# Rich Context:
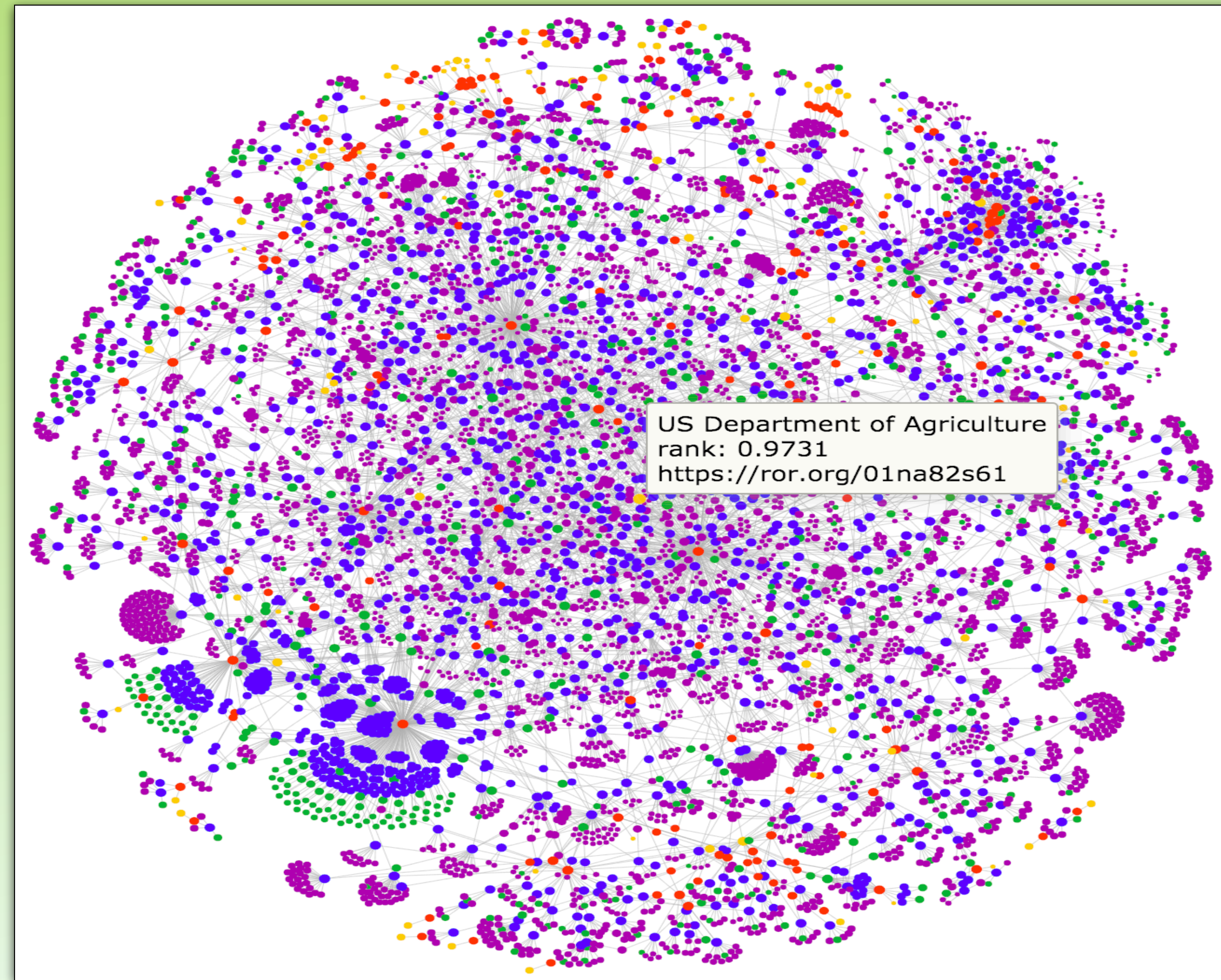## Rich Search and Discovery for Scholarly Datasets

Julia Lane

Paco Nathan

NYU Coleridge Initiative



US Department of Agriculture
rank: 0.9731
https://ror.org/01na82s61

# Context

# Context

# Context

# Evidence-Based Policymaking

- **Foundations for Evidence-based Policymaking Act** (2018)

- **Information Quality Act** (2001)

- **US Federal Data Strategy** (2019)

- **Year-1 Action Plan** (2019)

**2018 Government Innovation Awards**

# Data mashups at government scale

BY GCN STAFF | NOV 01, 2018

The Census Bureau — the government's original data agency — collaborated with the University of Chicago, University of Maryland and New York University to provide a secure cloud-based platform that allows government employees and academic researchers to take advantage of advanced data science tools.

The goal of the Administrative Data Research Facility is to give authorized users access to a secure supercomputer and thousands of datasets. It has received a Federal Risk and Authorization Management Program moderate certification.

John Abowd, an associate director and chief scientist at the bureau, said it is the first such platform in the federal

**Administrative Data Research Facility**

Census Bureau

See **article link**

# Approach: Technical and Human

**Technical**

- Create secure environment where data providers can share their data across agency and jurisdictional lines
- Census and USDA Authorization to Operate; HHS in process

**Operational**

- Link disparate data
- Analyze data

**Legal & Practical**

- Document value associated with the data linkage
  + Consistent with the agency mission
  + Useful enough to engage decision-makers

# Need: Recommenders for researchers/analysts

- Objective: provide better means of *search and discovery* for social science researchers and agency analysts.

- Collect workflow telemetry and query logs to augment the graph.

- Currently developing recommender systems based on the graph.

- This accelerates research and also assists training (e.g., onboarding agency analysts).

- Near-term goal: identify people with specific expertise.

- Long-term goal: learn workflow configurations to support AutoML meta-learning.

# Rich Context

- Focus on *socioeconomic impact*

- Funded by Schmidt Futures, Sloan, Overdeck

- Partnering with Bundesbank, USDA, etc.

- Collaboration with SAGE Pub, RePEc, ResearchGate, Digital Science, etc.



- providers
- datasets
- journals
- papers

# Rich Context

- Focus on *socioeconomic impact*

- Funded by Schmidt Futures, Sloan, Overde...

- Partnering with... USDA, etc.

- Collaboration w... Pub, RePEc, R... Digital Science, etc.

Challenges that empirical researchers face: for a given dataset, find out **who** has worked with the data before, **what methods and code** were used, and **what results** were produced.

providers
datasets
journals
papers

# Knowledge Graph – why?

- Allow flexibility for metadata representation

- Measure metadata quality

- Prepare features for ML models

- Build recommenders for *experts*, *topics*, *tools*, etc.

- Engage the public with automated data inventories

- Recommend configurations to new analysts

- Identify which datasets get used with others

- Quantify impact of datasets on policy

# Knowledge Graph – how?

- Manual data entry and curation of linked data

- Use persistent identifiers whenever possible:
  DOI, ISSN, ROR, ORCID, etc.

- Leverage ML models to infer missing metadata

- Federate queries of discovery services APIs

- Suggest corrections for metadata errors

- Use HITL to build feedback loops that engage experts,
  and provide convenient means for manual override

- Identify errors by using unit tests, ontology axioms,
  graph analytics, etc.

- **Collaborate with agency libraries!**

# KG process

- who are the expert people?
- which topics are emerging?
- how can methods be shared?

**activities** → **outputs** → **outcomes** → **impact**

curated datasets     research projects     published research     better science, government, education

*ML models infer new metadata links*

*how do we track the linkage??*

*how do we measure these behaviors??*

# KG construction and representation

```
providers
        \
         datasets
                  \
                   manually curate
                   or
journals           infer by ML models
        \          /
         publications
```

# KG construction and representation

# KG construction and representation



See github.com/Coleridge-Initiative/rclc/wiki/Corpus-Description

# KG construction and representation



federated queries

data catalog

providers
**ROR**

datasets

manually curate
or
infer by ML models

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

*scholarly
infra*

journals
**ISSN**

affiliations
**ROR**

publications
**DOI**

policy

authors
**ORCID**

methods

**URI**

topics

*subject headings*

*Library of Congress LCSH
PubMed MeSH
Wikidata + DBPedia
EuroVoc*

*metadata
updates*

# KG construction and representation



*federated queries*

*metadata updates*

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

*scholarly infra*

*data catalog*

providers
**ROR**

datasets

journals
**ISSN**

projects

affiliations
**ROR**

publications
**DOI**

policy

authors
**ORCID**

methods

**URI**
topics

*subject headings*

*Library of Congress LCSH
PubMed MeSH
Wikidata + DBPedia
EuroVoc*

**HITL**
*author/expert feedback*

# KG construction and representation

# Open Source Projects

- **RCGraph** – Rich Context knowledge graph management
  **github.com/Coleridge-Initiative/RCGraph**

- **richcontext.scholapi** – federated discovery services and
  metadata exchange across scholarly infrastructure APIs
  **pypi.org/project/richcontext-scholapi**

- **adrf-onto** – controlled vocabulary for ADRF and Rich Context
  using OWL, SKOS, DCAT, PAV, CITO, FaBiO, etc.
  **github.com/Coleridge-Initiative/adrf-onto**

- **RCLC** – ML leaderboard competition
  **github.com/Coleridge-Initiative/rclc**

See also:

**"Machine Learning Highlights for Rich Context"**

# Funded additions to Project Jupyter

Make datasets and projects top-level constructs, support metadata exchange and privacy-preserving telemetry from notebook usage:

- JupyterLab **Commenting** and real-time collab similar to Google Docs

- JupyterLab **Data Explorer**: register datasets within research projects

- JupyterLab **Metadata Explorer**: browse metadata descriptions, get recommendations through knowledge graph inference (via extension)

- **Data Registry** (original proposal)

- **Telemetry** (privacy-preserving, reports usage)

**Saul Shanabrook**
@SShanabrook

Replying to @choldgraf

Data Catalog Vocabulary and other related vocabularies are useful here. w3.org/TR/vocab-dcat/ We are building a way to explore metadata defined in JSON LD that uses these in JupyterLab github.com/jupyterlab/jup... cc @pacoid



12:46 AM · Oct 11, 2019 · Twitter Web App

# ML Leaderboard Competition

**github.com/Coleridge-Initiative/rclc**

- update from RCC competition in 2018-2019
- ongoing ML leaderboards (similar to **NLP-progress**)
- open source, hosted on GitHub
- highly curated test sets, all open-access publications
- teams collaborate via GH issues on corpus data quality, etc.
- focus on *precision* for ML model evaluation

**Current SOTA**

| source | precision | entry | code | paper | corpus | submitted | notes |
|---|---|---|---|---|---|---|---|
| LARC @philipskokoh | 0.7836 | ipynb | repo | RCC_1 | v0.1.5 | 2019-09-26 | RCLC baseline experiment using RCC_1 approach |
| KAIST @HaritzPuerto | 0.6319 | ipynb | repo | RCC_1 | v0.1.5 | 2019-11-01 | model trained a different dataset using DocumentQA and Ultra-Fine Entity Typing -- NB: this approach is able to identify new datasets |

# Human-in-the-loop

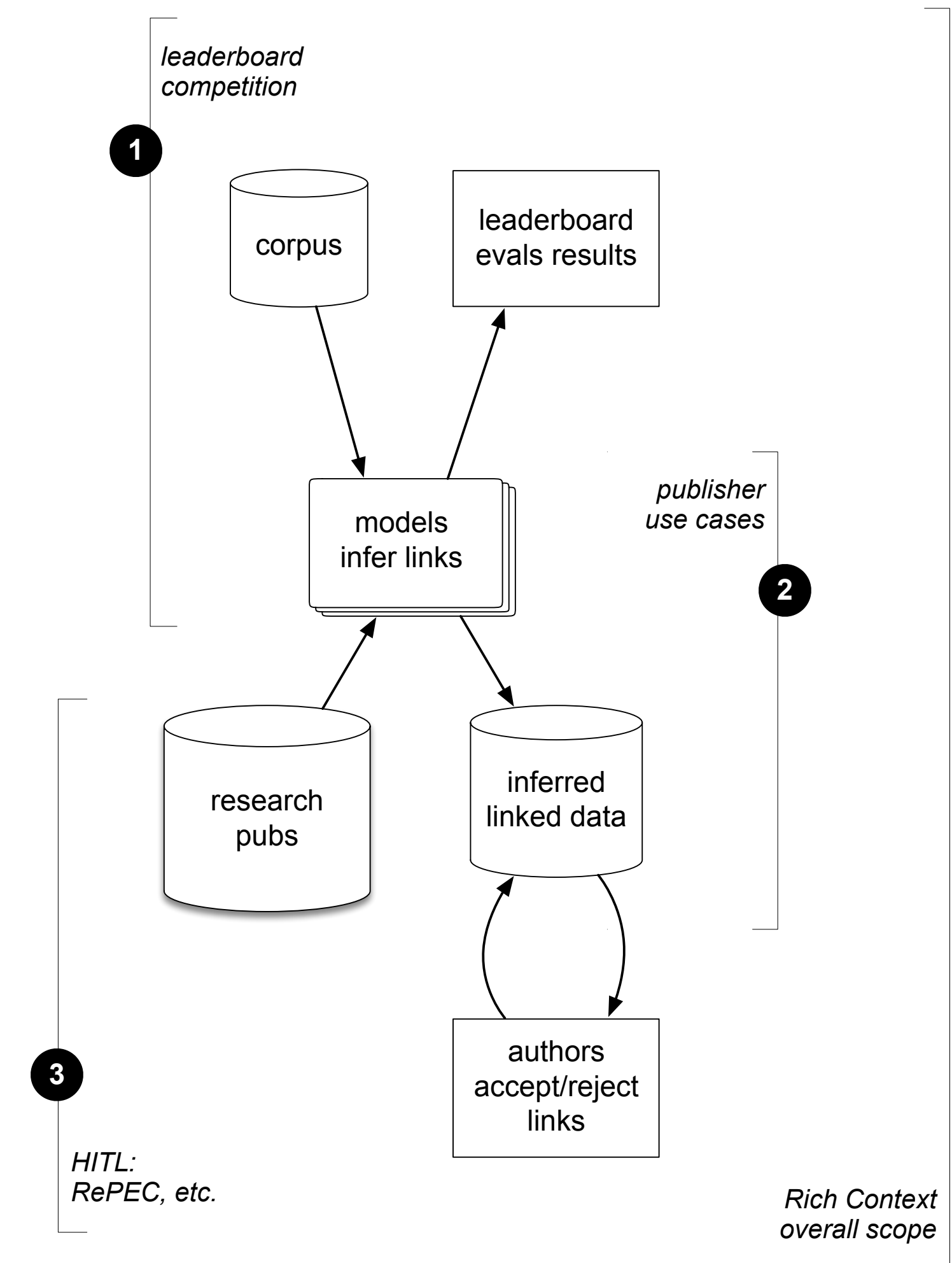- semi-supervised learning, aka "human-in-the-loop" – in progress via RePEc

- interact with authors to confirm metadata inferred by ML models

- feedback from experts improves the corpus metadata and the ML modeling

See also:
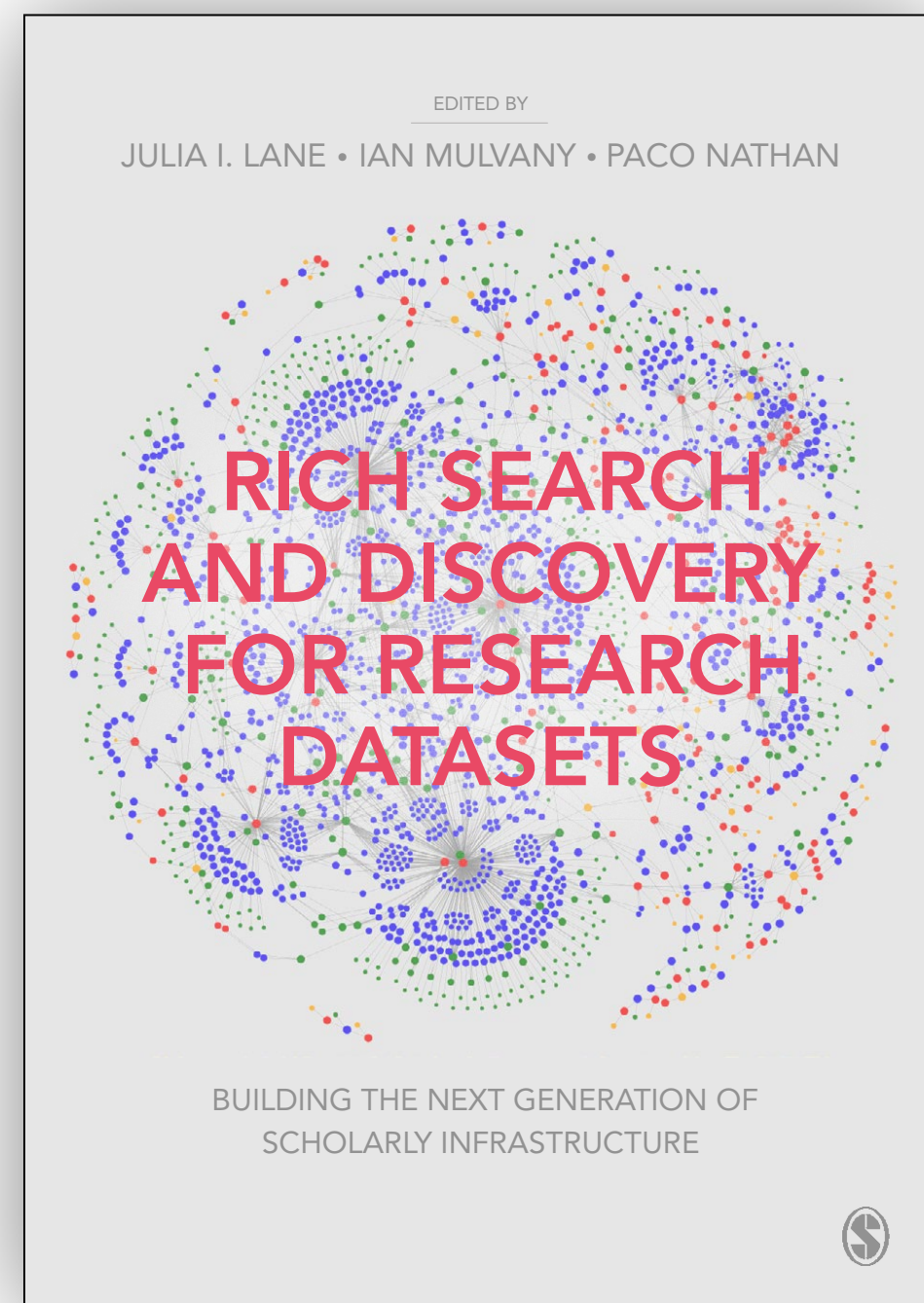
**"Human-in-the-loop AI for scholarly infrastructure"**

**"New initiative to help with discovery of dataset use in scholarly work"**, Christian Zimmerman

# Additional Information

Rich Context @ NYU Coleridge Initiative
**coleridgeinitiative.org/richcontext**

- **white paper**

- **upcoming book** (Jan 2020)

- **feedback/propose collaboration**



EDITED BY

JULIA I. LANE · IAN MULVANY · PACO NATHAN

RICH SEARCH
AND DISCOVERY
FOR RESEARCH
DATASETS

BUILDING THE NEXT GENERATION OF
SCHOLARLY INFRASTRUCTURE

SAGE

**"Empty rhetoric over data sharing slows science"**
*Nature* (2017-06-12)

**"Experiences of the Deutsche Bundesbank"**
Stefan Bender
*CEMLA* (2019-05-28)

**"Where's Waldo: Finding datasets in empirical research publications"**
Julia Lane
*AKBC* (2019-05-22)

**"Google data set search"**
Ian Mulvany
*ScholCommsBlog* (2019-11-19)

**"Impact for social science researchers"**
Ian Mulvany
*FORCE11* (2019-11-17)