# UCLA
## Publications

**Title**

From Data Creator to Data Reuser: Distance Matters

**Permalink**

**Authors**

Borgman, Christine L.

Groth, Paul T.

**Publication Date**

2024-04-09

**Copyright Information**

Peer reviewed

# From Data Creator to Data Reuser: Distance Matters

Presentation by Christine L. Borgman, Distinguished Research Professor in Information Studies, University of California, Los Angeles, based on joint work with Paul T. Groth, University of Amsterdam

John P. Schneider Honorary Colloquium
iSchool, University of Texas, Austin
April 9, 2024
https://www.ischool.utexas.edu/events/572

Sharing research data is complex, labor-intensive, expensive, and requires infrastructure investments by multiple stakeholders. Open science policies focus on data release, yet reuse is also difficult and may never occur. Investments in data management could be made more wisely by considering who might reuse data, how, why, for what purposes, and when. Drawing upon empirical studies of data sharing and reuse, we develop the theoretical construct of *distance* between data creator and data reuser, identifying six distance dimensions that influence the ability to transfer knowledge effectively: domain, methods, collaboration, curation, purposes, and time and temporality. These dimensions are primarily social in character, with associated technical aspects that can decrease – or increase – distances between creators and reusers. We identify ways that data creators, data reusers, data archivists, and funding agencies can make data sharing and reuse more effective.

Christine L. Borgman conducts research in scientific data practices and information policy. Her publications in information studies, computer science, communication, and law include three award-winning books from MIT Press and more than 250 journal articles, conference papers, and other scholarly products. A Fellow of the American Association for the Advancement of Science and the Association for Computing Machinery, she has held visiting posts at Oxford, Harvard, and several European institutions. Professor Borgman is a member of the Library of Congress Scholars Council and the Board of Directors of the Electronic Privacy Information Center. Her honors and awards include the Paul Evan Peters Award from the Coalition for Networked Information, Association for Research Libraries, and EDUCAUSE; Award of Merit and the Research in Information Science Award, both from the Association for Information Science and Technology; and a Legacy Laureate of the University of Pittsburgh.

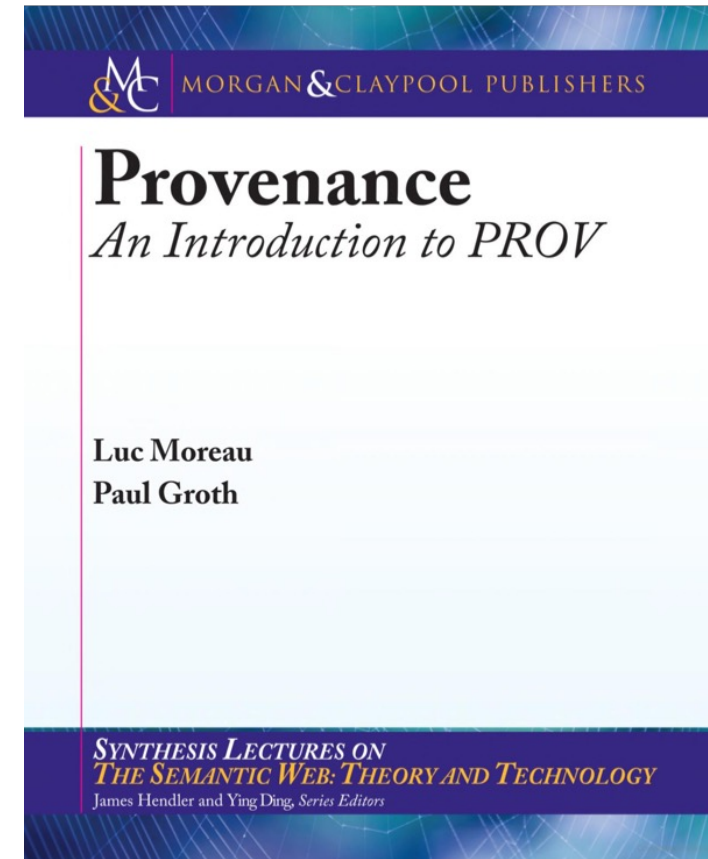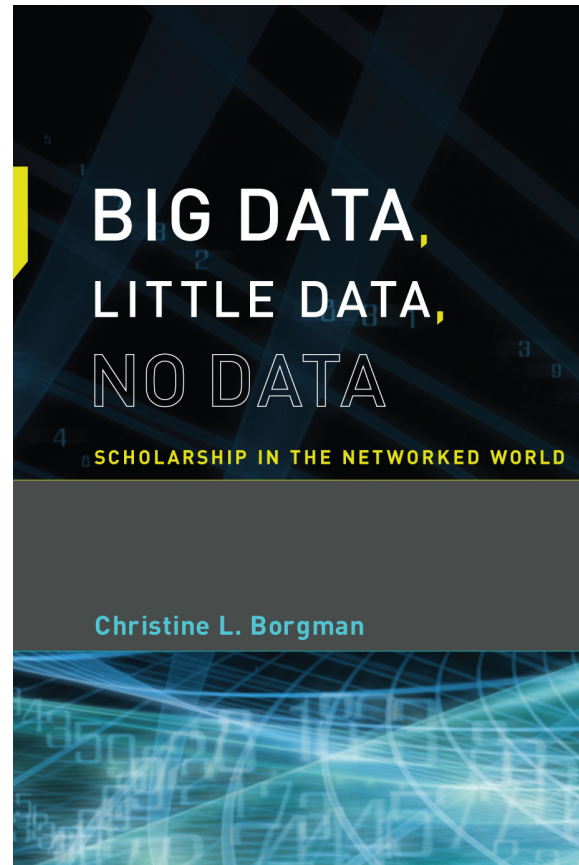# From Data Creator to Data Reuser: Distance Matters

## Christine L. Borgman

Distinguished Research Professor
University of California, Los Angeles

## & Paul T. Groth

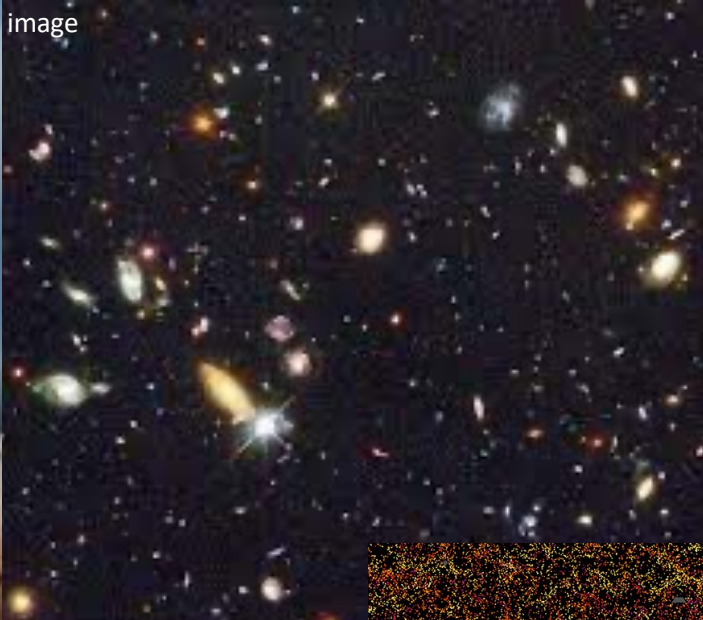Professor of Algorithmic Data Science
University of Amsterdam

Distinguished Schneider Lecture
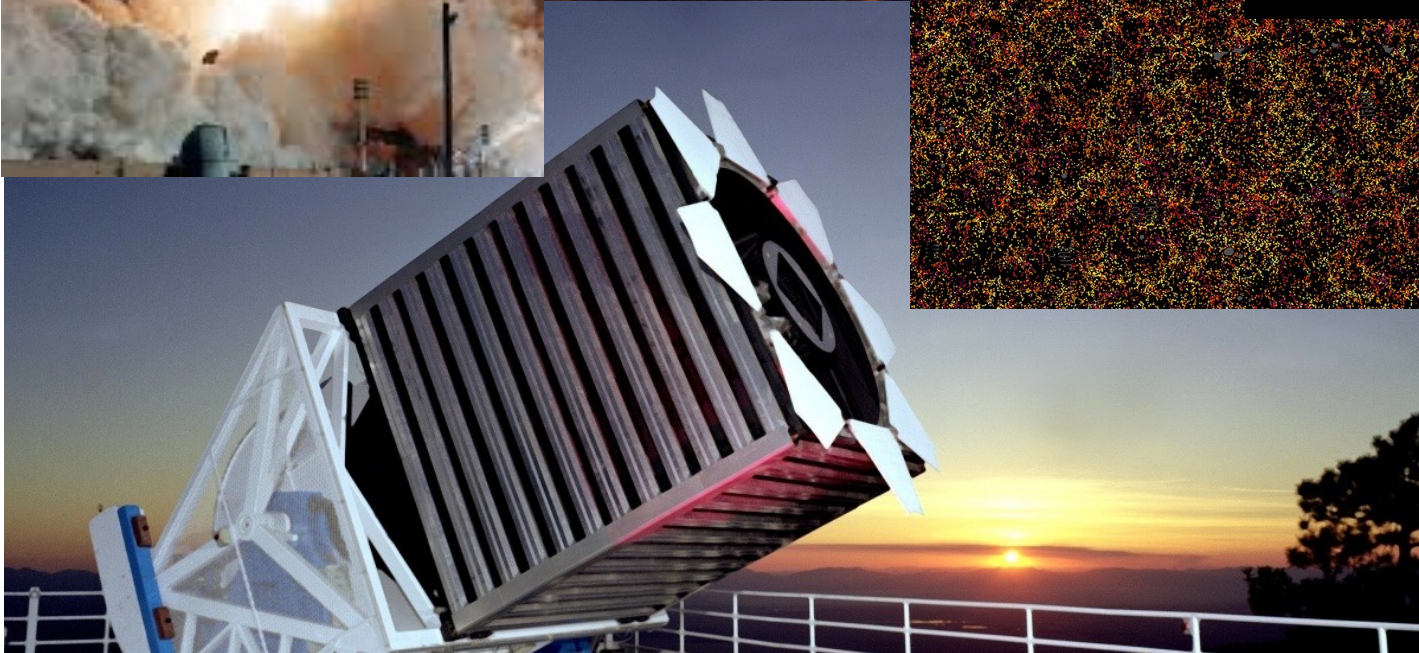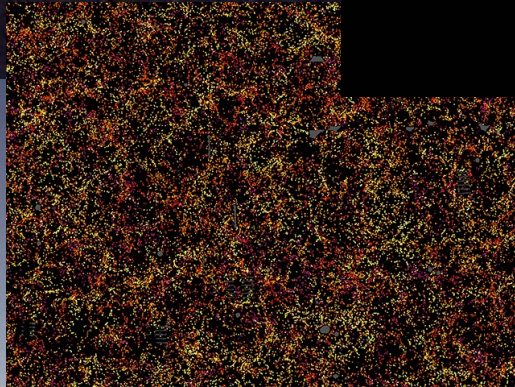University of Texas, Austin
April 9, 2024

# Decades to acquire data, decades to preserve



Hubble space telescope launch; deep field image

Stratospheric Observatory for Infrared Astronomy: Center of Milky Way Galaxy

2

# Research data for the public good

## Why share research data?

- Reuse
- Reproduce
- Transparent
- Educate
- Required
  - Funding agencies
  - Journals

## How to share research data?

- Deposit in data archive
- Publish data documentation
  - Research protocols
  - Codebooks
  - Software
  - Algorithms
- Link datasets to publication
- Cite data and software
- Develop instructional materials

# Data sharing vs. Data reuse



- What are best investments in data sharing?

- What data are most likely to be reused?

- What factors influence data reuse?
  - Social
  - Technical

- How can data sharing be more effective ?
  - Data creators
  - Data reusers

Photo by Mathieu Stern on Unsplash

Borgman, C. L., & Bourne, P. E. (2022). Why It Takes a Village to Manage and Share Data. *Harvard Data Science Review*, *4*(3).
Borgman, C. L., & Brand, A. (2022). Data blind: Universities lag in capturing and exploiting data. *Science*, *378*(6626), 1278–1281.

# How can data creators enable data reuse?

Goodman, A., et al. (2014). Ten Simple Rules for the Care and Feeding of Scientific Data. *PLoS Computational Biology*, *10*(4), e1003542. https://doi.org/10.1371/journal.pcbi.1003542

Wilkinson, M. D., et al. (2016). The FAIR Guiding Principles for scientific data management and stewardship. *Scientific Data*, *3*, 160018. http://dx.doi.org/10.1038/sdata.2016.18

# Dimensions of Distance from Data Creator to Reuser

Social and Technical Distances

1. Domain
2. Methods
3. Collaboration
4. Curation
5. Purposes
6. Time and temporality



**Figure 2**. Digital Curation Center Curation Lifecycle Model (Higgins, 2008). Reprinted with permission of the Digital Curation Centre, U.K.

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

# Domain Distance

# Domain Distance

# Data Creators' Advantage

**Comparative Data Reuse**

- Ground truthing: calibrate, compare, confirm

- Instrument calibration

- Frequent, routine practice

**Integrative Data Reuse**

- Analysis: identify patterns, correlations, causal relationships

- Novel statistical analyses

- Rare, emergent practice

Pasquetto, I. V., Borgman, C. L., & Wofford, M. F. (2019). Uses and reuses of scientific data: The data creators' advantage. *Harvard Data Science Review,* 1:2

Bret Kavanaugh, Unsplash

National Cancer Institute

# Collaboration Distance



Advisor / Student

Invisible Colleges

Community of interest

Formal Projects

Best research friend next door

Random Downloader

Interpersonal exchange

Collegial relationship

Occasional interaction

No known relationship

# Collaboration Distance

# Curation Distance

# Curation Distance

# Data production, knowledge production, and reuse

- Borgman, C. L., & Wofford, M. F. (2021). From Data Processes to Data Products: Knowledge Infrastructures in Astronomy. *Harvard Data Science Review*, *3*(3).
- Baker, K. S., & Mayernik, M. S. (2020). Disentangling knowledge production and data production. *Ecosphere*, *11*(7), e03191.

14

# Discussion and Implications

# Data creator ⟺ Data reuser

- Knowledge transfer
  - Direct, interpersonal
  - Intermediaries: archives
- Shorter distances
  - Easier, fuller transfer
  - Higher likelihood of occurrence
- Longer distances
  - Greater difficulty, expense
  - Lower likelihood of occurrence



Photo by Nadine Shaabana on Unsplash

# Social and Technical Distances

Domain  ⟷ *Lubrication / friction* ⟷  Infrastructure

Methods  ⟷ *Lubrication / friction* ⟷  Standards

Collaboration  ⟷ *Lubrication / friction* ⟷  Data formats

Curation  ⟷ *Lubrication / friction* ⟷  Documentation

Purposes  ⟷ *Lubrication / friction* ⟷  Interfaces

Time and temporality  ⟷ *Lubrication / friction* ⟷  Provenance

# Recommendations to Stakeholders

Providers: Who are your potential reusers?

Reusers: How close am I to the data creator?

Archivists: Know thy users and select carefully

Funders: It takes a village to share data

# Recommendations to Data Creators


Identify the audiences for your data


Focus on reusers early in the research process


Greater distance => greater investment in data

Infrastructure
Standards
Data formats
Documentation
Interfaces
Provenance

# Recommendations to Data Reusers

Identify your locus on each dimension

Assess investments needed to find reusable data

Assess investments needed to make data reusable

Infrastructure
Standards
Data formats
Documentation
Interfaces
Provenance

# Recommendations to Data Archivists

Know where your users reside on each dimension

Greater distance => more investment in curation

Greater distance => more investment user services

Infrastructure
Standards
Data formats
Documentation
Interfaces
Provenance

# Recommendations to Funding Agencies

**Support investigators in targeted data release**

**Support more research on**
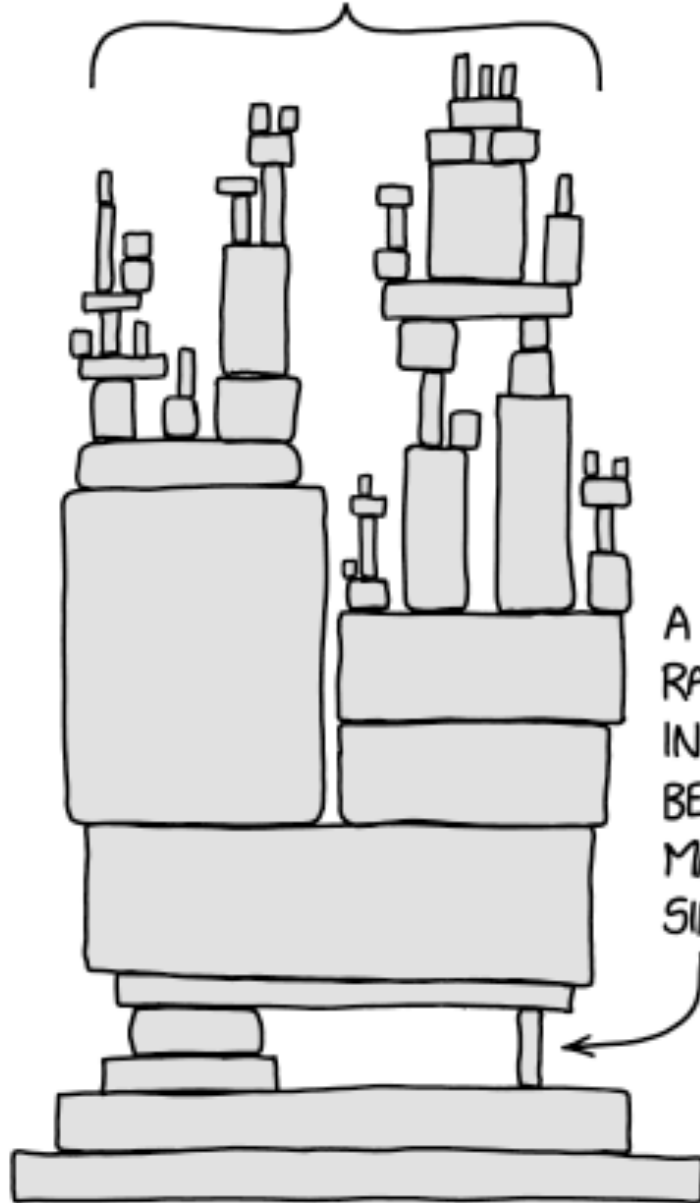
Infrastructure for data sharing and reuse

What data are reused, by whom, when, why, how

How to facilitate data reuse

How to develop 'human infrastructure' for research

How to improve knowledge exchange at distance

https://www.explainxkcd.com/wiki/index.php/2347:_Dependency

# *May all your problems be technical*

Jim Gray, Turing Award Winner