# UCLA
## Presentations

**Title**
Creating, Collaborating, and Celebrating the Diversity of Research Data

**Permalink**
https://escholarship.org/uc/item/7cs811mj

**Author**
Borgman, Christine L.

**Publication Date**
2015-10-26

**Copyright Information**

October 26, 2015

# Creating, Collaborating, and Celebrating the Diversity of Research Data

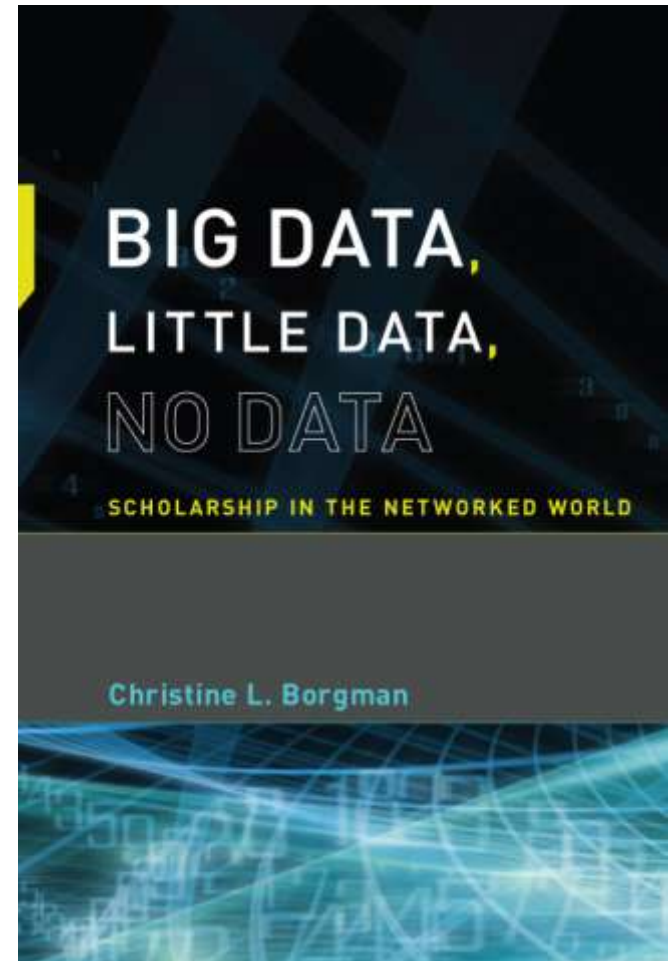Christine L Borgman, *University of California, Los Angeles*

# Creating, Collaborating, and Celebrating the Diversity of Research Data

## Christine L. Borgman

Distinguished Professor
and Presidential Chair in Information Studies
University of California, Los Angeles
@SciTechProf


Seminar Presentation
Graduate School of Library, Information, and Media Studies
University of Tsukuba, Japan
October 26, 2015



BIG DATA, LITTLE DATA, NO DATA

SCHOLARSHIP IN THE NETWORKED WORLD

Christine L. Borgman

# PHILOSOPHICAL

## TRANSACTIONS:

### GIVING SOME

# ACCOMPT

#### OF THE PRESENT

Undertakings , Studies , and Labours

#### OF THE

# INGENIOUS

#### IN MANY

#### CONSIDERABLE PARTS

#### OF THE

# WORLD

*Vol I.*

For *Anno* 1665, and 1666.

# Data <–> Publications

Publications are arguments made by authors, and data are the evidence used to support the arguments.

PHILOSOPHICAL
TRANSACTIONS:
GIVING SOME
ACCOMPT
OF THE PRESENT
Undertakings , Studies , and Labours
OF THE
INGENIOUS
IN MANY
CONSIDERABLE PARTS
OF THE
WORLD

Vol I.
For Anno 1665, and 1666.

In the SAVOY,
Printed by T. N. for John Martyn at the Bell, a little with-
out Temple-Bar , and James Allestry in Duck-Lane,'
Printers to the Royal Society.

PHILOSOPHICAL
TRANSACTIONS B

Celebrating 350 years of Philosophical Transactions life sciences papers

**Theme issue 'Celebrating 350 years of Philosophical Transactions: life sciences papers' compiled and edited by Linda Partridge**
19 April 2015; volume 370, issue 1666

# Open access policies

- Australian Research Council
  - Code for the Responsible Conduct of Research
  - Data management plans
- National Science Foundation
  - Data sharing requirements
  - Data management plans
- U.S. Federal policy
  - Open access to publications
  - Open access to data
- European Union
  - European Open Data Challenge
  - OpenAIRE
- Research Councils of the UK
  - Open access publishing
  - Provisions for access to data

Australian Government
National Health and Medical Research Council

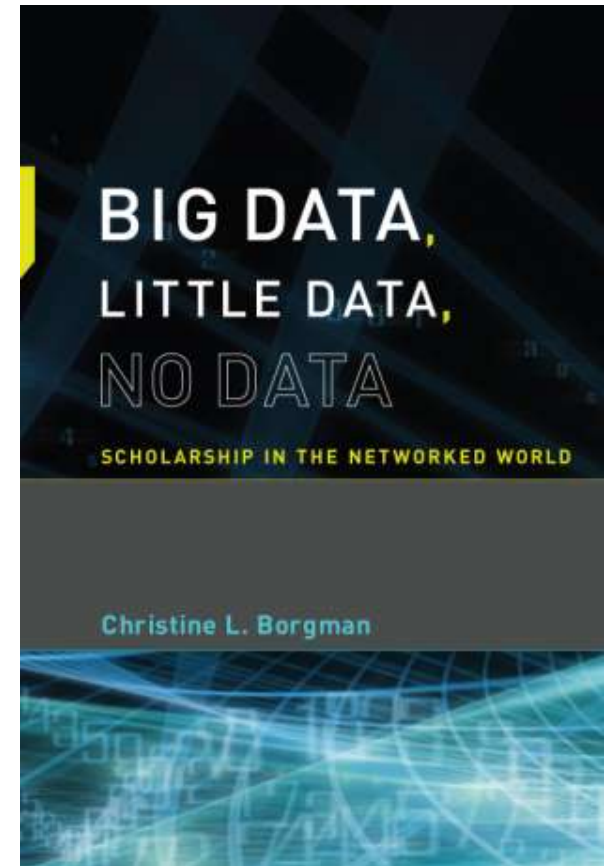National Science Foundation
WHERE DISCOVERIES BEGIN

Supported by
wellcometrust

Policy RECommendations for Open Access to Research Data in Europe

RECODE

5

# Big Data, Little Data, No Data: Scholarship in the Networked World

- Part I: Data and Scholarship
  - Ch 1: Provocations
  - Ch 2: What Are Data?
  - Ch 3: Data Scholarship
  - Ch 4: Data Diversity
- Part II: Case Studies in Data Scholarship
  - Ch 5: Data Scholarship in the Sciences
  - Ch 6: Data Scholarship in the Social Sciences
  - Ch 7: Data Scholarship in the Humanities
- Part III: Data Policy and Practice
  - Ch 8: Releasing, Sharing, and Reusing Data
  - Ch 9: Credit, Attribution, and Discovery
  - Ch 10: What to Keep and Why
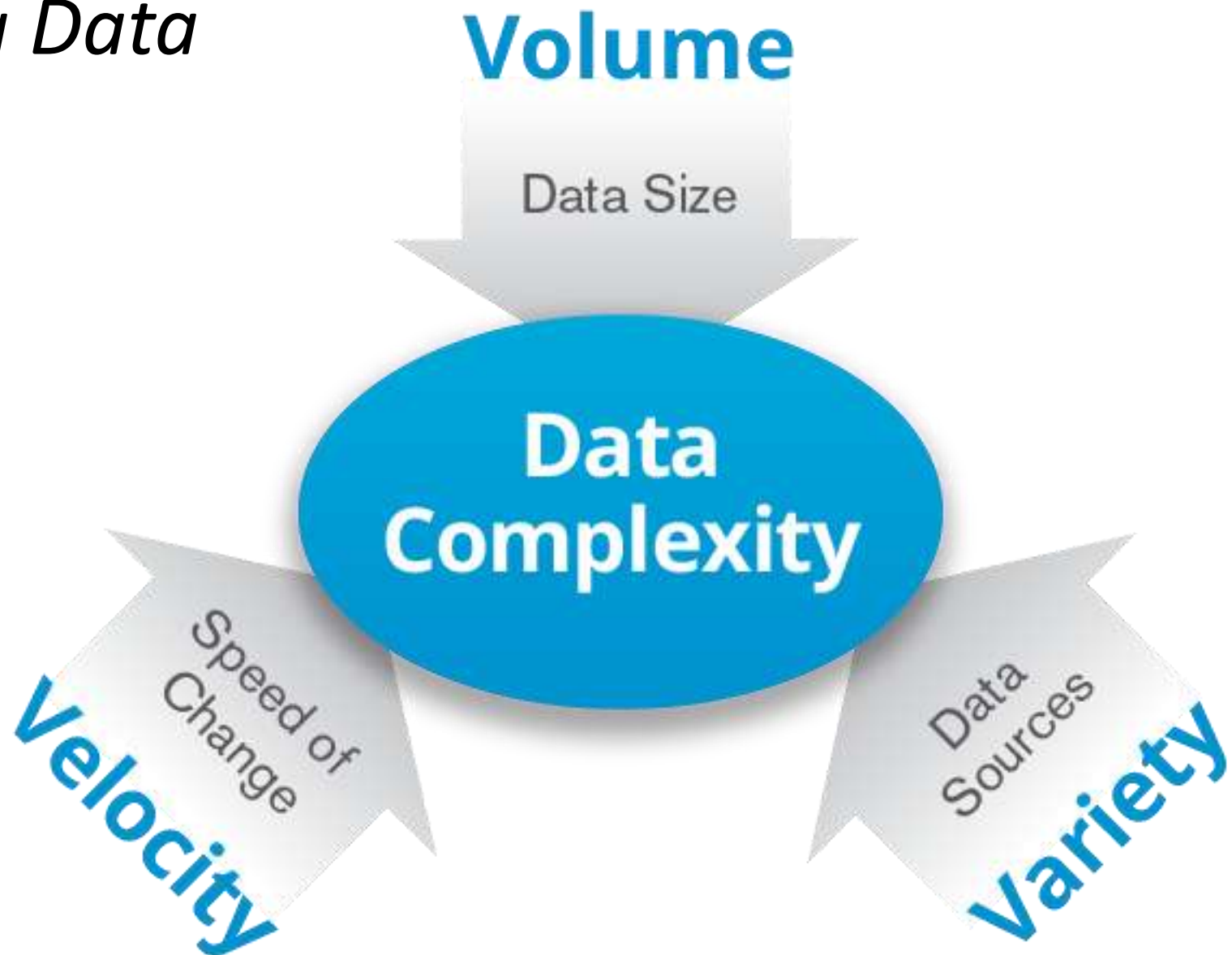

BIG DATA, LITTLE DATA, NO DATA

SCHOLARSHIP IN THE NETWORKED WORLD

Christine L. Borgman

# Celebrating the diversity of data

- Defining data

- Creating data

- Collaborating with data

- Consolidating data value

7

Data

# *Big Data*
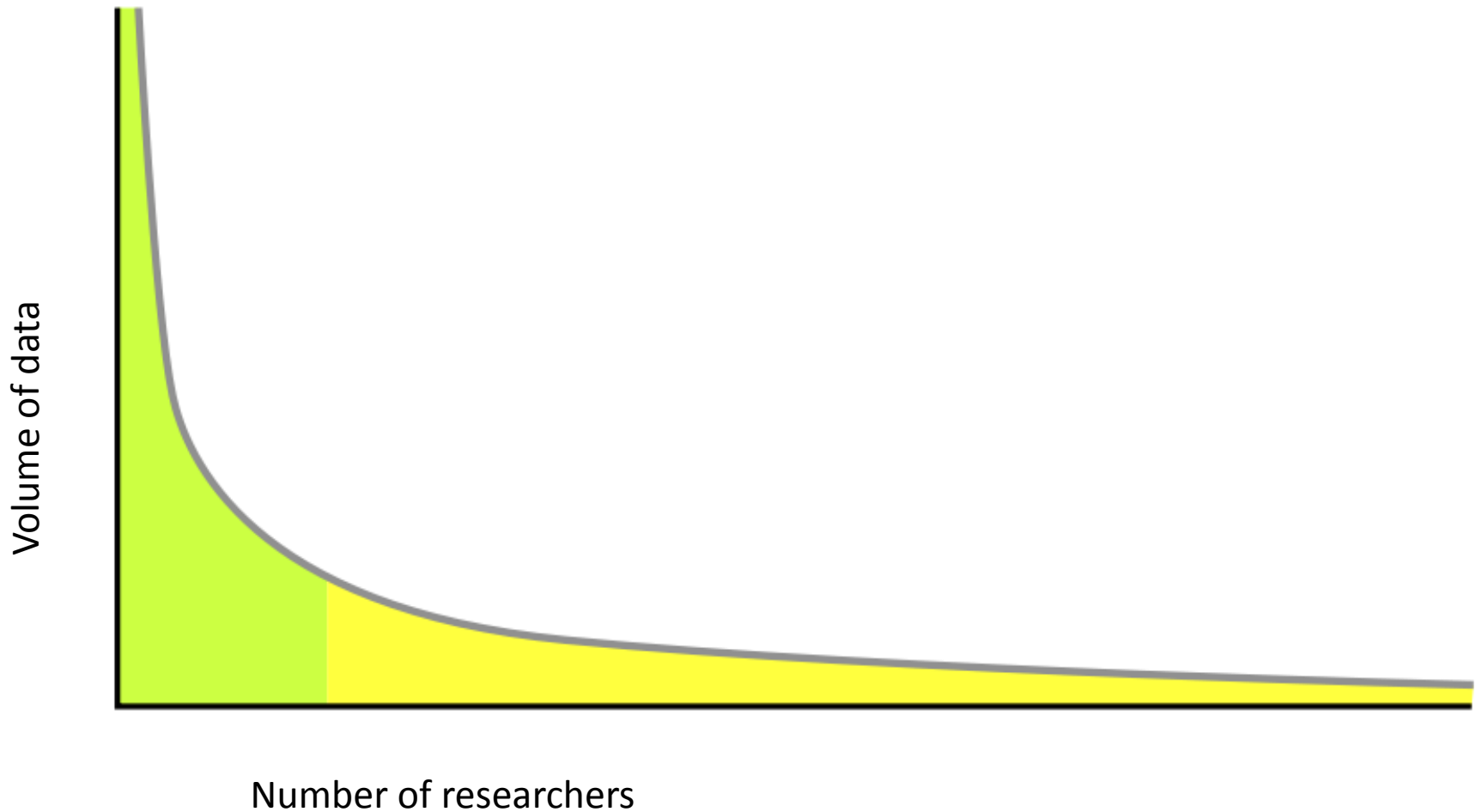
# Long tail of data



Volume of data

Number of researchers

Slide: The Institute for Empowering Long Tail Research

# Open Data: Free

- A piece of data or content is open if anyone is free to use, reuse, and redistribute it — subject only, at most, to the requirement to attribute and/or share-alike
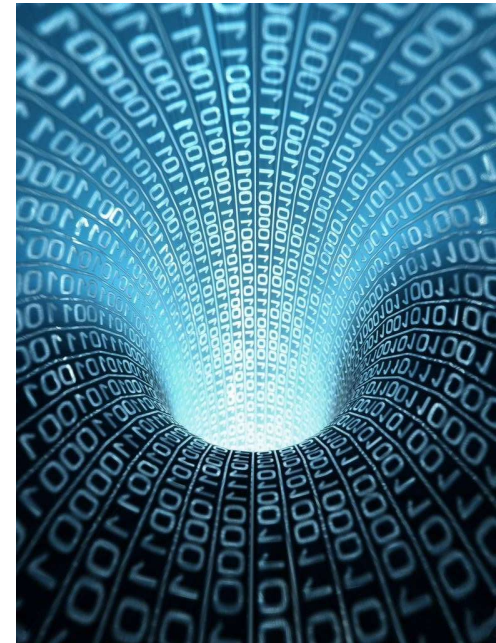
Open Data Commons. (2013).



State Library and Archives of Florida, 1922. Flickr commons photo

```
/FontMatrix matrix def
/FontBBox[2048 -1164 1 index div -628 2 index div 4096 3 index div 2062 5 -1 rol
/sfnts [<
74727565000090000000000000063767420000000000000009C000007DA6670676D00000000000000878
019901AC01C101C501C901E101F601F601F60222022202280236023F024302460267028502850294
01160125011800EA00EA00AE0000003E05BB008A04D70053003FFF8CFFD500150028002200990062
B200402F2B59B002602D2C21B0C051580C6423648BB81555621BB200802F2B59B002602D2C0C6423
B8FF80B30809341CB8FF8040E80809343609352F4A0959045809A7090626062B082C0B280C281328
D4401E091202550C0C047F180118401112025518401D1D02551810251E073C2C04002FED3FEDC42B
1700030209131415002159020718250901E0505262500180C0D0D06551802101006551800CB8FFF8
1F1E4A62A1A4FB432E7900020044FFE40405058E001F002D024BB1020243545840262F4010100255
2F2B2BDD2B2BC01112392F2BCDD0CD003FED5DC45D5D2B323FED12392F5D5D5DCD31301B4019125F
11242517012B50100110302A2912110608070A101B011BB80152B38F2D012DB802F6B2012A2BB801
025D5972103C103C1112173911123939111239011112393912393931304379407A4B573A4524351C
2D602D702D04802D902D02B02D01002D102DC02DD02D042D6037A67F182B10F65D5D71723CFD3C10
FFEAB40C0C02550CB8FFE2B40D0D02550CB8FFD6B4101002550CB8FFDEB50F0F02550C05BD03E200
2F2F2F3FDDCD3F3F10ED10ED313001B00D4B5458BF0046FFE80045FFE8002FFFE80030FFE8B51A18
1010065529B8FFF2B70F0F0655292935341112392F2B2B2BDD2B2B2BC01112392F2B2B2BCD2B2B2B
3236353427260200D07E6B76CF7FCF7A677DCC53356B429F82617E694703AF9E87AF7BFC80A58BAD
961A9C1E9621982A9E2BA816A61AAB1CAD2BB916BE2BCD2BDA2BEC2BFB2B2A202D732573288F1397
201A01601A701A021A120B003FC45D5DED5D5D2F3FC45DED5D5D5D1217393F012F2B2BCD2F2BCDD4
363702902126775C4656201F5F92CBBD75546C2115170D21211C9E62455761FEDE2D2D9B7B364D33
393FDD5DCD31301B40350127600D5D362027602770277B02704340B371F3A20481F4820051A08134F
FFE0400A1339082013391B201339012B2B2B2B002B012B2B595959132115232206151417133363
5C0B215391AB8FFF040131539360815392830143929301439110816390 9B8FFE0401B163929401139
000000090001000000000000000100000721FE4500571000FB74FADF100000010000000000000000
0034036300440363000403638B2242F1FBA034E006D0800400E1F7F027F037F047F05043044011 2BF
014A001F000D0126400B1F0DC61F0D571F0D371F0D410D019E0141000D00420141000D001E014100
021E00244552585B90024021E4459594BB8020153205C58B9010F00224544B1222245445958B90C00
2B2B2B2B2B2B2B2B2B2B2B2B2B0001737500737300456944007373017374 2B2B2B2B2B732B00732B
00>] def
/CharStrings 27 dict dup begin
/.notdef 0 def
            /space 1 def
                    /comma 2 def
                            /period 3 def
                                    /W 4 def
                                            /a 5 def
                                                    /b 6 def
                                                            /c 7
```

# Open Data: Useful

- Openness, flexibility, transparency, legal conformity, protection of intellectual property, formal responsibility, professionalism, interoperability, quality, security, efficiency, accountability, and sustainability.
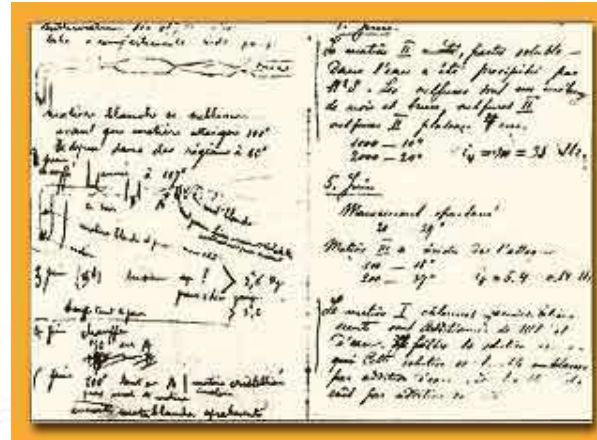


Organization for Economic Cooperation and Development. (2007).
*OECD Principles and Guidelines for Access to Research Data from Public Funding.*
http://www.oecd.org/dataoecd/9/61/38500813.pdf

# What are data?

Marie Curie's notebook aip.org

Pisa Griffin

hudsonalpha.org

Figure 2. Numeric Change in Resident Population for the 50 States, the District of Columbia, and Puerto Rico: 1990 to 2000

http://www.census.gov/population/cen2000/map02.gif

Monthly Mean: f17_ssmis_201207v7.nc

ncl.ucar.edu

Date:1/2.07.75    Place:Sakaltutan
Zafor
He will grow old in his present house; new house is for sons - 5 sons. Not sure they want to live in village. He will only build another if they want him to. eS came from Germany and did the plastering. He arranged the carpentry in Kayseri. Çok para gitti. {much money went} Has a tractor.

Date:July1980    Place:Sakaltutan
Zafor:
Household now Zafor and wife; Nazif Unal and wife and youngest son, still a boy. They run two dolmuß; one with a driver from Süleymanli. Goes in and out once a day. He gets 8,000 a month. Zafor then said, keskin de©il. { not sharp - i.e.? not profitable} I said he did very well on 8,000 TL with only two journeys a day. Nazif Unal has "bought" a Durak {dolmuß stop} from Belediye and works all day in Kayseri.

http://onlineqda.hud.ac.uk/Intro_QDA/Examples_of_Qualitative_Data.ph

hudsonalpha.org

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

# Creating research data



Sloan Digital Sky Survey Telescope,
Apache Point, New Mexico



Sensor networks

WIRED MAGAZINE: 16.07

SCIENCE : DISCOVERIES

# The End of Theory: The Data Deluge Makes the Scientific Method Obsolete

By Chris Anderson ✉   06.23.08

Illustration: Marian Bantjes

The
FOURTH
PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

Hey, Tansley & Tolle (eds.) (2009)

# Tools for Astronomical Big Data
## Tucson, Arizona, March 9-11, 2015

**Links:**
Home
Program
Participants (129)

**Scientific Organizing Committee:**
Eric Feigelson (Penn State)
David Hogg (NYU)
John Kececioglu (Arizona)
Tod R. Lauer (NOAO, Chair)
Dara Norman (NOAO)
Chris Smith (NOAO)

**Local Organizing Committee:**
Tod R. Lauer (NOAO)
Shelley Weintraub (NOAO)

**Current Weather for**
**Tucson, AZ**

**62.0° F**
Feels like: 62° F
Fair

Humidity: 70%
Wind: Southeast at 6.9 mph
19 March, 2015

SHARE

## Program

Invited speakers in **bold**.

| Monday, March 9, 2015 | |
|---|---|
| 8:00-9:00 | Registration/Continental Breakfast |
| 9:00-9:15 | Introductory Remarks |
| 9:15-9:45 | **Alyssa Goodman (Harvard)** *Wide Data vs. Big Data* |
| 9:45-10:15 | **Carlos Scheidegger (University of Arizona)** *How do you look at a billion data points? Exploratory Visualization for Big Data* |
| 10:15-11:00 | Break |
| 11:00-11:20 | Joshua Peek (STScI) *Machine Vision Methods for the Diffuse Universe* |
| 11:50-2:00 | Lunch |
| 2:00-2:20 | Elisabeth Mills (NRAO) *Visualization and Analysis of Rich Spectral-Line Datasets* |
| 2:20-2:40 | Brian Bue (JPL) *Leveraging Annotated Archival Data with Domain Adaptation to Improve Data Triage in Optical Astronomy* |

19

# Research process

- Models and theories
- Research questions
- Methods
  - Practices
  - Data sources
  - Software
  - Instruments
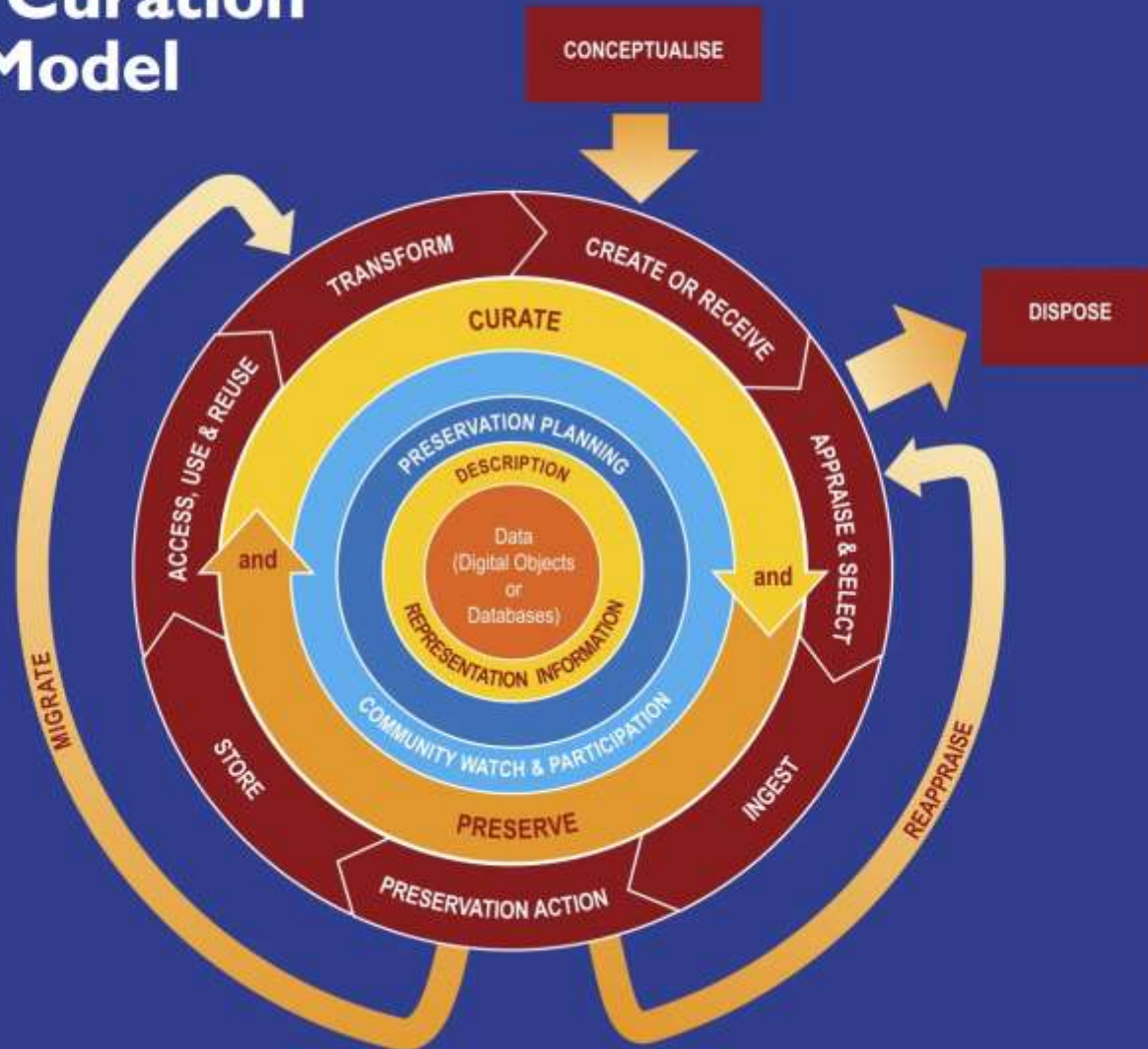  - Infrastructure
  - Domain expertise





http://blogs.agu.org/wildwildscience/files/2009/12/Screen-shot-2009-12-13-at-04.10.146.jpg

http://www.livescience.com/20767-dinosaur-weight-estimates.html

Pepe, A., Mayernik, M. S., Borgman, C. L. & Van de Sompel, H. (2010). From Artifacts to Aggregations: Modeling Scientific Life Cycles on the Semantic Web. Journal of the American Society for Information Science and Technology, 61(3): 567–582.

The DCC Curation Lifecycle Model

# Random walk

# Collaborating with data

24

# Big Science <–> Little Science

- Large instruments
- High cost
- Long duration
- Many collaborators
- Distributed work
- Domain expertise

- Small instruments
- Low cost
- Short duration
- Small teams
- Local work
- Domain expertise



Sloan Digital Sky Survey



Sensor networks for science

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

# LETTERS

# A role for self-gravity at multiple length scales in the process of star formation

Alyssa A. Goodman[1,2], Erik W. Rosolowsky[2,5], Michelle A. Borkin[1]†, Jonathan B. Foster[2], Michael Halle[1,4], Jens Kauffmann[1,2] & Jaime E. Pineda[3]

Self-gravity plays a decisive role in the final stages of star formation, where dense cores (size ~0.1 parsecs) inside molecular clouds collapse to form star-plus-disk systems[1]. But self-gravity's role at earlier times (and on larger length scales, such as ~1 parsec) is unclear; some molecular cloud simulations that do not include self-gravity suggest that 'turbulent fragmentation' alone is sufficient to create a mass distribution of dense cores that resembles, and sets, the stellar initial mass function[2]. Here we report a 'dendrogram' (hierarchical tree-diagram) analysis that reveals that self-gravity plays a significant role over the full range of possible scales traced by $^{13}$CO observations in the L1448 molecular cloud, but not everywhere in the observed region. In particular, more than 90 per cent of the compact 'pre-stellar cores' traced by peaks of dust emission[3] are projected on the sky within one of the dendrogram's self-gravitating 'leaves'. As these peaks mark the locations of already-forming stars, or of those probably about to form, a self-gravitating cocoon seems a critical condition for their existence. Turbulent fragmentation simulations without self-gravity—even of unmagnetized isothermal material—can yield mass and velocity power spectra very similar to what is observed in clouds like L1448. But a dendrogram of such a simulation[4] shows that nearly all the gas in it (much more than in the observations) appears to be self-gravitating. A potentially significant role for gravity in 'non-self-gravitating' simulations suggests inconsistency in simulation assumptions and output, and that it is necessary to include self-gravity in any realistic simulation of the star-formation process on subparsec scales.
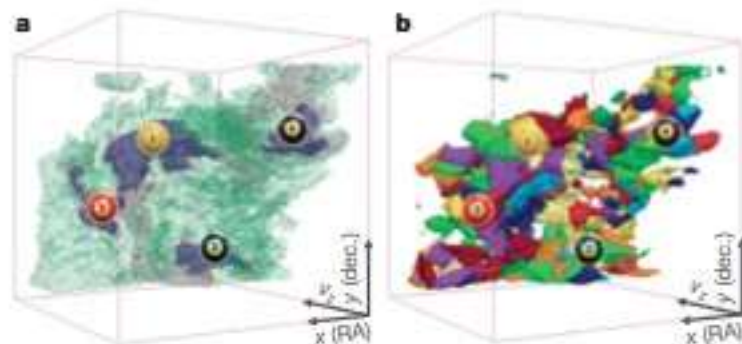
Spectral-line mapping shows whole molecular clouds (typically tens to hundreds of parsecs across, and surrounded by atomic gas) to be marginally self-gravitating[5]. When attempts are made to further break down clouds into pieces using 'segmentation' routines, some self-gravitating structures are always found on whatever scale is sampled[6,7]. But no observational study to date has successfully used one spectral-line data cube to study how the role of self-gravity varies as a function of scale and conditions, within an individual region.

Most past structure identification in molecular clouds has been explicitly non-hierarchical, which makes difficult the quantification of physical conditions on multiple scales using a single data set. Consider, for example, the often-used algorithm CLUMPFIND[8]. In three-dimensional (3D) spectral-line data cubes, CLUMPFIND operates as a watershed segmentation algorithm, identifying local maxima in the position–position–velocity (p–p–v) cube and assigning nearby emission to each local maximum. Figure 1 gives a two-dimensional (2D) view of L1448, our sample star-forming region, and Fig. 2 includes a CLUMPFIND decomposition of it based on $^{13}$CO observations. As with any algorithm that does not offer hierarchically nested or
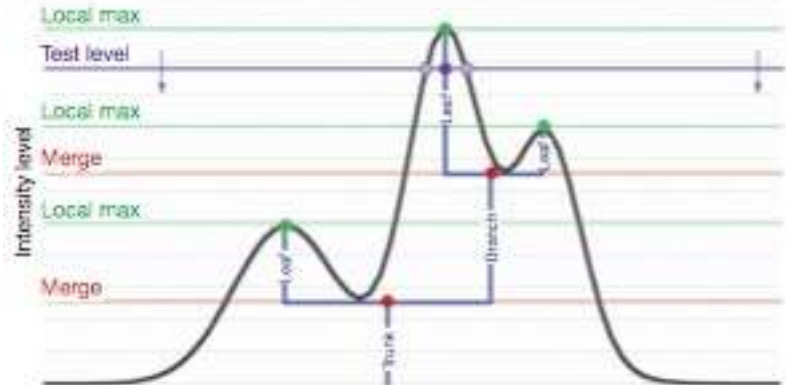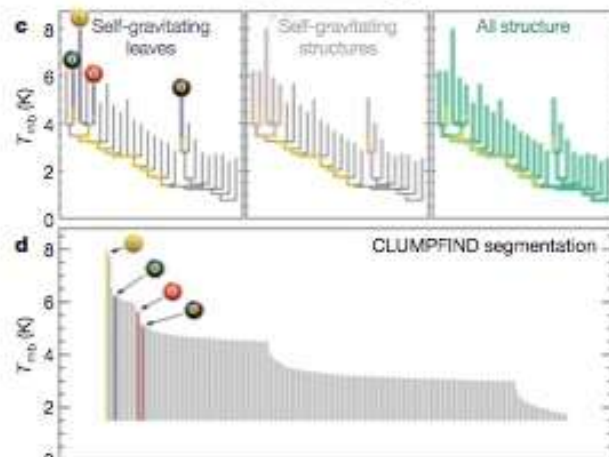
overlapping features as an option, significant emission found between prominent clumps is typically either appended to the nearest clump or turned into a small, usually 'pathological', feature needed to encompass all the emission being modelled. When applied to molecular-line



**Figure 1 | Near-infrared image of the L1448 star-forming region with contours of molecular emission overlaid.** The channels of the colour image correspond to the near-infrared bands $J$ (blue), $H$ (green) and $K$ (red), and the contours of integrated intensity are from $^{13}$CO(1–0) emission[9]. Integrated intensity is monotonically, but not quite linearly (see Supplementary Information), related to column density[10], and it gives a view of 'all' of the molecular gas along lines of sight, regardless of distance or velocity. The region within the yellow box immediately surrounding the protostars has been imaged more deeply in the near-infrared (using Calar Alto) than the remainder of the box (2MASS data only), revealing protostars as well as the scattered starlight known as 'Cloudshine'[11] and outflows (which appear orange in this colour scheme). The four billiard-ball labels indicate regions containing self-gravitating dense gas, as identified by the dendrogram analysis, and the leaves they identify are best shown in Fig. 2a. Asterisks show the locations of the four most prominent embedded young stars or compact stellar systems in the region (see Supplementary Table 1), and yellow circles show the millimetre-dust emission peaks identified as star-forming or 'pre-stellar' cores[3].
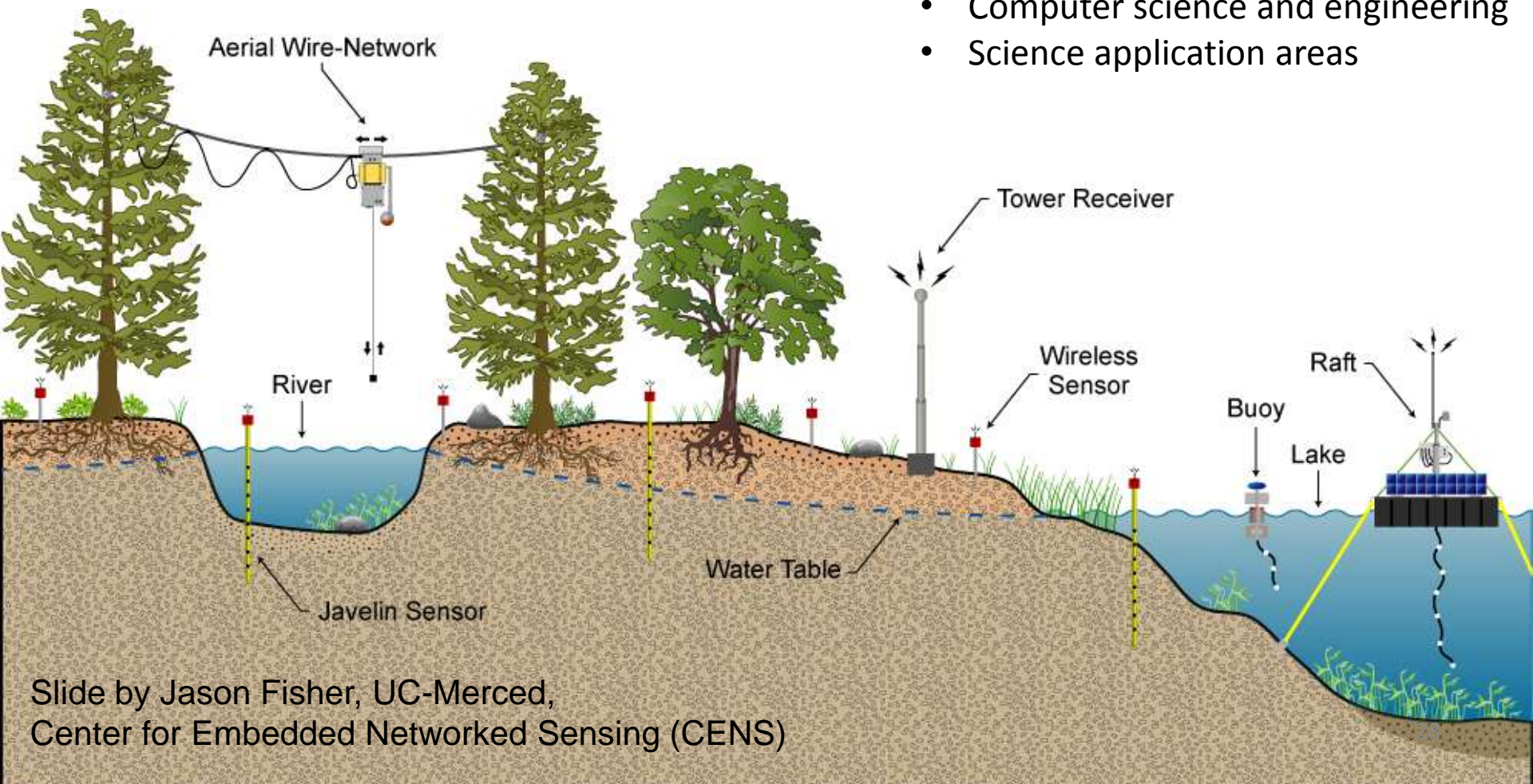
Initiative in Innovative Computing at Harvard, Cambridge, Massachusetts 02138, USA. [2]Harvard-Smithsonian Center for Astrophysics, Cambridge, Massachusetts 02138, USA. [3]Department of Physics, University of British Columbia, Okanagan, Kelowna, British Columbia V1V 1V7, Canada. [4]Surgical Planning Laboratory and Department of Radiology, Brigham and Women's Hospital, Harvard Medical School, Boston, Massachusetts 02115, USA. †Present address: School of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts 02138, USA.

63





**Figure 3 | Schematic illustration of the dendrogram process.** Shown is the

# 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 <–> 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
Iodp.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



30

# Self-descriptions C-DEBI scientists

Aquatic chemistry

Aquatic microbial ecology

Astrobiology

Biochemistry

Bioenergetics

Biogeochemistry

Biogeography

Bioinformatics

Biology

Chemical ecology

Chemical oceanography

Deep sea biogeochemistry

Deep sea microbiology

Ecology

Ecophysiology

Environmental chemistry

Environmental microbiology

Genomics

Geobiology

Geochemistry

Geomicrobiology

Geophysics

Hydrogeology

Hydrology

Hydrothermal microbiology

Inorganic chemistry

Marine ecology

Marine geology

Marine microbial biogeochemistry

Marine microbial ecology

Marine microbiology

Metabolomics

Metagenomics

Microbial biogeochemistry

Microbial biogeography

Microbial ecology

Microbial oceanography

Microbial physiology

Microbiology

Mineralogy

Molecular biogeochemistry

Molecular biology

Molecular microbial ecology

Molecular microbiology

Molecular physiology

Paleoceanography

Paleoclimatology

Paleogeomicrobiology

Petrology

Physiology

Plant biochemistry

Sedimentary biogeochemistry

Sedimentary geochemistry

Sedimentology

Peter Darch, UCLA

Left column (prefixes):
Aqua-
Astro-
Bio-
Chem-
Climat-
Deep sea-
Eco-
Energe-
Environmenta
l-
Gen-
Geo-
Geograph-
Hydro-
-ics
Informat-
Inorganic-
Marine-
Meta-
Metabol-
Micro-
Mineral-
Molecular-
Ocean-
-omics
Paleo-
Petro-
Physic-
Physiol-
Plant-
Sediment-
Therm-

Right column (fields):
Aquatic chemistry
Aquatic microbial ecology
Astrobiology
Biochemistry
Bioenergetics
Biogeochemistry
Biogeography
Bioinformatics
Biology
Chemical ecology
Chemical oceanography
Deep sea biogeochemistry
Deep sea microbiology
Ecology
Ecophysiology
Environmental chemistry
Environmental microbiology
Genomics
Geobiology
Geochemistry
Geomicrobiology
Geophysics
Hydrogeology
Hydrology
Hydrothermal microbiology
Inorganic chemistry
Marine ecology
Marine geology
Marine microbial biogeochemistry
Marine microbial ecology
Marine microbiology
Metabolomics
Metagenomics
Microbial biogeochemistry
Microbial biogeography
Microbial ecology
Microbial oceanography
Microbial physiology
Microbiology
Mineralogy
Molecular biogeochemistry
Molecular biology
Molecular microbial ecology
Molecular microbiology
Molecular physiology
Paleoceanography
Paleoclimatology
Paleogeomicrobiology
Petrology
Physiology
Plant biochemistry
Sedimentary biogeochemistry
Sedimentary geochemistry
Sedimentology

Peter Darch, UCLA

Browse Data  >  Conservatives report, but liberals display, greater happiness  >  Wojcik et al - Behavioral Happiness - Study 1 data

# Conservatives report, but liberals display, greater happiness: Wojcik et al - Behavioral Happiness - Study 1 data

Principal Investigator(s) :  Wojcik, Sean; Hovasapian, Arpine; Graham, Jesse; Motyl, Matt; Ditto, Peter;

[f Share] [ 0 ]

**wojcik-et-al-behavioral-happiness-study-1-data-26097.dta** (application/x-stata) 55839

Extended Properties

⊙ **Download this dataset**

Browse variables:  sesladdercountry (Subjective SES) ▾

## sesladdercountry:Subjective SES

| Category | N=1433 |
|---|---|
| **Valid Cases** | **Valid N=1431** |
| 0 : Lowest | 408 |
| 1 | ←3 |
| 2 | ←24 |
| 3 | ←56 |
| 4 | ←122 |
| 5 | ←112 |
| 6 | ←180 |
| 7 | 219 |
| 8 | ←214 |
| 9 | ←70 |
| 10 : Highest | ←23 |
| **Invalid Cases** | **Invalid N=2** |
| System Missing | ←2 |

33

# The Pisa Griffin Project

The aim of this project is to perform a comparative study of three artworks (bronze casts of Islamic provenance), to discover evidence of similarities and to get new insight on their origin.

Probably produced within the Islamic Mediterranean in the eleventh century, the Griffin has incised on its body a long inscription in Arabic expressing good wishes. Captured by the Pisans, it underwent an extraordinary transformation: for centuries it was a terrifying, sound-producing guardian figure on top of the roof of Pisa Cathedral. The present project is focused on the Griffin but also includes alongside it other bronze animal sculptures such as a Lion and a Falcon. It is hoped that the interdisciplinary study of the Griffin will shed light on the significance of such objects in a global Mediterranean culture.

## Videos

The Pisa Griffin: an introduction

http://vcg.isti.cnr.it/griffin/

Arte islamica, ippogrifo, XI sec 03, own work

# 6 BEASTS THAT ROARED: THE PISA GRIFFIN AND THE NEW YORK LION

*Anna Contadini, Richard Camber and Peter Northover*

## The Pisa Griffin
*Anna Contadini*

My interest in the Pisa Griffin (pl. 6.1) goes back to my childhood, when my parents took me to visit Pisa for the first time. I am still as impressed by the beast as I was then, but I am now equally intrigued by the mystery that [...]

[...] further substantiation. [...] to be learnt about its [...] one might call a gyn[...] through the opening [...] discovered that it had [...] Griffin and attached [...] with slightly everted [...] the animal (pl. 6.3). [...]

## The Composition of the Lion and the Griffin
*Peter Northover*

The discussion of the compositions of the two sculptures will be made in the order in which the analyses were carried out, that is with the Lion first, followed by the Griffin. All the analyses have been made by electron probe microanalysis (EPMA) using wavelength dispersive spectrometry; this method has been well standardised against other current techniques so the results will be broadly comparable with those from other laboratories.

[permission to take samples for analysis, Mr Milliam Chiarini, or the Archivio dell'Opera del Duomo, who assisted in the study of the documents and facilitated the photography; Mons. Mario Baroncini of the Archivio Capitolare; and the staff of the Pisa ...]

[problems... with ... though Ralph has ... he has followed out ... discussing various ...]

### Analysis of the Lion and Griffin

| sample | Object | Part | Fe | Co | Ni | Cu | Zn | As | Sb | Sn | Ag | Bi | Pb | Au | S |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 24853.1 | Lion | belly | | | | | | | | | | | | | |
| 24853.2 | | | | | | | | | | | | | | | |
| 24853.3 | | | | | | | | | | | | | | | |
| 24892.1 | Lion | right fore leg | | | | | | | | | | | | | |
| 24892.2 | | | | | | | | | | | | | | | |
| 24892.3 | | | | | | | | | | | | | | | |
| 24873.1 | Lion | seamus vessel | | | | | | | | | | | | | |
| 24873.2 | | | | | | | | | | | | | | | |
| 24873.3 | | | | | | | | | | | | | | | |
| 24903.1 | Lion | right hind leg | | | | | | | | | | | | | |
| 24903.2 | | | | | | | | | | | | | | | |
| 24903.3 | | | | | | | | | | | | | | | |
| 24903.1 | Lion | right fore leg | | | | | | | | | | | | | |
| ... | Griffin | left wing, drill hole | | | | | | | | | | | | | |
| ... | Griffin | right wing, drill hole | | | | | | | | | | | | | |
| ... | Griffin | body, brazing | | | | | | | | | | | | | |
| ... | Griffin | vessel, upper | | | | | | | | | | | | | |
| ... | Griffin | vessel, lower | | | | | | | | | | | | | |
| ... | Head on mud | weld metal | | | | | | | | | | | | | |
| ... | Head on mud | shaft | | | | | | | | | | | | | |
| ... | Head on mud | left hind hoof | | | | | | | | | | | | | |

Precondition:

# Researchers share data

# Lack of incentives to share data



- Labor to document data

- Benefits to unknown others

- Competition

- Control

- Confidentiality...

Image source: www.buildingsrus.co.uk/.../ target1.htm

# Lack of incentives to reuse data

- Identify useful data
  - Documentation
  - Interpretation
  - Software
- Cleaning
- Trust
- Credit
- Licensing...

http://fyi.uiowa.edu/wp-content/uploads/2011/10/utopia_in_four_movements_filmstill5_utopiasign.jpg

# Consolidating value in data



July 19, 1922. State Library and Archives of Florida.
Flickr commons



Page 105 of "The Street railway journal" (1884);
Flickr Commons

Australian National Data Service

# Ways to pool data

- Centralized data production
  - Top down investments in data
  - Pooled data resources for the community

- Decentralized data production
  - Bottom up investments in data
  - Local data resources pooled later

# Sloan Digital Sky Survey

Social Science Surveys

# Discovery and Interpretation

- Identify the form and content
- Identify related objects
- Interpret
- Evaluate
- Open
- Read
- Compute upon
- Reuse
- Combine
- Describe
- Annotate…

DATA

METADATA

# Describing and attributing data



- Compound objects
  - Observations
  - Software
  - Protocols…
- Attribution
  - Investigators
  - Data collectors
  - Analysts…
- Ownership, responsibility



Mary Jane Rathbun (1860-1943), working with crab specimens

# Metadata

- Metadata is structured information that describes, explains, locates, or otherwise makes it easier to retrieve, use, or manage an information resource.
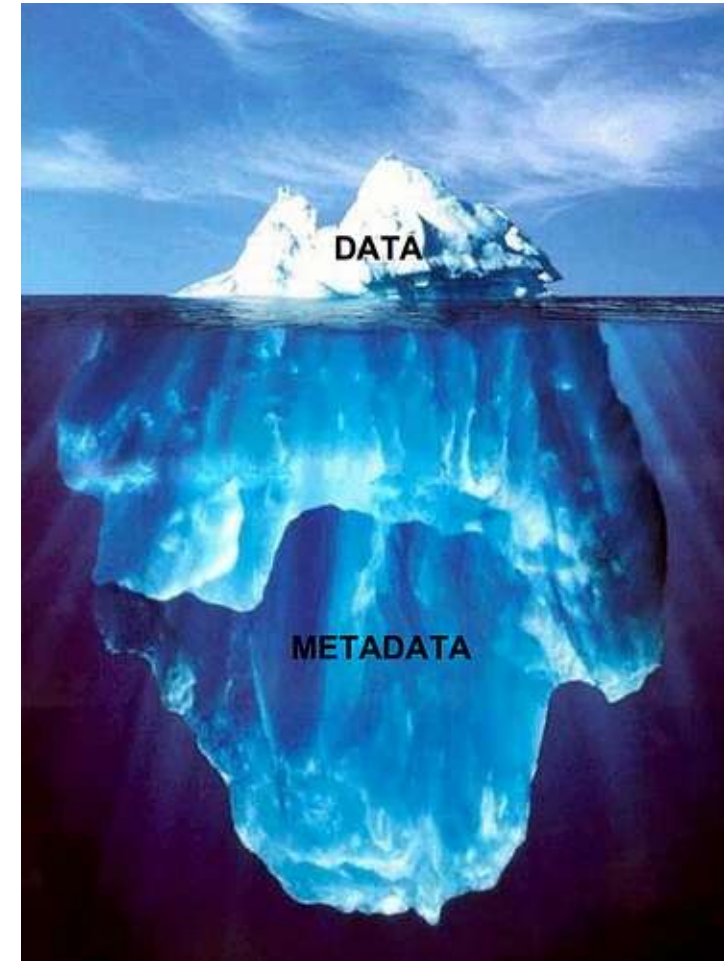  - descriptive
  - structural
  - administrative

photo by @kissane

# Provenance



- Libraries: Origin or source
- Museums: Chain of custody
- Internet: Provenance is information about entities, activities, and people involved in producing a piece of data or thing, which can be used to form assessments about its quality, reliability or trustworthiness. (World Wide Web Consortium (W3C) Provenance working group)

British Library, provenance record: Bestiary - caption: 'Owl mobbed by smaller birds'

# Reuse across place and time

- Reuse by investigator
- Reuse by collaborators
- Reuse by colleagues
- Reuse by unaffiliated others
- Reuse at later times
  - Months
  - Years
  - Decades
  - Centuries

Image from Soumitri Varadarajan blog. Iceberg image © Ralph A. Clevenger. Flickr photo

48

# MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21th 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
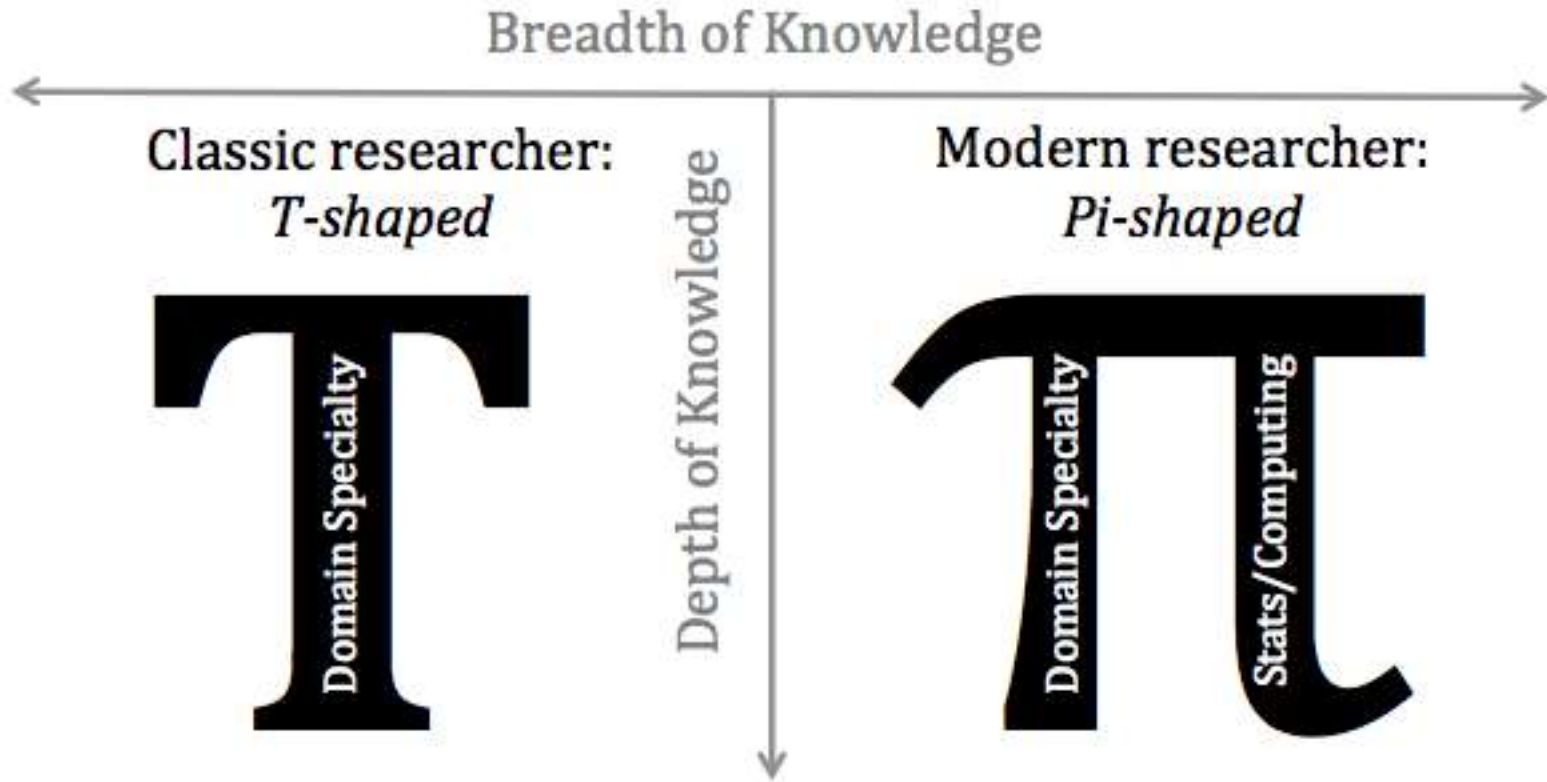
https://github.com/okul bilisim/awesome-datascience

# Data Curation and Stewardship
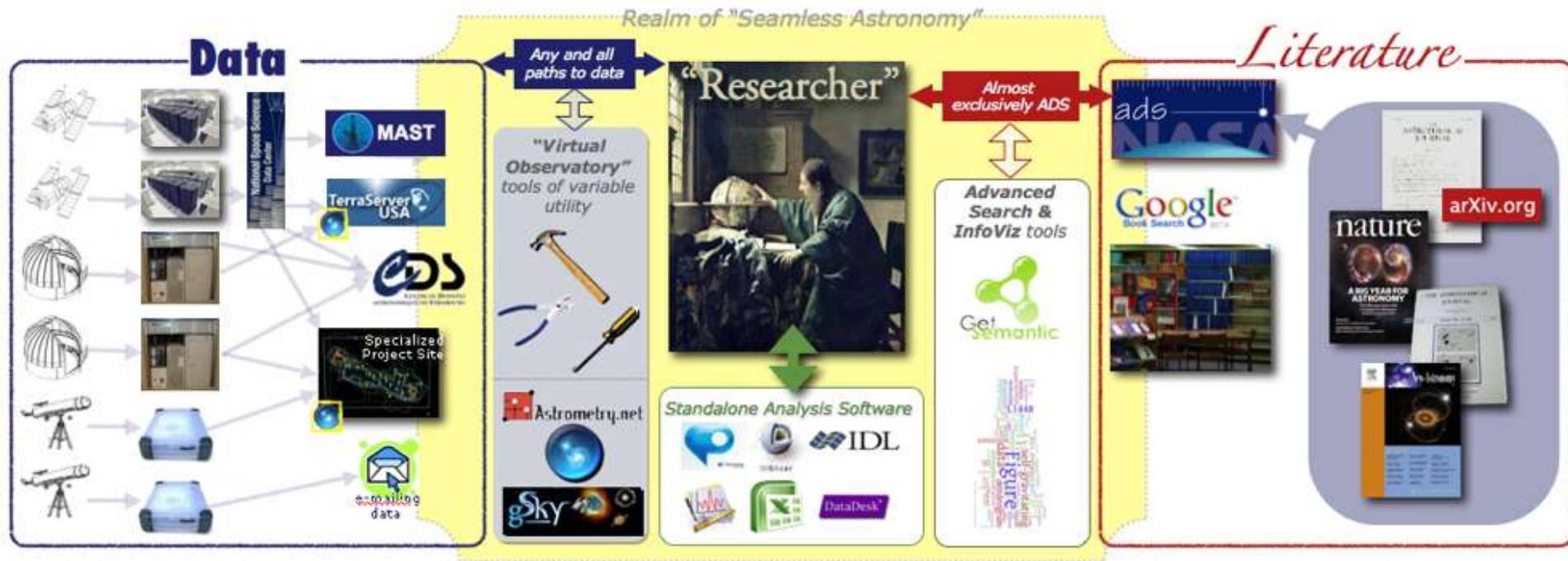
- Services and tools
- Data management planning
- Selection and appraisal
- Metadata, provenance
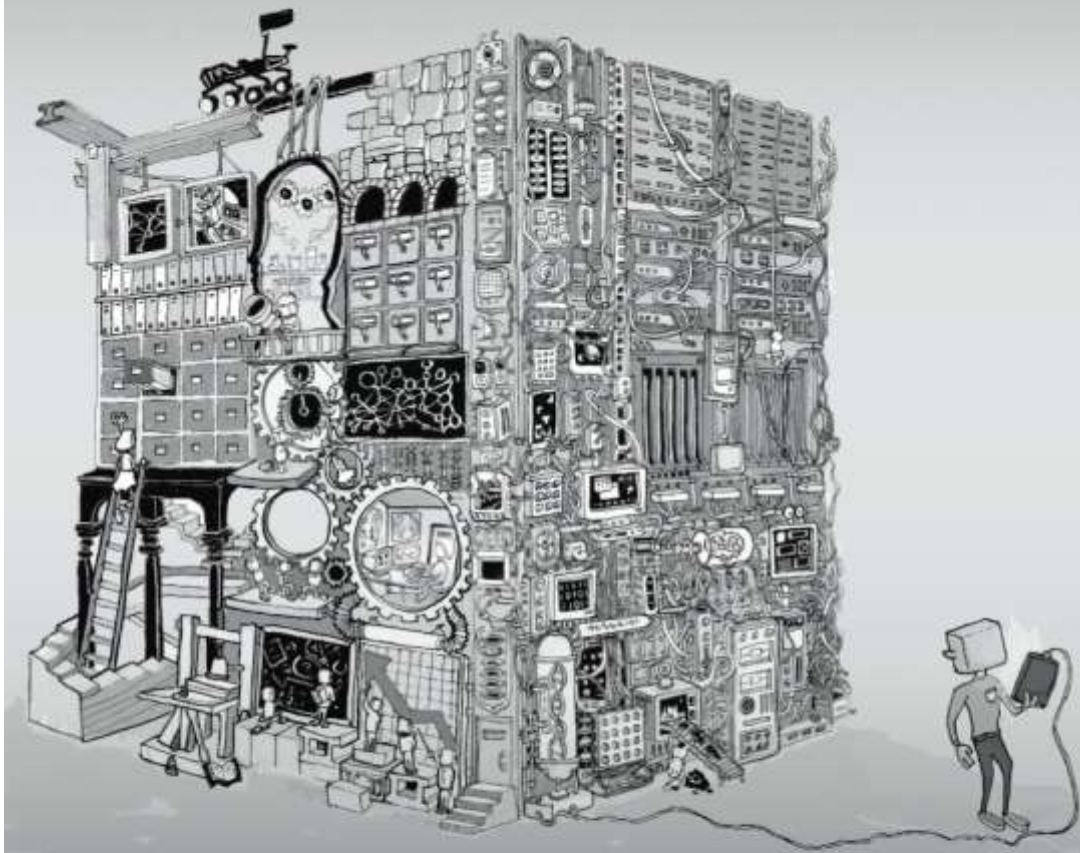- Migration
- Economics
- Infrastructure



ZOMBIE LIBRARIANS DISREGARD DATA!

# Research workforce



Breadth of Knowledge

Classic researcher:
*T-shaped*

Modern researcher:
*Pi-shaped*

Depth of Knowledge

Domain Specialty

Domain Specialty

Stats/Computing

# Knowledge Infrastructures



Image: Alyssa Goodman, Harvard Astronomy

Knowledge Infrastructures:
Intellectual Frameworks and Research Challenges

Report of a workshop sponsored by the National Science Foundation and the Sloan Foundation

University of Michigan School of Information, 25-28 May 2012

http://knowledgeinfrastructures.org

# Economics of the Knowledge Commons

| | | Subtractability / Rivalry | |
|---|---|---|---|
| | | Low | High |
| Exclusion | Difficult | **Public Goods** General knowledge Public domain data | **Common-pool resources** Libraries Data archives |
| | Easy | **Toll or Club Goods** Subscription journals Subscription data | **Private Goods** Printed books Raw or competitive data |

Adapted from C. Hess & E. Ostrom (Eds.), *Understanding knowledge as a commons: From theory to practice*. MIT Press.

# Data Repositories



This chart is based on the number of repositories in each Continent. However, some organisations have two or more repositories - over 20 in some cases - and this arguably skews the results.

For a different viewpoint, please see the equivalent chart for Repository Organisations, in which each organisation only counts once, regardless of how many repositories it hosts.
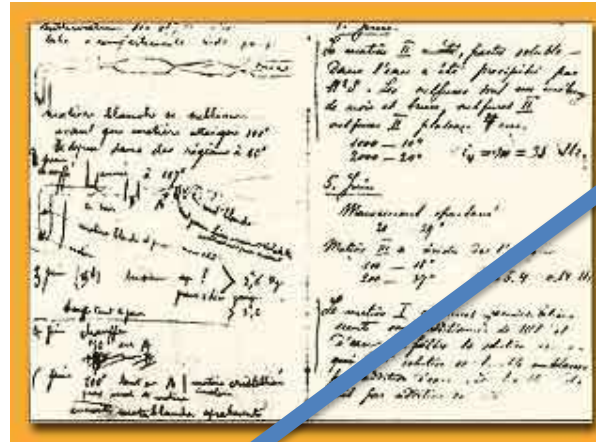
For further data, please see the corresponding table of repositories sorted by country.

Show embedding code

Show legacy chart and embedding code

# No Data

- Data not available

- Data not released

- Data not usable

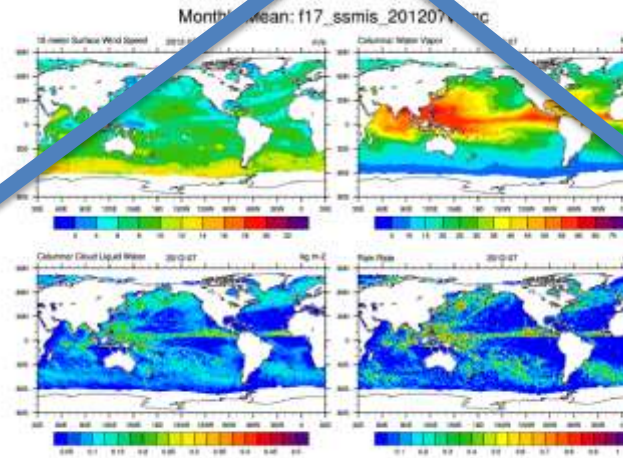Marie Curie's notebook aip.org

Pisa Griffin

hudsonalpha.org

Figure 2. Numeric Change in Resident Population for the 50 States, the District of Columbia, and Puerto Rico: 1990 to 2000
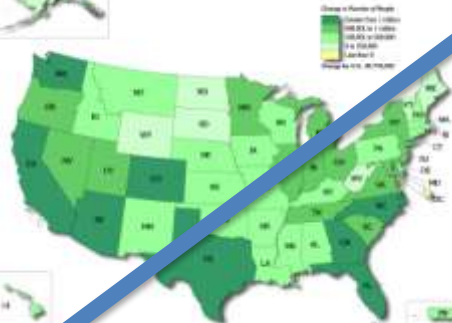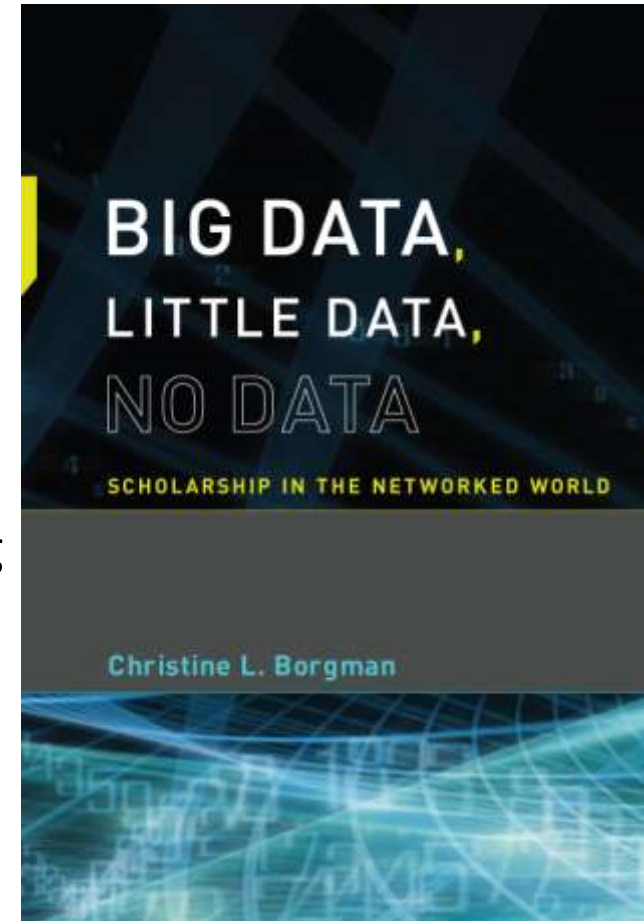
Monthly Mean: f17_ssmis_201207v.nc

Date:1/2.07.75        Place:Sakaltutan
Zafor
He will grow old in his present house; new house is for sons – 5 sons. Not sure they want to live in village. He will only build another if they want him to. eS came from Germany and did the plastering. He arranged the carpentry in Kayseri. Çok para gitti. {much money went} Has a tractor.

Date:Jul 1980        Place:Sakaltutan
Zafor:
Household now Zafor and wife; Nazif Unal and wife and youngest son still a boy. They run two dolmuß; one with a driver from Süleymanli. Goes in and out once a day. He gets 8,000 a month. Zafor then said, keskin değil. { not sharp - ? not profitable} I said he did very well on 8,000 TL with only two journeys a day. Nazif Unal has "bought" a Durak {dolmuß stop} from Belediye and works all day in Kayseri.

http://www.census.gov/population/cen2000/map02.gif

ncl.ucar.edu

http://onlineqda.hud.ac.uk/Intro_QDA/Examples_of_Qualitative_Data.p

56

# Conclusions

- Defining data
  - Representations used as evidence
  - One person's signal is another's noise
- Creating data
  - Models, questions, methods
  - Domain expertise
  - Data science expertise
- Collaborating with data
  - Documentation, description, identity, linking
  - Incentives for release and reuse
  - Curation and stewardship expertise
- Consolidating data value
  - Infrastructure and workforce investments
  - Value propositions
  - Trust fabric



BIG DATA,
LITTLE DATA,
NO DATA

SCHOLARSHIP IN THE NETWORKED WORLD

Christine L. Borgman

# Acknowledgements