

# Project Timelines

# Scale factors

| Research site | Scale factors |
|---|---|
| Astronomy sky surveys | Uncertainty due to long temporal frame; paradigm shifts |
| Deep subseafloor biosphere | Scarce data are sparse data; high variety; difficult to standardize |
| Craniofacial research | High variety in genomes studied, models, methods, duration of analysis; difficult to standardize |
| Computational sciences | High variety in data, methods, tool expertise; difficult to standardize |

# Research Question 3

How does the degree of *centralization of data collection and analysis* influence use, reuse, curation, and project strategy?

| Domain |
| --- |
| Astronomy sky surveys |
| Deep subseafloor biosphere |
| Craniofacial research |
| Computational science |

UCLA Center for Knowledge Infrastructures

# Centralization factors

| Research Site | Centralization factors |
|---|---|
| Astronomy sky surveys | Centralized data collection and initial processing; decentralized use and analysis |
| Deep subseafloor biosphere | Common data source, shared repositories of cores; decentralized analysis |
| Craniofacial research | Decentralized data collection; efforts to integrate data for centralized analysis reveal lack of commonalities |
| Computational sciences | Decentralized data collection; efforts to integrate data for centralized analysis reveal lack of commonalities |

# Conclusions so far

- General
  - Data sharing is not one problem, but many
  - Factors interact: domain mixtures, scale, centrality
- Research themes
  - Domains consist of subdomains with fluid boundaries
  - Volume might be least important scale factor
  - Centrality contradictions
    - Centralized data collections become decentralized in analysis
    - Decentralized data collections are hardest to integrate for analysis

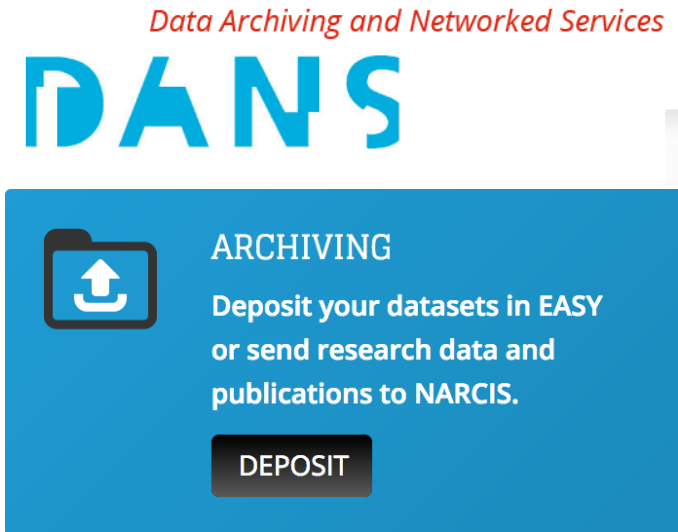UCLA Center for Knowledge Infrastructures

# DANS Users and Uses

# Why do users put data in DANS?

- Meet legal requirements

- Preserve data for long term

- Get credit for data

- Control access to data

- Use as background service

- Motivate citizen science participants

Borgman, C. L., Van de Sompel, H., Scharnhorst, A., van den Berg, H., & Treloar, A. (2015). Who Uses the Digital Data Archive? An Exploratory Study of DANS. In *Proceedings of Association for Information Science and Technology*. http://doi.org/10.1002/pra2.2015.145052010096

# Recommendations so far

- Identify practices of subdomains and interactions

- Seek right level of abstraction for data sharing, integration, curation, reuse

- Invest in data curation early in project design

- Promote infrastructure solutions
  - Shared tools and services
  - Data discovery mechanisms
  - Iterative stewardship

**UCLA** Center for Knowledge Infrastructures

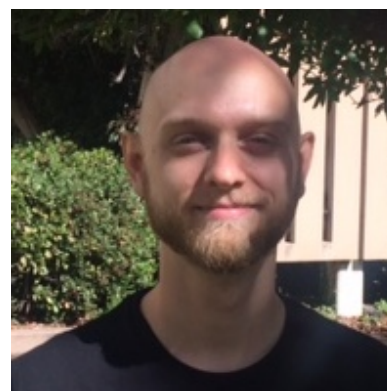# Acknowledgements



Christine Borgman

Peter Darch

Ashley Sands

Irene Pasquetto

Bernie Randles

Milena Golshan

Pietro Santachiara

*Data Archiving and Networked Services*

DANS

UCLA Center for Knowledge Infrastructures