

Status and Trajectory of Supporting Tools and Systems

Julia Lane, Paco Nathan and Sophie Rand

The Coleridge Initiative

BIPARTISAN POLICY CENTER
Foundations for Evidence-Based Policymaking Act of 2018

The bipartisan Foundations for Evidence-Based Policymaking Act of 2018 builds off the work of the U.S. Commission on Evidence-Based Policymaking to strengthen data privacy protections, improve secure access to data, and enhance the federal government's capacity for producing and using evidence.

Strengthens Privacy Protections

Maintains Strong Confidentiality Protections for Sensitive Data. Reauthorizes the Confidential Information Protection and Statistical Efficiency Act (CIPSEA), an existing law that gives the American public strong privacy safeguards and legal protections for appropriate uses of confidential data.

Institutes Processes to Assess Data Risks. Strengthens efforts to protect confidentiality while making data accessible for evidence building and transparent to the public by requiring comprehensive risk assessments for certain publicly released data.

Enhances Public Trust in Data. Improves public trust in statistical activities by codifying language directing certain agencies to establish procedures to protect trust in data activities by appropriately maintaining objectivity, independence,

Makes Administrative Records Available for Evidence Building. Under a strong set of confidentiality protections encourages that government data can and should generate evidence about policies and programs, and otherwise restricted by law.

Creates a Common Portal for Researcher Access to Restricted Data. Reduces burden on researchers applying to access government data by establishing a common application system for qualified individuals to access confidential data for approved projects.

Facilitates Continuous Feedback about Data Use. Promotes the use of data for evidence building by government advisory committee to review existing and availability of data.

Enhances Government's Evidence Building Capacity.

Data Science

and Engineering

The National Academies of
 SCIENCES • ENGINEERING • MEDICINE

CONSENSUS STUDY REPORT

DATA SCIENCE

FOR UNDERGRADUATES

Agency Actions

1. Identify Data Needs to Answer

	Foundations for Evidence-Based Policymaking Act of 2018 and Associated OMB Guidance	Executive Order on Maintaining American Leadership in Artificial Intelligence	Improving Implementation of the Information Quality Act (M-19-15)
1. Identify Data Needs to Answer	✓		✓
Governance Body	✓		✓
Infrastructure Maturity			✓
Increase Staff Data Skills	✓		
Support for Agency	✓	✓	✓
Inventory	✓		✓



Elena Semenova 9:09 PM

HI DOC data gurus! Do you know what the following indicates in reality? A person admitted first time in >= 2008 year with no previous incarcerations for lower offence class (1-3) being in jail for a few days but has sentence and custody dates goes back >=10 years. Does it mean that he/she was hiding from law enforcement all those years? How does custody date could go back like that in such situations? Is it just a bad data?



Vivek Ananda 11:27 PM

It mostly is bad data please email me the doc number so we can verify in the system

Learn more at bipartisanpolicy.org/evidence



ving tricky to fund and host.

Experience

- Biologists used the myGrid and Taverna workbench tools to conduct bioinformatics experiments and found that “manually, the processes undertaken by the workflows developed here could take at least 2 days, while the workflows achieve the same output in approximately an hour” (Stevens et al., 2004).
- A group at the Otto-von-Guericke-Universität Magdeburg working on Diffusion Tensor Imaging imaging for functional MRI (fMRI) used the Grid Workflow Execution System (GWES), “a grid implementation of the algorithm with slice based parallelization,” which, when applied to their image processing tasks, reduced “the processing down to 10% compared to a local cluster and 20% compared to sequential processing on the grid” (Solomonides, 2009).
- Another group, also in neuroimaging, used the SHIWA Platform to create “integrated workflow systems [that] have been shown to reduce experiment times from weeks to minutes” (Terstevanszky et al., 2014).

The screenshot shows the XSEDE User Portal website. The header includes the XSEDE logo and 'USER PORTAL' text, with a search bar and a 'SIGN IN' button. The main navigation bar contains links for MY XSEDE, RESOURCES, DOCUMENTATION, ALLOCATIONS, TRAINING, USER FORUMS, HELP, and ABOUT. A secondary navigation bar includes links for Get Started, Manage Data, User Guides, Community Codes, News, Usage Policy, Knowledge Base, and MFA. The main content area features the title 'Getting Started with XSEDE' with a sub-header 'Last update: October 31, 2019'. Below this is a section titled 'What is XSEDE?' which describes the XSEDE as a virtual collaboration funded by the National Science Foundation. A video player titled 'What Is XSEDE?' is visible at the bottom of the page.

NeuroImaging Tools & Resources Collaboratory (NITRC) Project: This resource for neuroinformatics software and data enables researchers to share software and data. Researchers using this resource saw a reduction in costs of 25%, and a reduction in time to process neuroimaging data by 85%. Users stated that having "access to several projects data in one central place.... saved us several weeks, and helped [them] select the best tools and conduct research of better quality."

SOC
ARXIV SocArXiv Papers

Enhancing and accelerating social science via automation.
Challenges and opportunities

AUTHORS

Tal Yarkoni, Dean Eckles, James Heathers, Margaret Levenstein, Paul Smaldino, Julia Lane

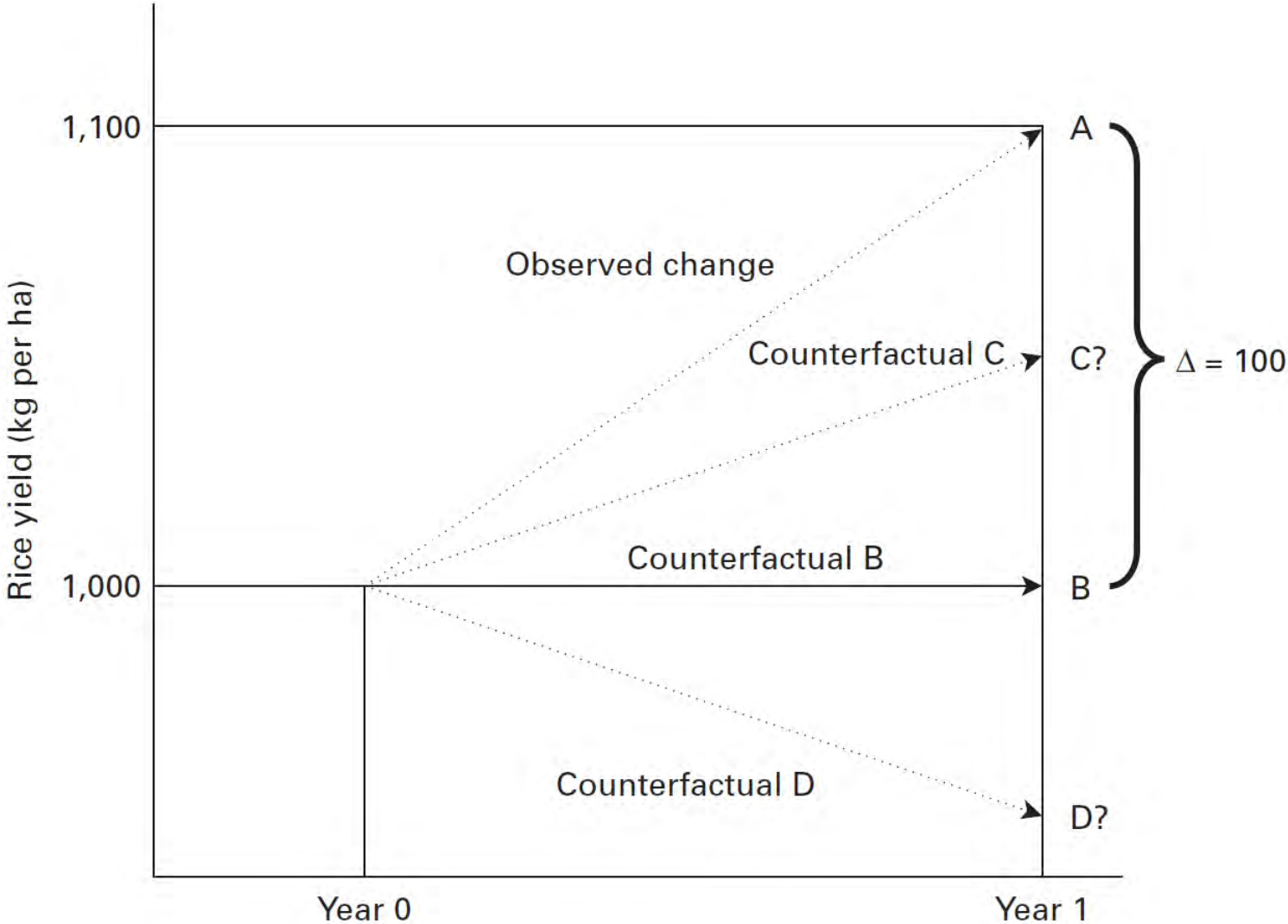
SUBMITTED ON

January 24, 2019

LAST EDITED

January 25, 2019

Figure 3.3 Before-and-After Estimates of a Microfinance Program



Note: Δ = Change in rice yield (kg); ha = hectares; kg = kilograms.

Figure

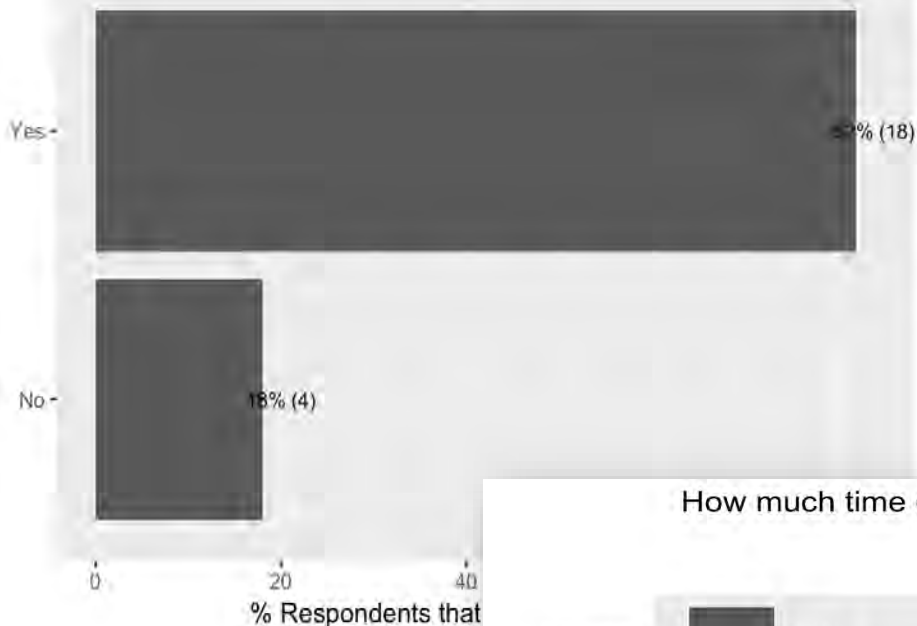
INP

Final
huma
other re
mobili
sup
activ

Bu
sta
other
res

Do you seek out experts on a dataset to help you learn about the data?

22 total respondents

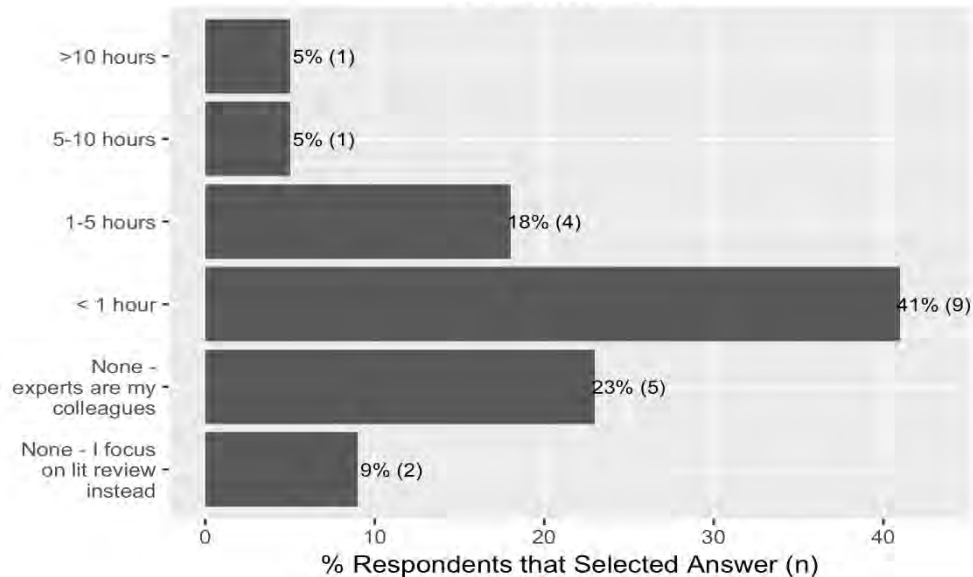


Source: Coleridge

state?

How much time do you spend searching for experts to consult with?

22 total respondents

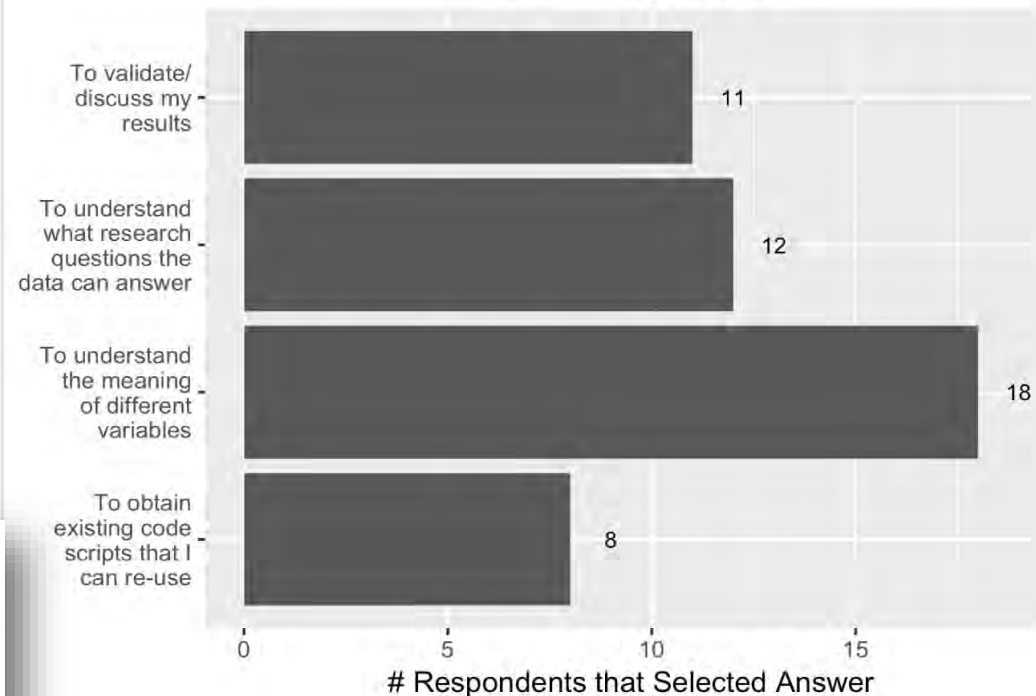


Source: Coleridge Initiative, ADA Training Program Participants

What motivates your search for experts?

(Respondents can select multiple answers)

21 total respondents



Source: Coleridge Initiative, ADA Training Program Participants

Rich Context

The goal of the Rich Context project is to create a new platform that enables empirical analysts to search for and discover datasets.

The **challenge** that empirical researchers face is that, for a given dataset, it is difficult to find out **who** has worked with the data before, **what methods and code** were used, and **what results were produced**.

Leaderboard competition

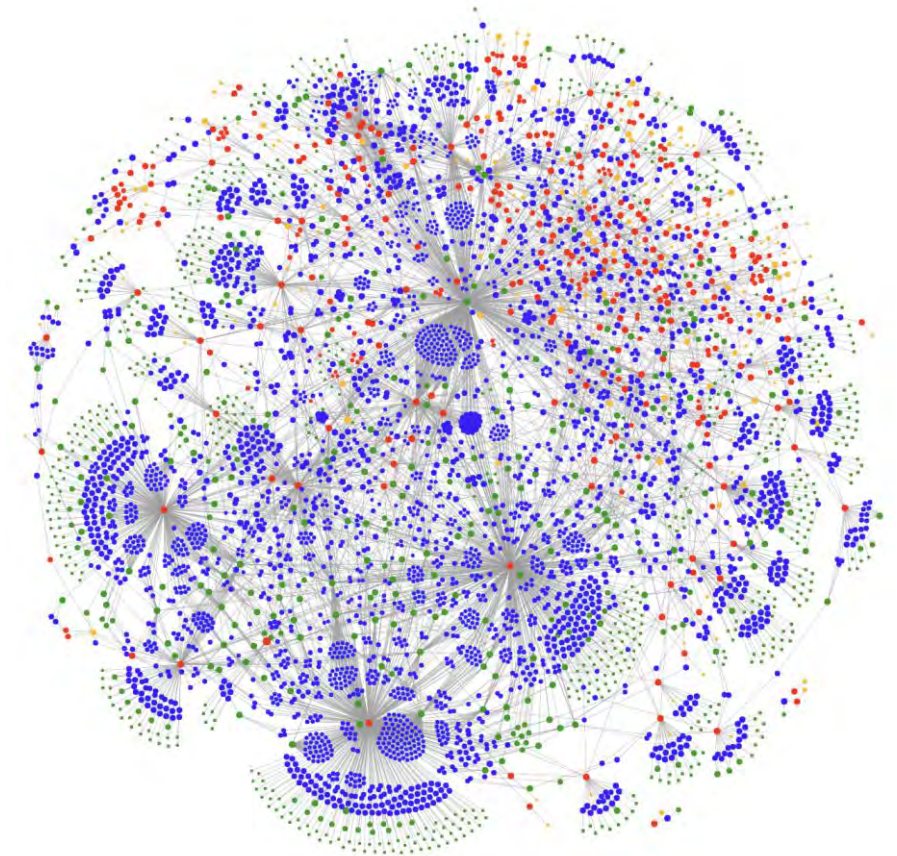
github.com/Coleridge-Initiative/rclc

See also:

[“Human-in-the-loop AI for scholarly infrastructure”](#)

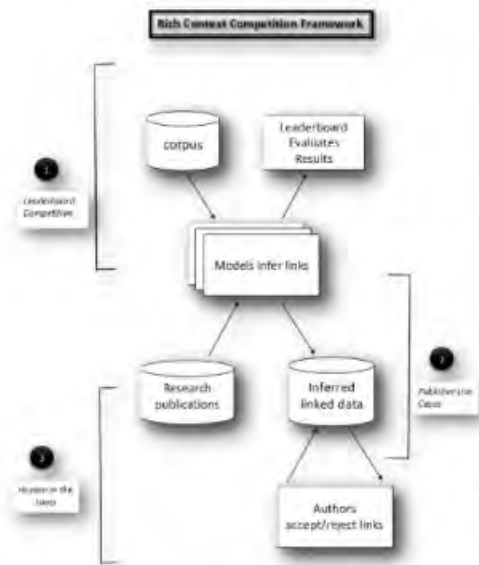
[“New initiative to help with discovery of dataset use in scholarly work”](#), Christian

Zimmerman



- providers
- datasets
- journals
- papers

1. A Python library [richcontext-scholapi](#) which provides API integrations for federating exchange across multiple scholarly infrastructure providers. Additional APIs are cons
2. A **Knowledge Graph** of known links between datasets, research publications, resear sciences.



3. A **leaderboard competition** in models to automatically detect publications.

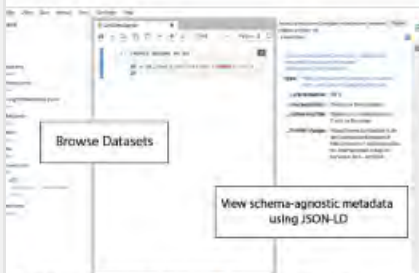
Model development relies on research publications and the The corpus continues to expand which are validated by domain

Teams use a subset of the cor models. Using another subset and improve their models. Our performance; scores are update model versions. Read more at [Competition Wiki](#).

Teams are also contributing to extraction. The competition is

4. A set of Jupyter extensions w Features like the **Data Registry** **Commenting** enable users to browse metadata on their data collaborate in real-time on sha goal is to enrich the knowledg comments, collaboration and c

and developing a recommendation engine for datasets and code.



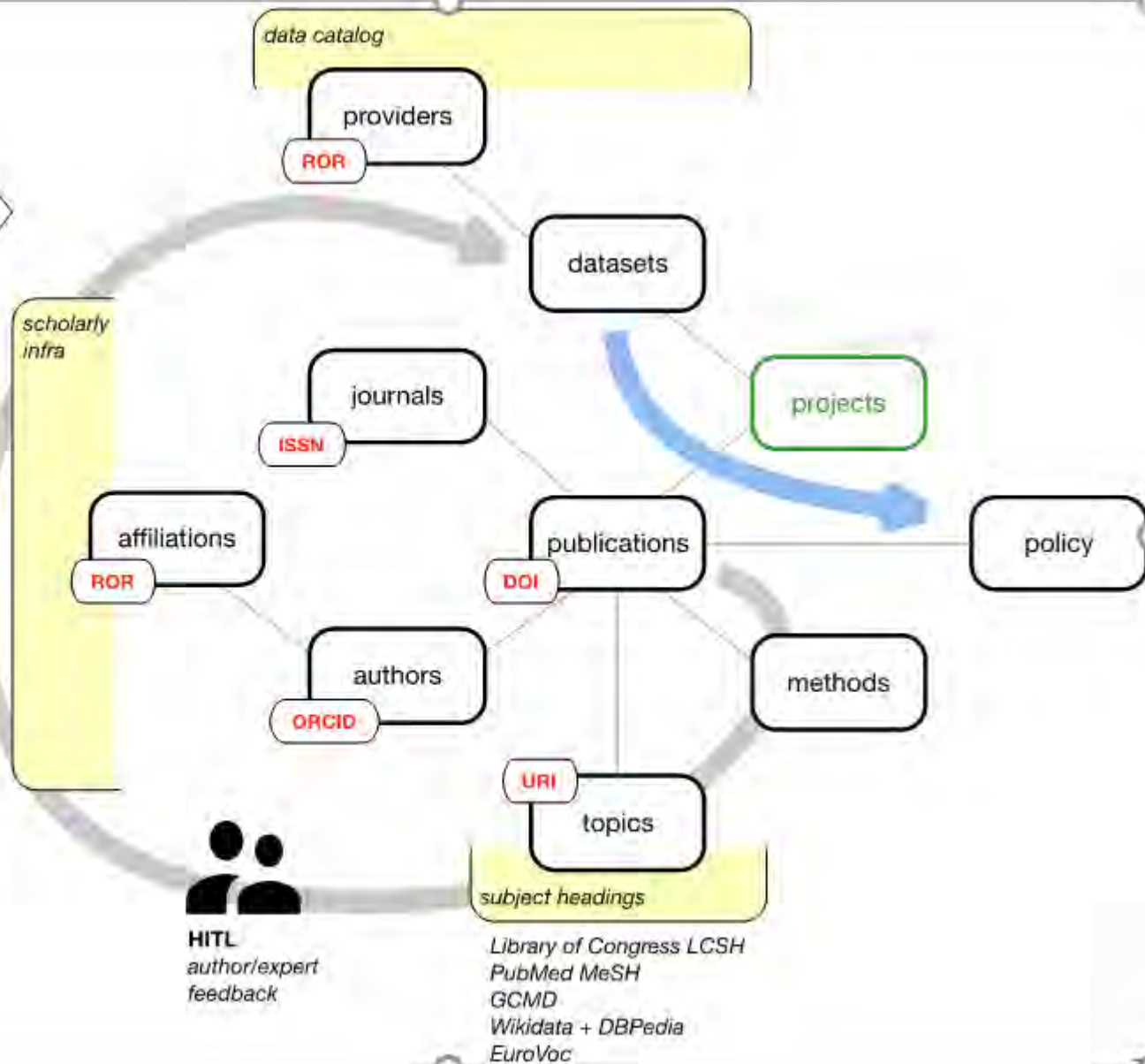
Registry (left) and Metadata Browser (right) in terLab. *Project Jupyter*

Demo - Linked Data generated by user appears in Jupyter Data Registry. *Project Jupyter*

Discovery Services
 Unpaywall
 Dimensions
 RePEc
 ResearchGate
 Crossref
 DataCite
 ORCID
 OpenAIRE
 PubMed
 EuropePMC
 Semantic Scholar
 ScholeXplorer
 dissemin
 Elsevier
 SSRN
 etc.

federated queries

metadata updates



Result

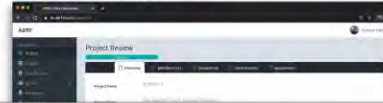
Rich Context in Secure Facility

Managing, Monitoring, and Measuring the way datasets are used in research projects.

MANAGING

Research Project Request Workflow

Researchers can request access to datasets via the project request workflow. The appropriate Data Stewards are notified when their datasets are requested. They can then review and approve access.



MONITORING

Projects and Dataset Access

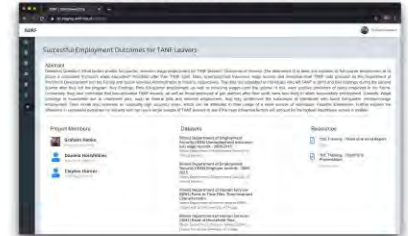
The Data Stewardship Dashboard gives Data Stewards a high-level overview as to which of their datasets researchers have access to as a part of their projects.



MEASURING

Research Papers and Presentations

Once a research project has finished, data products such as papers or presentations can be uploaded and associated with the datasets used, providing context and knowledge sharing with future researchers.



Rich Context as an open platform

The screenshot shows a web browser at rc.coleridgeinitiative.org/?radius=2&entity=IRI%20Infoscans. The interface features a search bar with 'IRI Infoscans' and a 'radius: 2' filter. On the left, there are several filter categories: datasets (Food Security Survey Module, IRI Infoscans, IRI Consumer Network), providers (IRI), publications (WIC and the Retail Price of Infant Formula, Spatial and Temporal Variation in the Value of the Women, In..., Tobacco Marketing at SNAP- and WIC-Authorized Retail Food St..., Household food security: Perceptions, behavior and nutrition..., One's Nutrition and Food Security in Food Deserts?), authors (Gundersen, Craig, Laska, Melissa Nelson, Powell, Lisa M., Yen, Steven T., Smallwood, David M., Wu, Ping-Chao, Zhen, Cheng), journals (Food Policy, Nutr. Rev., J Health Econ, Am J Agric Econ, Prev Med Rep, Health Econ, SSRN Electronic Journal), and topics (children, food insecurity, women, wic, food, food security, special supplemental nutrition program). The main area displays a network graph with nodes of various colors (blue, purple, red, green) and sizes, connected by lines. The footer contains: ©2020, Coleridge Initiative • Feedback • Research • Configure • OpenAPI • ML Competition.

It also provide high-level metrics on which and by whom so that Data Providers can get the most value.



For example, it's a completed project.

Export Metrics

Dataset usage information can be exported from the application so that Data Stewards can generate reports to share amongst key stakeholders at their agency.

Dataset usage metrics can be exported or generated from the tool and used.

<https://coleridgeinitiative.org>

Contact and more information

- Website: <https://coleridgeinitiative.org/>
- White paper: <https://tinyurl.com/rcwp2019>
- Book: <http://34.82.145.119:8080/>
- Email
 - dataanalytics@coleridgeinitiative.org
 - julia.lane@nyu.edu
- GitHub organization: <https://github.com/Coleridge-Initiative>



BILL & MELINDA
GATES foundation



Alfred P. Sloan
FOUNDATION

