

# データサイテーションマイニングによる 科学データの利活用分析

## Data Citation Analysis Framework for Open Science Data

是津 耕司, 村山 泰啓, 渡邊 堯

Koji Zettsu, Yasuhiro Murayama, Takashi Watanabe

国立研究開発法人 情報通信研究機構

National Institute of Information and Communications Technology



National Institute of Information and Communications Technology

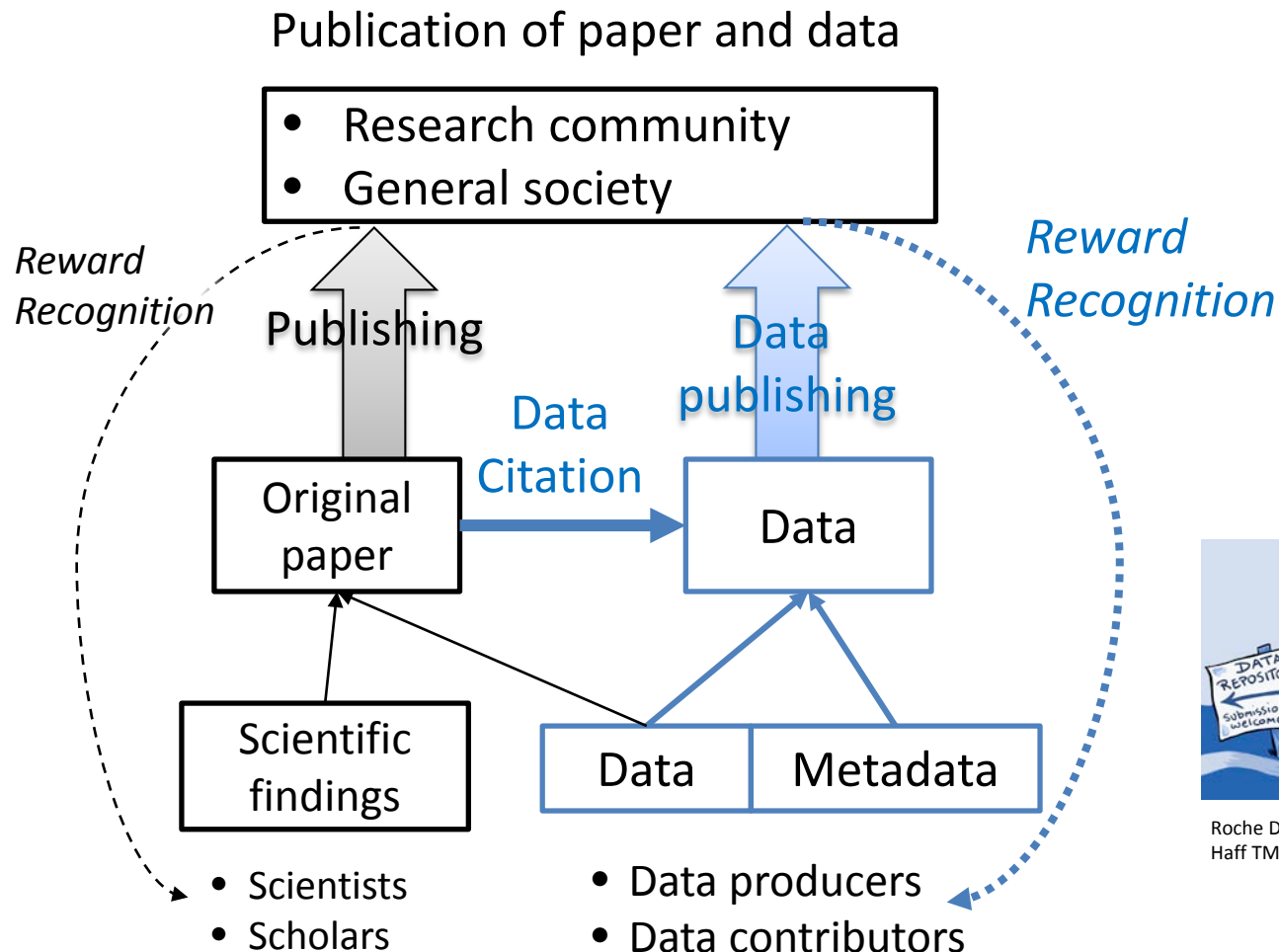
第2回オープンサイエンスデータ推進ワークショップ

2015年12月7-8日



# *Out of Cite, Out of Mind*

CODATA-ICSTI Task Group on Data Citation Standards and Practices (2013) Out of Cite, Out of Mind: The Current State of Practice, Policy, and Technology for the Citation of Data. Data Science Journal 12, CIDCR1–CIDCR75.



Roche DG, Lanfear R, Binning SA, Haff TM, Schwanz LE, et al. (2014)

Reference: Society of Geomagnetism, Earth, Planetary and Space Sciences, [http://sgepps.org/sgepps/shorai/SGEPSS\\_syorai\\_Jan2013.pdf](http://sgepps.org/sgepps/shorai/SGEPSS_syorai_Jan2013.pdf) [accessed on January 2013].

Data Publisher	Description	Data Citation Example
<b>PANGAEA:</b> The Publishing Network for Geo-scientific and Environmental Data <a href="http://www.Pangaea.de/">http://www.Pangaea.de/</a>	Open access library and data publisher for earth and environmental science	<i>Gershanovich, DE; Zinkovskiy, AB (1987): Distribution of particulate matter and particulate organic carbon in waters of the Caspian Sea. doi:10.1594/PANGAEA.756520</i>
<b>ICPSR:</b> The Inter-university Consortium for Political and Social Research <a href="https://www.icpsr.umich.edu/">https://www.icpsr.umich.edu/</a>	International consortium of about 700 academic institutions and research organizations that maintains and provides access to social science data	<i>Escarce, Jose J., Nicole Lurie, and Adria Jewell. RAND Center for Population Health and Health Disparities (CPHHD) Data Core Series: Pollution, 1988-2004 [United States]. ICPSR27864-v1. Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributor], 2011-10-21. <a href="http://doi.org/10.3886/ICPSR27864.v1">http://doi.org/10.3886/ICPSR27864.v1</a></i>
<b>Dryad</b> <a href="http://datadryad.org/">http://datadryad.org/</a>	International data repository of peer reviewed scholarly literature specialized in bioscience data	<i>López-Rodríguez MJ, Tierno de Figueroa JM (2012) Data from: Life in the dark: on the biology of the cavernicolous stonefly <i>Protonemura gevi</i> (Insecta, Plecoptera). The American Naturalist <a href="http://dx.doi.org/10.5061/dryad.8m8r1">http://dx.doi.org/10.5061/dryad.8m8r1</a></i>

... and more

*Source: CODATA-ICSTI Task Group on Data Citation Standards and Practices: Out of Cite, Out of Mind: The Current State of Practice, Policy, and Technology for the Citation of Data, Data Science Journal, Vol. 12, pp. CIDCR1-CIDCR75 (2013)*

- Data Citation Index (DCI) [Thomson Reuters]
  - Harvests citations to research data from papers indexed in the Web of Science



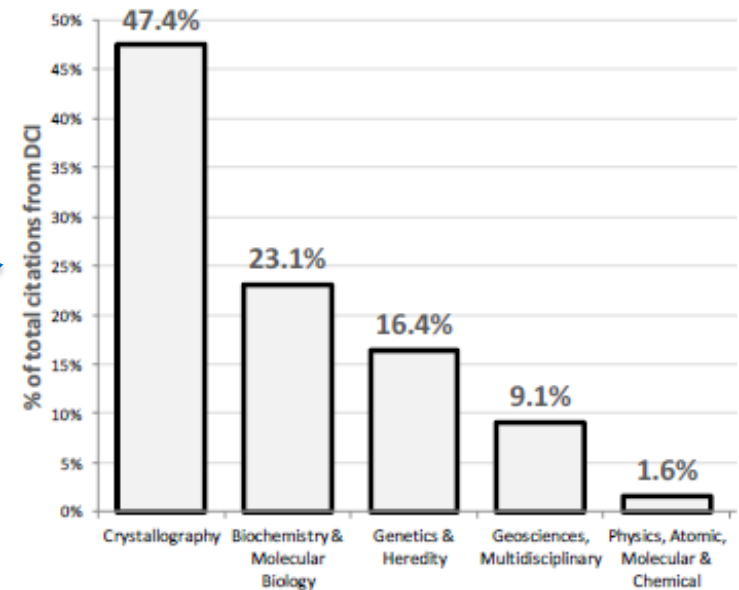
## A. DCI overall statistics

	All Document Types	Data repositories	Data studies	Data sets
Total Citations	404,211	3,265	106,895	294,051
Total Records	2,623,528	90	154,674	2,468,736
Uncited Records	2,311,553	63	126,428	2,185,062
% Uncited	88.11	40.0	81.74	<u>88.51</u>
Citation Average	0.15	36.28	0.69	0.12
Standard Deviation	3.06	336.07	9.56	0.36

## B. DCI statistics by area of data studies

	Total Records	% Records	Total Citations	% Citations	Citation Average
Engineering & Technology	1,545	0.06	890	0.30	0.58
Humanities & Arts	44,588	1.81	1	0.00	0.00
Science	2,004,449	81.19	293,193	99.71	0.15
Social Sciences	424,952	17.21	7	0.00	0.00

## C. Distribution of citations by subject



**Source:** Robinson-Garcia, N. et. al: Analyzing data citation practices using the Data Citation Index, *Journal of the Association for Information Science and Technology*, DOI: 10.1002/asi.23529 (June, 2015)

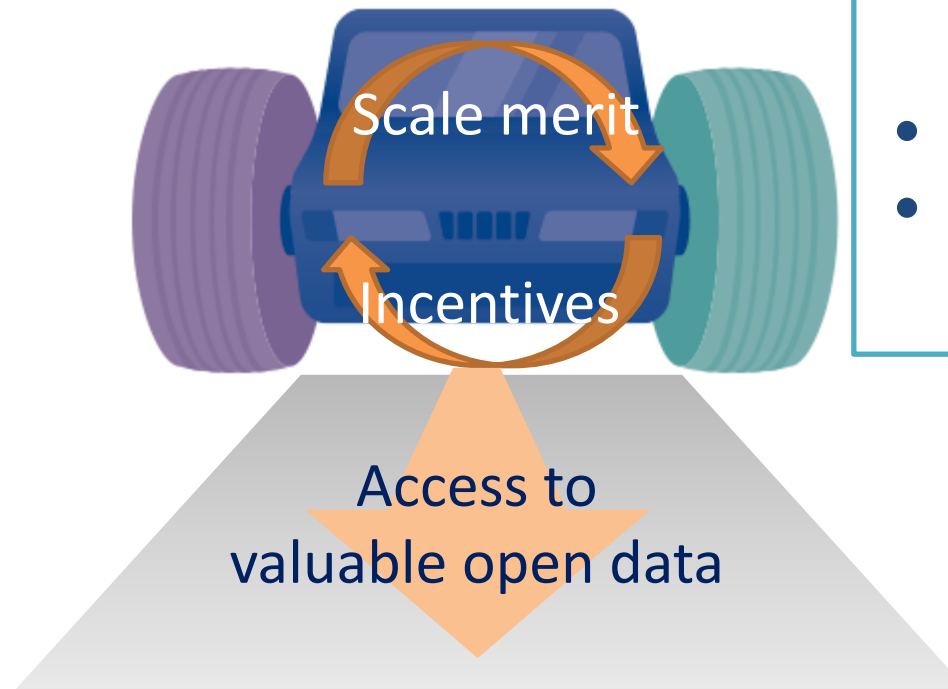
## Data Citation

### Provision

- Principles
- Standards & practices

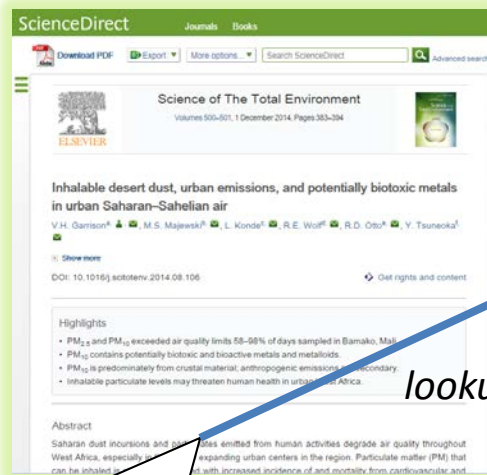
### Usage

- Data evaluation
- Data discovery



## Open Science

## Online Journal



Garrison, VH et al. (2014): Particulate matter (PM<sub>2.5</sub> and PM<sub>10</sub>) in the air in Bamako, Mali (2012-2013). doi:10.1594/PANGAEA.834195

## Data citation



forward

lookup

- DataCite
- CrossRef
- JaLC, etc.

## Data Publisher

Landing page



Metadata

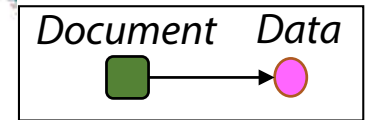
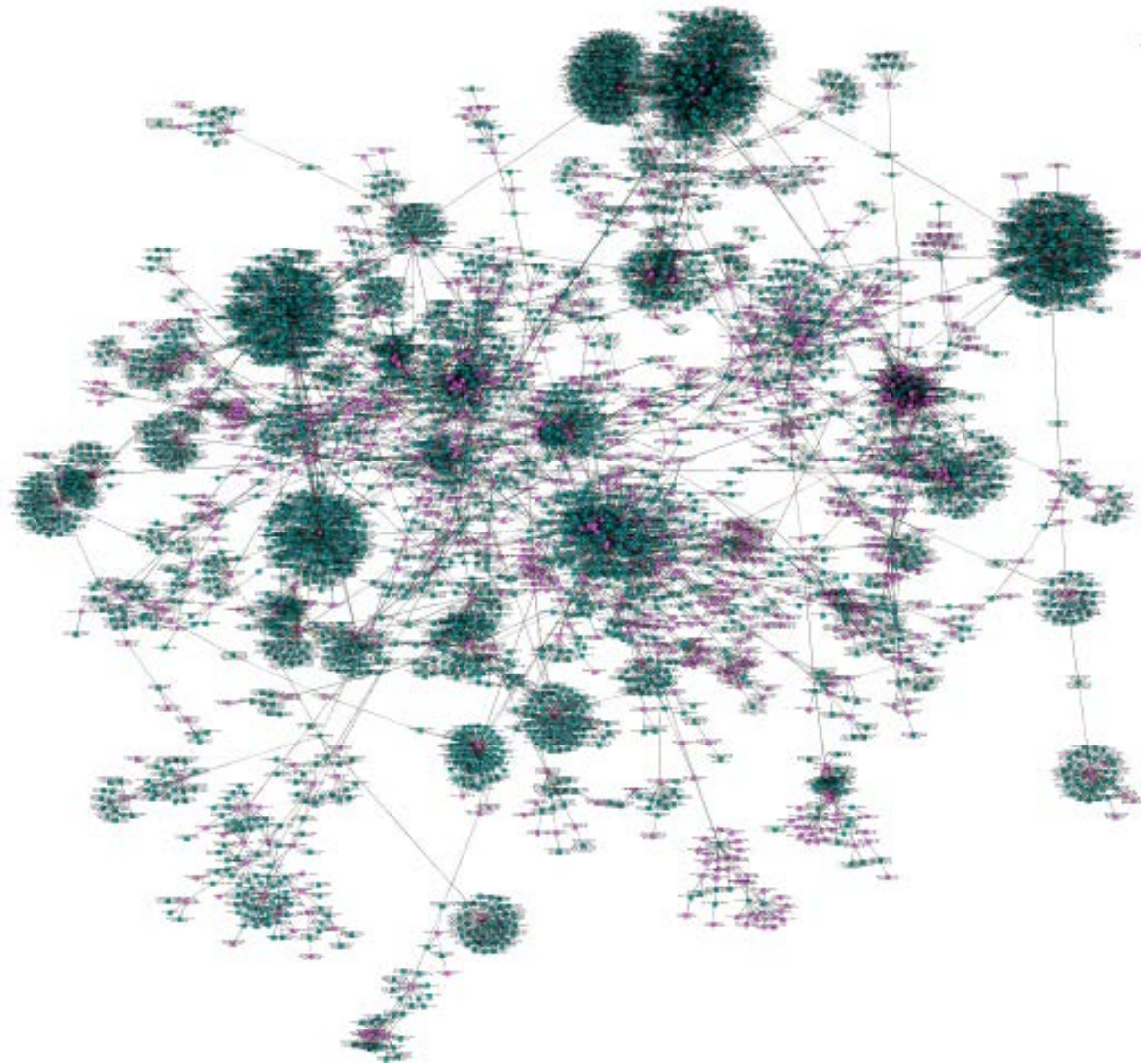
DOI

10.1594/PANGAEA.834195

Date/Time	PM <sub>2.5</sub> (µg/m <sup>3</sup> )	PM <sub>10</sub> (µg/m <sup>3</sup> )	TPP (µg/m <sup>3</sup> )
2012-04-10	39.5	139.9	
2012-04-14	20.0	147.8	
2012-04-16	16.7	146.4	
2012-04-18	17.7	98.3	131.0
2012-04-17	17.6	82.8	
2012-04-18	11.2	56.8	83.4
2012-04-19	15.7	83.5	
2012-04-20	14.2	67.7	91.1
2012-04-21	18.4	72.7	118.2
2012-04-22	22.8	89.2	166.1
2012-04-23	29.4	82.8	96.6
2012-04-24	29.3	32.6	146.1
2012-04-25	26.7	116.0	198.9
2012-04-26	37.6	154.3	

Data

- PANGAEA, ICPSR, Dryad, etc.



**Inter-university Consortium for Political and Social Research (ICPSR)**  
115,154 citations

- **Macro analysis:**

Analyze structure of data citation network

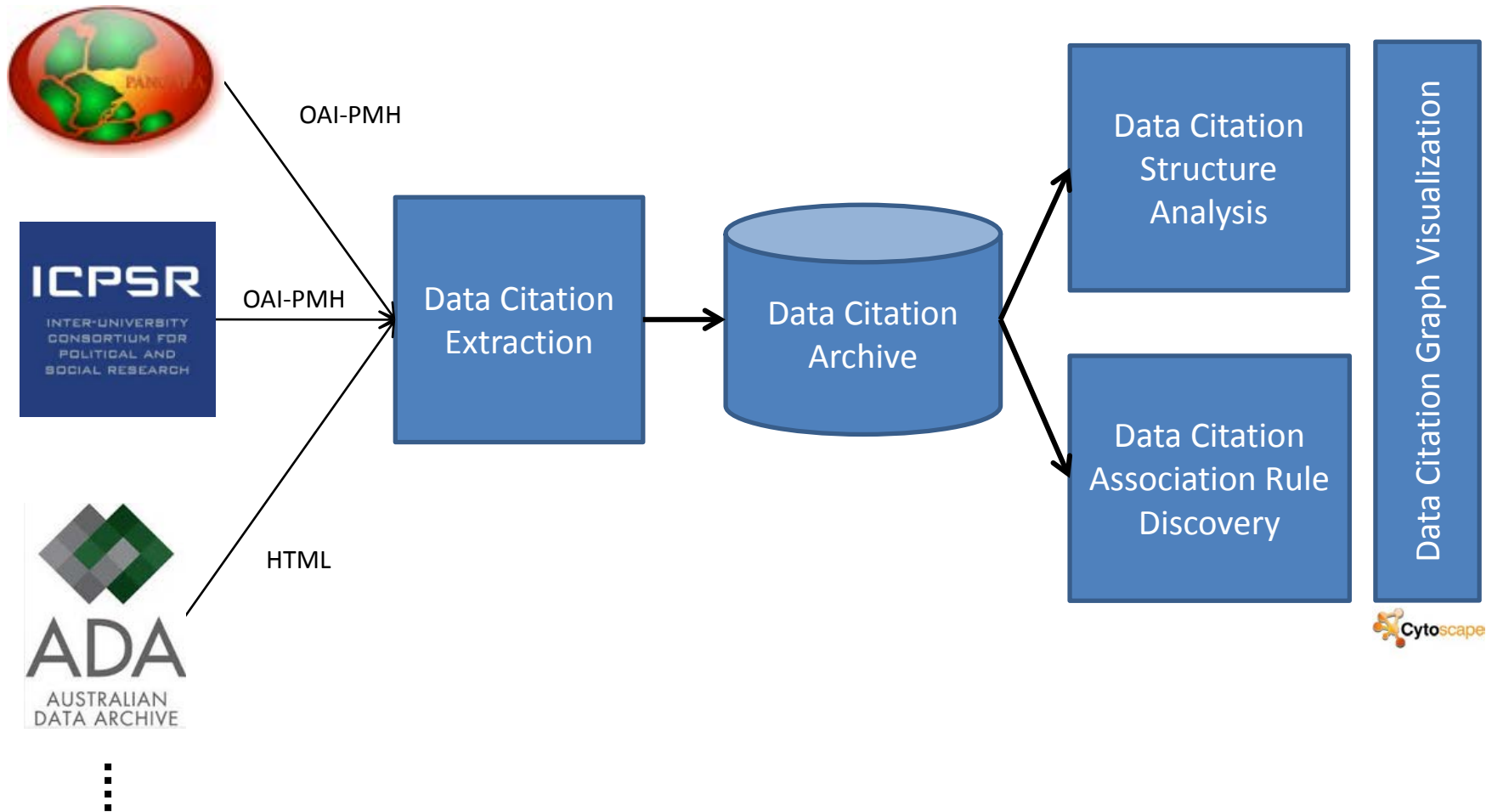
- Discover communities of data citation, and characterize data by citations in a community

- **Micro analysis:**

Analyze associations between document and data

- Discover typical associations (*i.e.* **association rules**) between documented knowledge and evidential data

## Data Publisher



## Data Citation Metadata

**PANGAEA**  
Data Publisher for Earth & Environmental Science

**Data Description**

**Citation:** Garrison, VH et al. (2014): Particulate matter (PM2.5 and PM10) in the air in Bamako, Mali (2012-2013). doi:10.1594/PANGAEA.834195.  
In Supplement to: Garrison, Virginia H; Majewski, Michael S; Konde, Lassana; Wolf, Ruth E; Otto, RD; Tsuneoka, Y (2014): Inhalable desert dust, urban emissions, and potentially biotoxic metals in urban Saharan-Saharan air. *Science of the Total Environment*, 500-501, 383-394, doi:10.1016/j.scitotenv.2014.08.105

**Coverage:** Latitude: 12.428330 \* Longitude: -8.320000  
Date/Time Start: 2012-09-13T00:00:00 \* Date/Time End: 2013-07-09T00:00:00

**Location:** Bamako \* Latitude: 12.428330 \* Longitude: -8.320000 \* Elevation: 325 ft asl  
\* Location: Mali \* \* Device: Air chemistry laboratory \*  
\* [View map](#) [Download Earth](#)

Parameter(s)	Name	Short Name	Unit	Principal Investigator	Method	Comment
DATE/TIME %	Date/Time					Records
Particulate matter < 2.5 PM2.5	µg/m³	Garrison, Virginia H	Gravimetric analysis %			
Particulate matter < 10 PM10	µg/m³	Garrison, Virginia H	Gravimetric analysis %			
Total suspended particulates	TSP	µg/m³	Garrison, Virginia H	Gravimetric analysis %		

**Download Data**  
Download dataset as tab-delimited text (use the following character encoding: ISO)  
View dataset as HTML

**Citation:** Garrison, VH et al. (2014): Particulate matter (PM2.5 and PM10) in the air in Bamako, Mali (2012-2013). doi:10.1594/PANGAEA.834195.  
**In Supplement to: Garrison, Virginia H; Majewski, Michael S; Konde, Lassana; Wolf, Ruth E; Otto, RD; Tsuneoka, Y (2014): Inhalable desert dust, urban emissions, and potentially biotoxic metals in urban Saharan-Saharan air. *Science of the Total Environment*, 500-501, 383-394, doi:10.1016/j.scitotenv.2014.08.105**

**Related to:**

**Bers, A Valeria; Momo, Fernando; Schloss, Irene R; Abele, Doris (2013):** Analysis of trends and sudden changes in long-term environmental data from King George Island (Antarctica): relationships between global climatic oscillations and local system response. *Climatic Change*, **116**, 789-803 [↗](#)

**Klöser, Heinz; Ferreyra, Gustavo A; Schloss, Irene R; Mercuri, Guillermo; Laternus, Frank; Curtosi, Antonio (1993):** Seasonal variation of algal growth conditions in sheltered Antarctic bays: the example of Potter Cove (King George Island, South Shetlands). *Journal of Marine Systems*, **4**, 289-301, doi:10.1016/0924-7963(93)90025-H [↗](#)

**Schloss, Irene R; Abele, Doris; Moreau, Sébastien; Demers, Serge; Bers, A Valeria; Gonzáles, Oscar; Ferreyra, Gustavo A (2012):** Response of phytoplankton dynamics to 19-year (1991-2009) climate trends in Potter Cove (Antarctica). *Journal of Marine Systems*, **92**, 53-66 [↗](#)

**Schloss, Irene R; Ferreyra, Gustavo A (2002):** Primary production, light and vertical mixing in Potter Cove, a shallow bay in the maritime Antarctic. *Polar Biology*, **25**, 41-48, doi:10.1007/s003000100309 [↗](#)

**Schloss, Irene R; Ferreyra, Gustavo A; Ruiz-Pino, Diana (2002):** Phytoplankton biomass in Antarctic shelf zones: a conceptual model based on Potter Cove, King George Island. *Journal of Marine Systems*, **36**, 129-143, doi:10.1016/S0924-7963(02)00183-5 [↗](#)

## Related Publications

- 1988 Long, Larry E. [Migration and Residential Mobility in the United States](#). New York: Russell Sage.  
Export Options: [RIS/EndNote](#)
- 1964 Eldridge, Hope T., Thomas, Dorothy Swaine. [Demographic Analyses and Interrelations. Population Redistribution and Economic Growth, United States, 1870-1950 series, vol. 3.](#) Philadelphia, PA: American Philosophical Society.  
Export Options: [RIS/EndNote](#)
- 1960 Kuznets, Simon Smith, Miller, Ann Ratner, Easterlin, Richard A. [Analyses of Economic Change. Population Redistribution and Economic Growth: United States, 1870-1950 series, vol. 2.](#) Philadelphia, PA: American Philosophical Society.  
Export Options: [RIS/EndNote](#)

Landing pages

**ICPSR Find & Analyze Data**

Find Data Search/Compare Variables Find Publications Resources for Students Get Help

**Population Redistribution and Economic Growth in the United States: Population Data, 1870-1960 (ICPSR 7753)** [↗](#)

**Principal Investigator(s):** Kuznets, Simon; Thomas, Dorothy Swaine

**Summary:**  
Detailed demographic characteristics of the population of the United States from 1870 to 1960 are contained in this data collection. Included are state-level estimates of the nation's inhabitants by sex, race, nativity and age, as well as intercessal migration calculated by age, race, and sex. The basic information recorded in this collection was obtained from the decennial censuses of the United States or estimated by the principal investigators from material collected by the decennial cen... [\(more info\)](#)

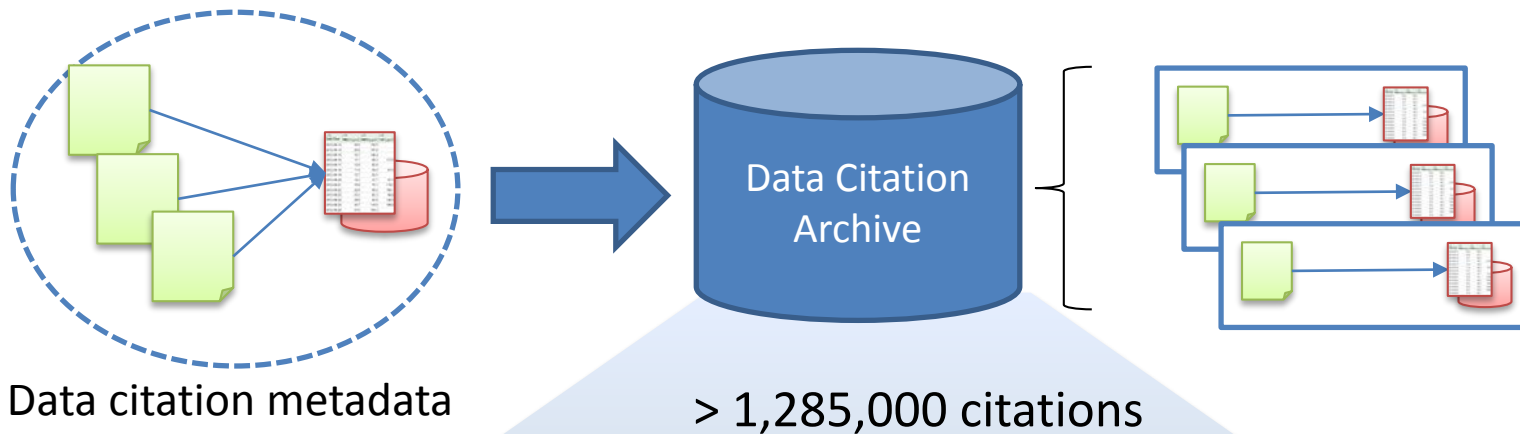
**Access Notes**

- These data are available only to users at ICPSR member institutions. Because you are not [logged in](#), we cannot verify that you will be able to download these data.

**Dataset(s)**

**DS0: Study-Level Files**  
Documentation: [Codebook.txt](#)

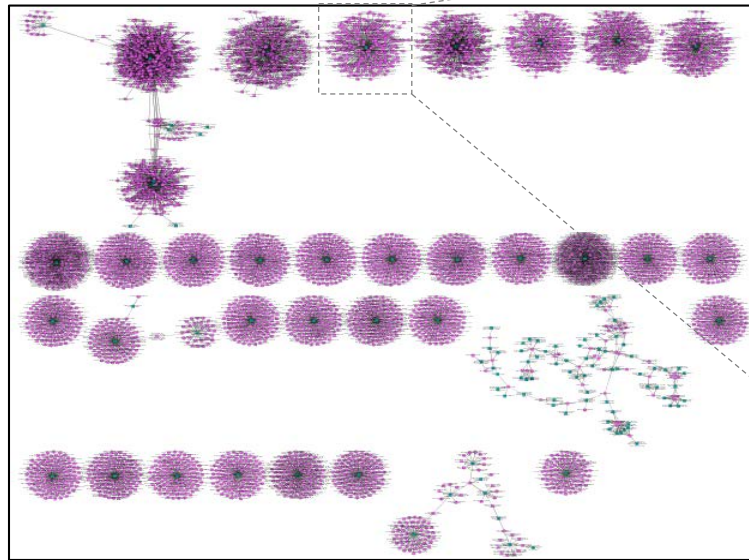
**DS1: Native White Population by Age and Sex, 1870-1960** [Download All Files \(0.7 MB\)](#)  
Data: [ASCII](#) [SAS Setup](#) [SPSS Setup](#) [Stata Setup](#)



Data Site	Domain	# of DC
<a href="#"><u>Pangaea</u></a>	Earth & Environment	322,477
<a href="#"><u>ICPSR</u></a>	Social Science	114,815
<a href="#"><u>DataCite</u></a>	(any)	773,173
<a href="#"><u>DRYAD</u></a>	Bioscience	1,556
<a href="#"><u>ADA</u></a>	Social Science	16,062
<a href="#"><u>ESDS</u></a>	Economic & Social Science	59,471

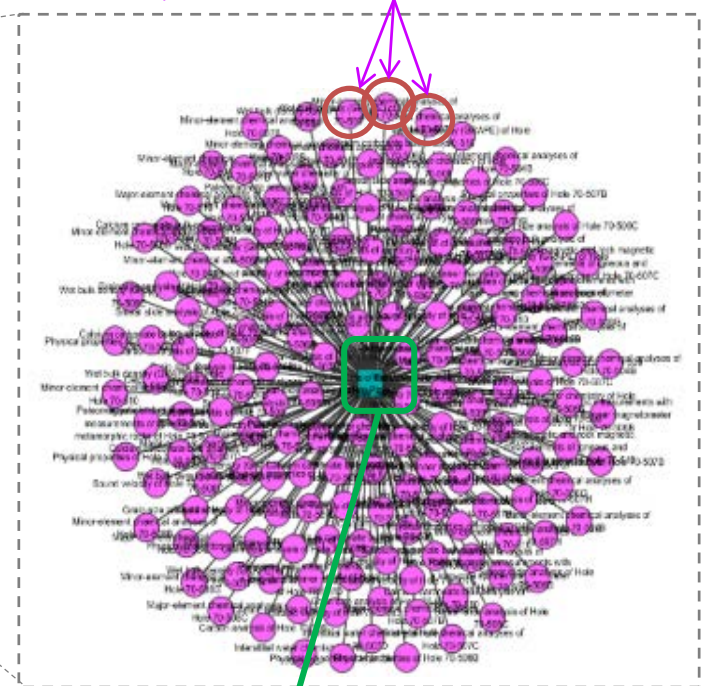
- Data Collection Community

Data collection

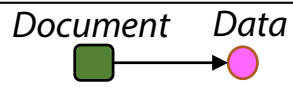


Pangaea

*“Physical properties of Hole #”*



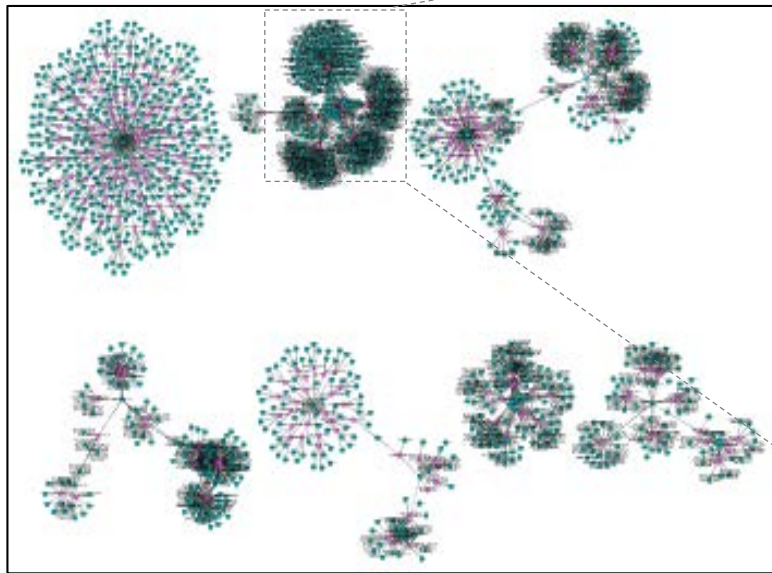
*“Reports of the Deep Sea Drilling Project”*



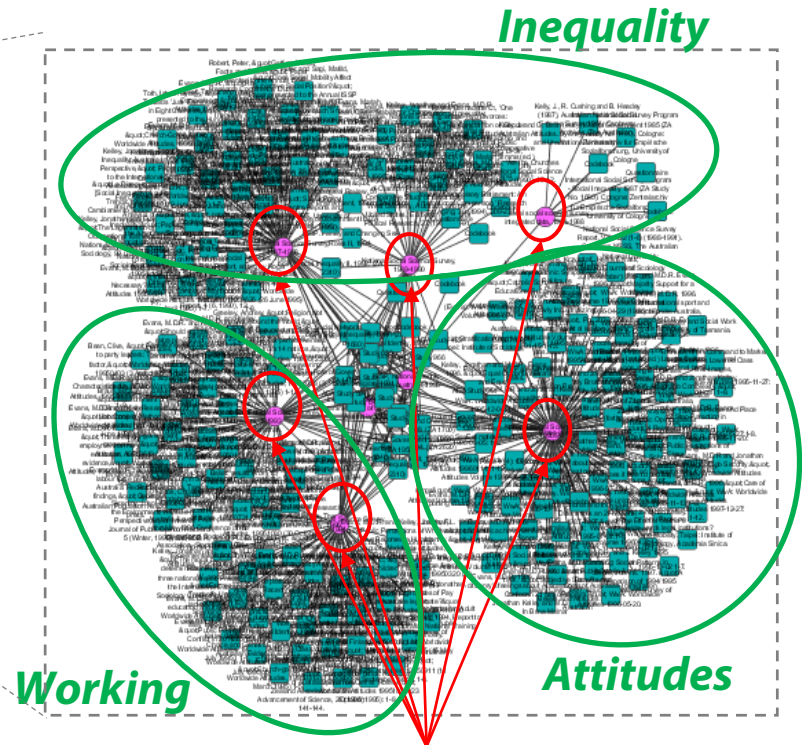
Catalogue document

- Data Sharing Community

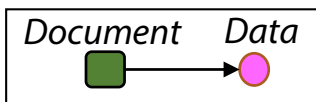
Data-sharing document clusters



Australian Data Archive (ADA)

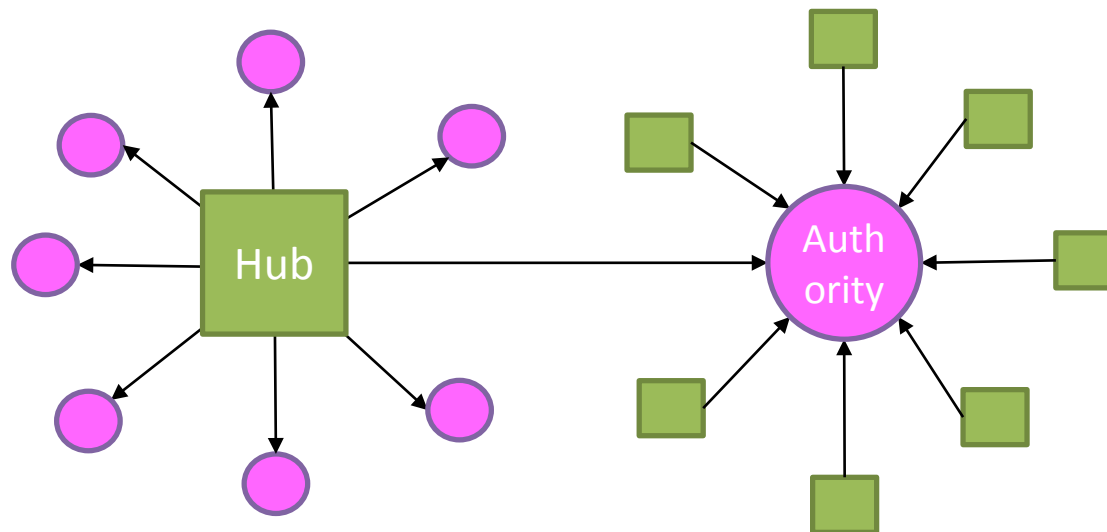


National Social Science Survey



Shared data

- HITS algorithm [Kleinberg 99]
  - A good **hub** links to many good authorities
    - Hub score:  $H(x) = \sum_{y \leftarrow x} A(y)$
  - A good **authority** is referenced by many good hubs
    - Authority score:  $A(x) = \sum_{y \rightarrow x} H(y)$
  - Discovery from result set





# Community Discovery Demo

NICT society Search Advanced Search Style Tools

Hub Authority

1. Culture Shift in Advanced Industrial Society  
Hub: 0.82563  
Repository: ICPSR  
Author: Inglehart, Ronald  
Date: 1990
2. The new political culture: Changing dynamics of support for the welfare state and other policies in postindustrial societies  
Hub: 0.20182  
Repository: ICPSR  
Author: Clark, Terry  
NicholsHoffmann-Martinot, VincentInglehart, Ronald  
Date: 1998
3. Economics and the vulnerability of the pan-European institutions  
Date: 1998

Document Data

Advanced Search

Search:

NICT sediment Search Advanced Search Style Tools

Hub Authority

1. Report and preliminary results of Meteor Cruise 23/2, Rio de Janeiro-Recife, 27.02.-19.03.1993  
Hub: 0.99664  
Repository: PANGAEA  
Author: Bleil, Ulrich  
Date: 1993
2. Paläo-Ozeanographie des Europäischen Nordmeeres an Hand stabiler Kohlenstoff- und Sauerstoffisotope  
Hub: 0.0162  
Repository: PANGAEA  
Author: Vogelsang, Elke  
Date: 1990
3. Zur Sedimentationsgeschichte von biogenem Opal im nördlichen Nordatlantik und dem Europäischen Nordmeer  
Hub: 0.0444

Document Data

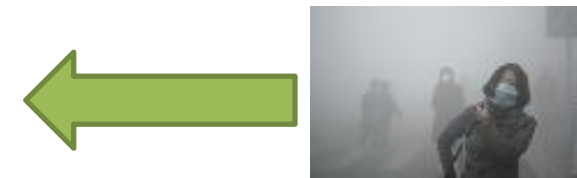
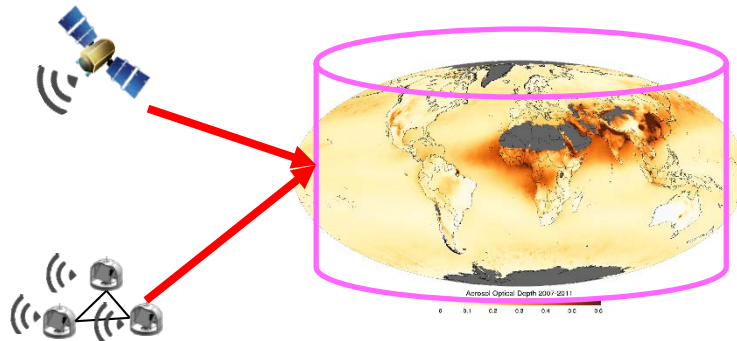
Advanced Search

Search:

Provider's view

User's view

“aerosol optical depth”



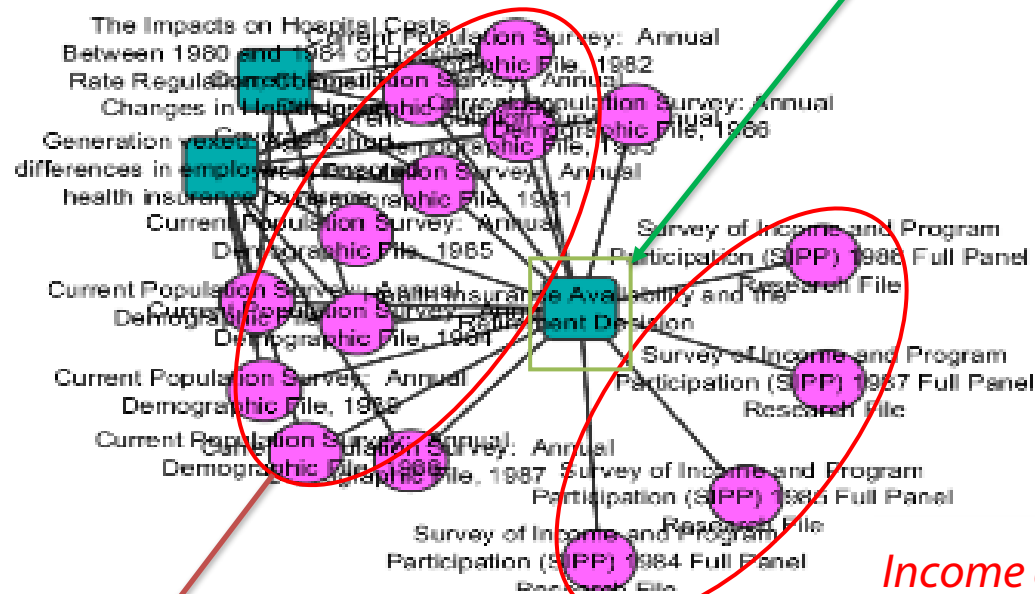
“air pollution data”

- Data usage is often different from the data content

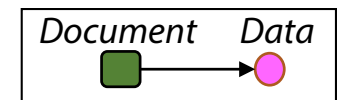
Health insurance

Population data

Income data

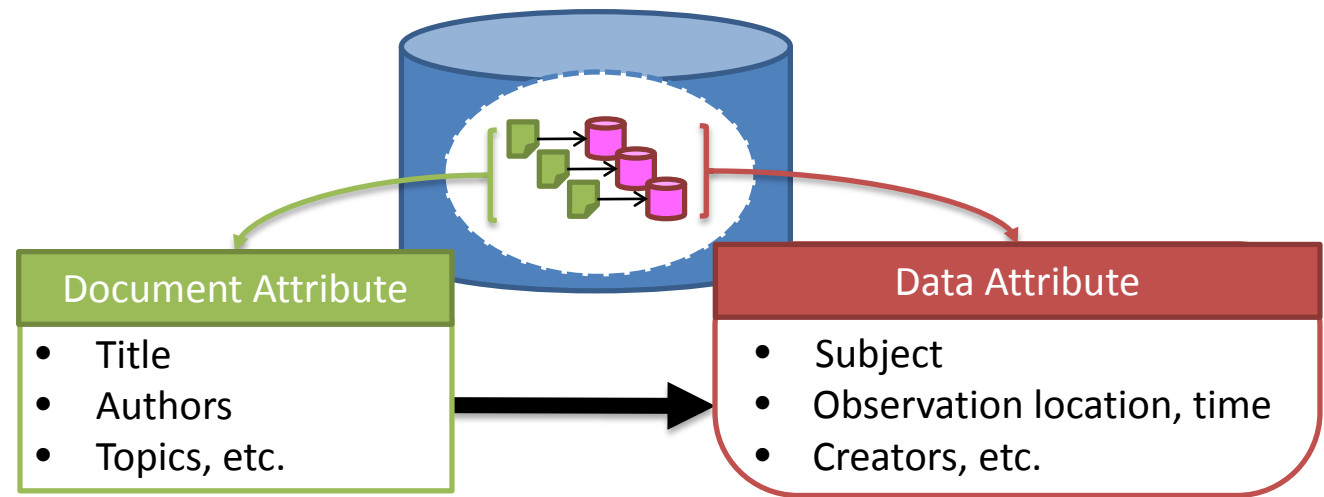


Both population data and income data are referenced by health insurance document



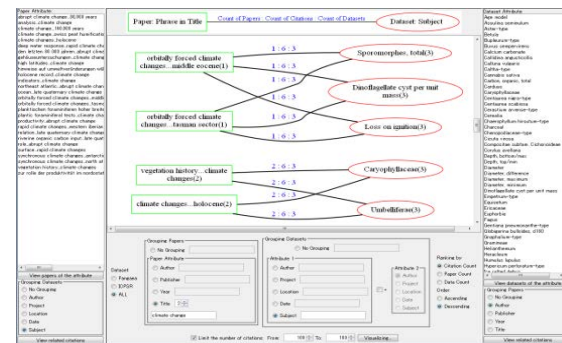
- Discover typical combinations of document and data attributes frequently co-occurring in data citations (referential contexts)

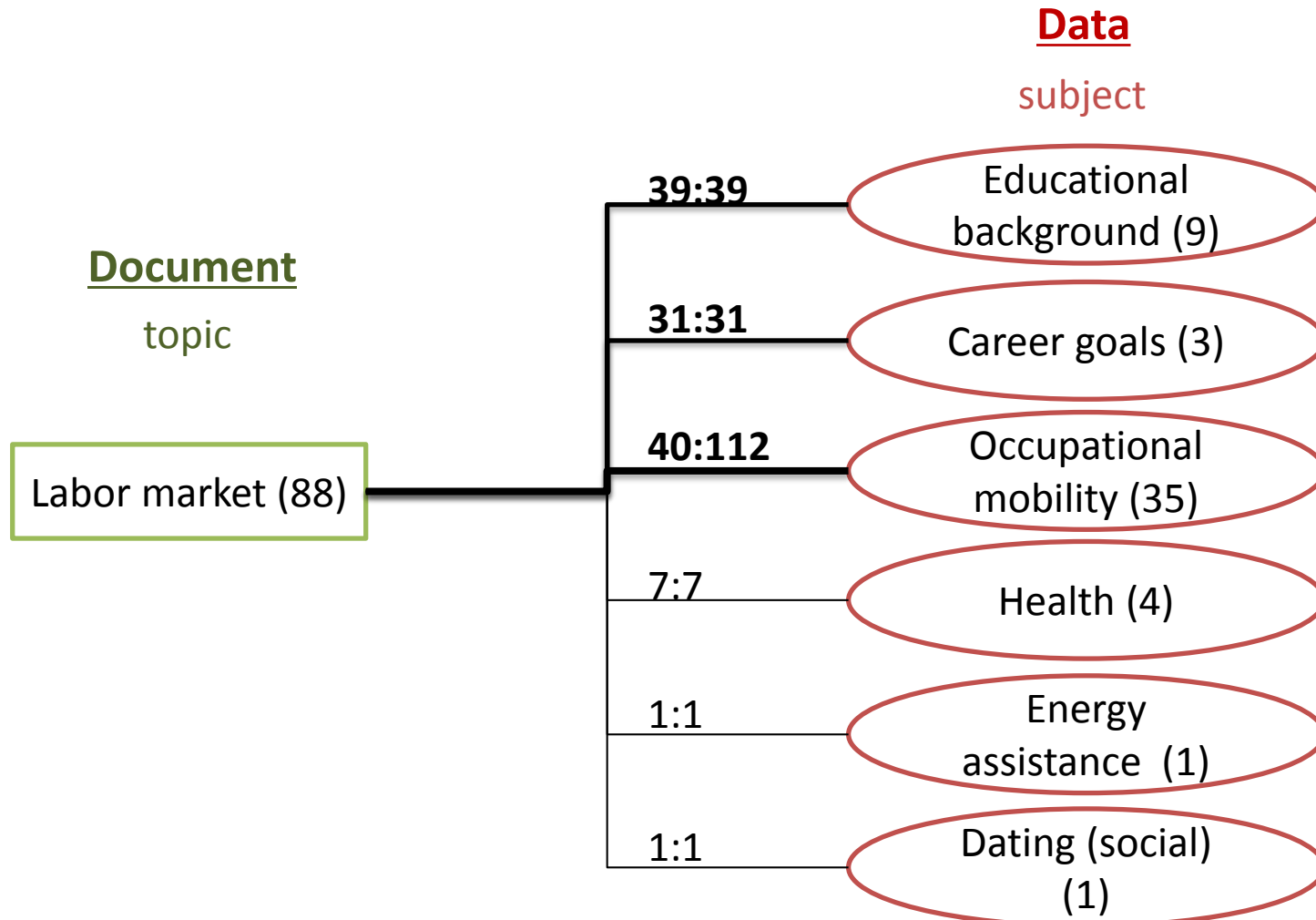
## Data Citation Archive



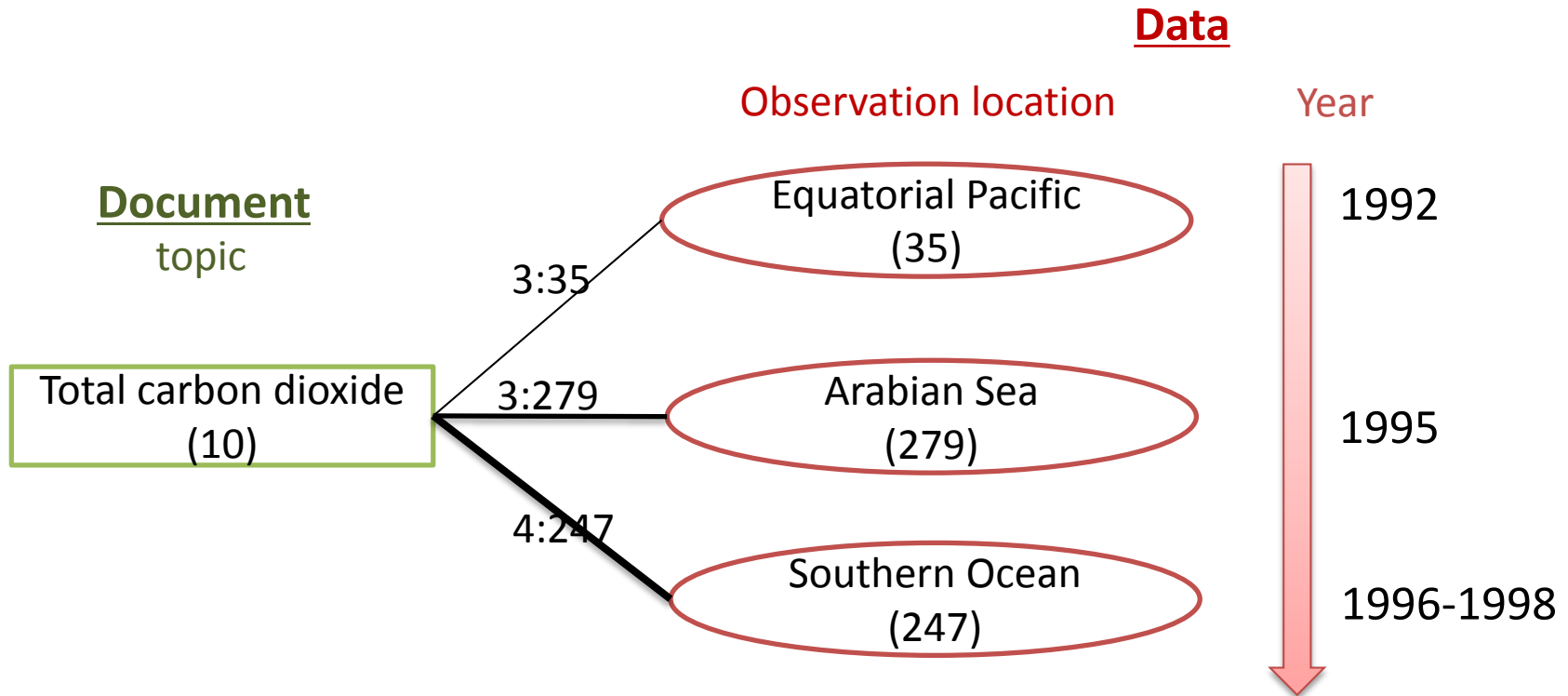
Referential Context:

Association rule discovery tool





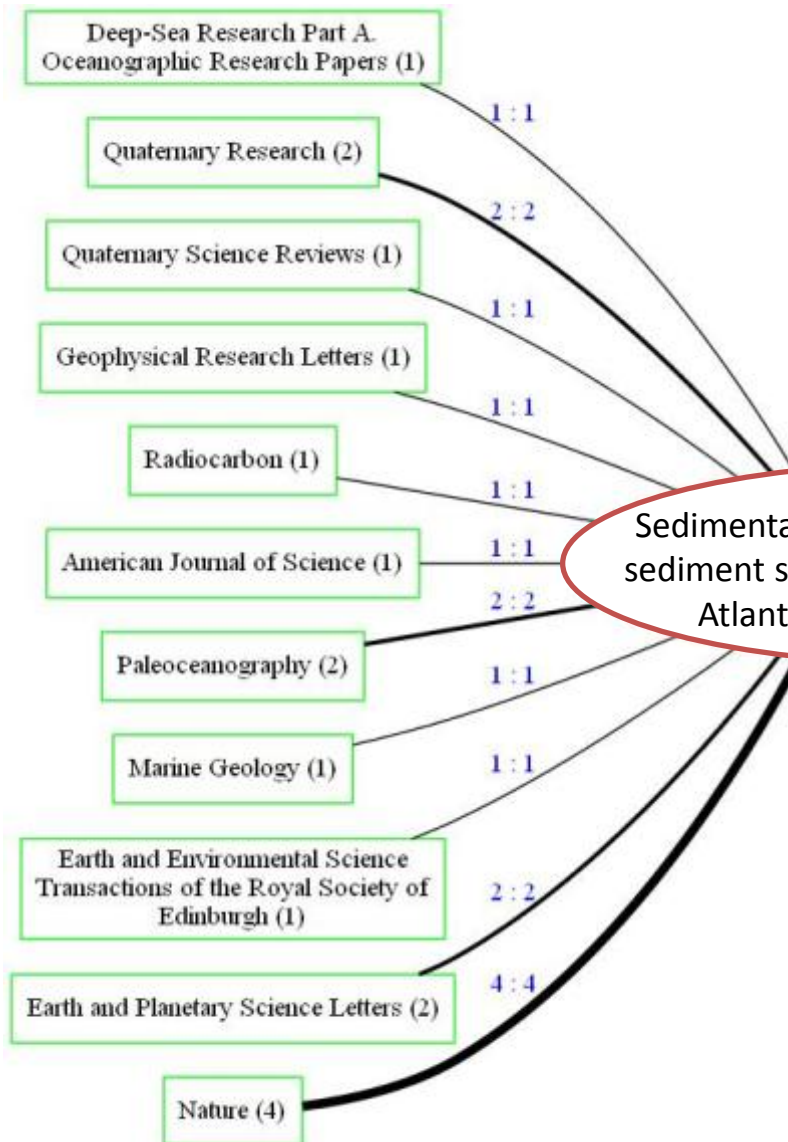
**Data citations** (#source: #target)



*Observation data for carbon dioxide research goes south over years.*



## Document



**Data**  
title

Data collected from 1963 ~ 1981 all over the world



Sedimentation rates calculated on surface sediment samples from different site of the Atlantic and Pacific Oceans, 1991.

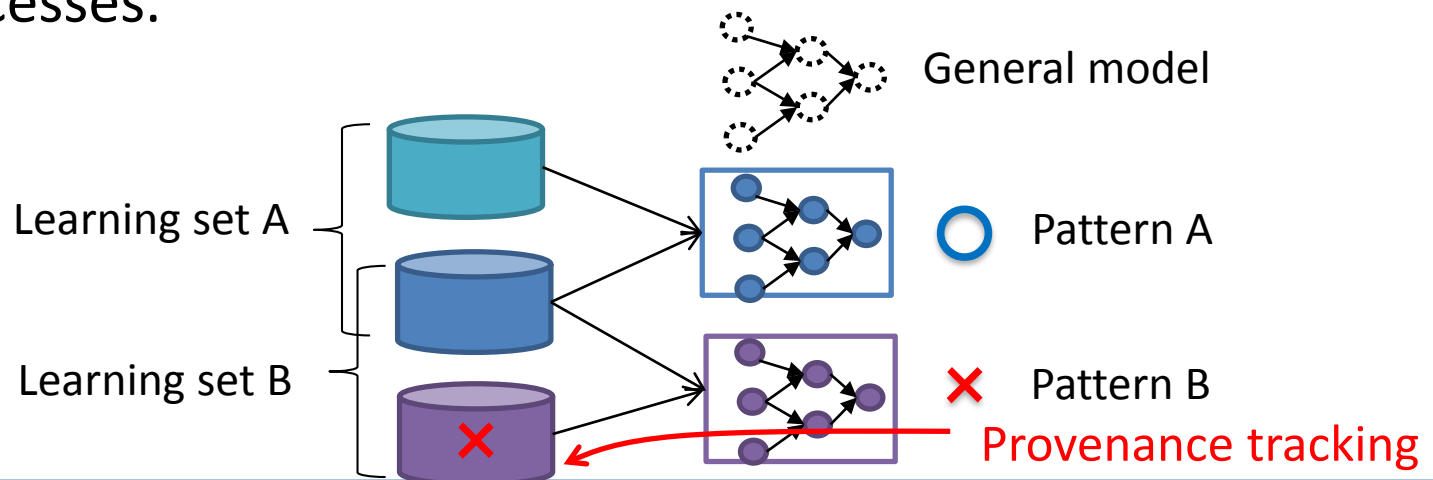


Created by  
**Wallace Smith Broecker** (1931 - )

- Discover data-intensive community
  - Collaborate with research communities having same or similar data [*researcher*]
  - Survey the data common to a research community [*data repository*]
- Evaluate reputation of data
  - Reward data (creators) based on its popularity and/or authority [*funding organization*]
  - Manage quality of data [*data creator*]
- Provide superior discoverability for better reuse of data
  - Search data by both content and context keywords [*data repository*]
  - Discover related data for interdisciplinary research [*data curator*]
  - Find research publications actionable for reuse of data [*researcher, publisher*]

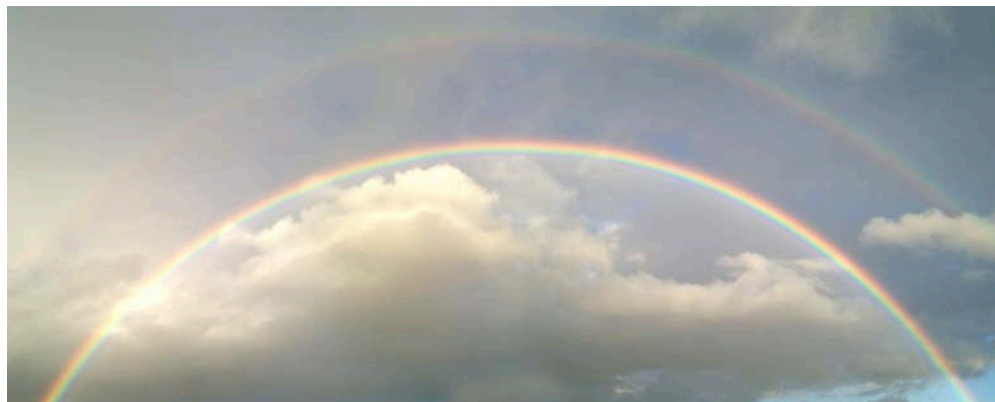
- Unstable metadata of data citation
  - Semantic compatibility among heterogeneous data citation metadata
    - E.g.) *relate to, supplement to (PANGAEA), related publications (ICPSR), is referenced by (Dryad), related materials (ADA), ....*
  - Up-to-date?
- More citations for better analysis
  - Unified and/or centralized access to distributed information of data citation
  - Citation-creating applications with harnessing citation analysis
    - E.g.) Data search with citation-based ranking (*more citations, more exposure*) → data citation optimization

- Analyzing *dynamic* citations
  - **Behavior analysis** based on user-to-data citations obtained from data (DOI) access log
    - E.g.) “Users Who Took This Item Also Took” (social filtering)
  - **Intention analysis** based on keyword-to-data citations obtained from search query log
  - **Quality analysis of data-intensive AI models** based on process-to-data citations obtained from machine learning processes.



- Data citation = link from documented knowledge to evidential data
  - Instead of knowledge-to-knowledge link by document citation
- Analysis of ‘Web of data citation’
  - Data citation structure analysis (macro analysis)
  - Data citation association rule discovery (micro analysis)
- For better reuse & reward of data
  - Discover data-intensive community
  - Evaluate reputation of data
  - Provide superior discoverability

# THANK YOU



Data Citation Mining  
are going to  
open source software project

Contact: [zettsu@nict.go.jp](mailto:zettsu@nict.go.jp)