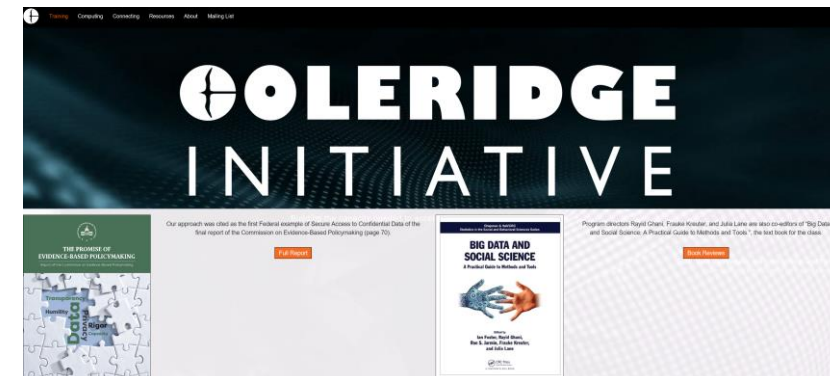


# Administrative Data Research Facility and Metadata

Julia Lane

New York University



# Key challenges to be solved with metadata – particularly for federal statistical system

- Limited internal capacity
- Security
- Legal mandates surrounding access and use
- Data sharing issues
  - cost
  - burden
  - data quality
  - data documentation
  - risk of bad analysis



## H.R. 1831: Evidence-Based Policymaking Commission Act of 2016

Introduced: **Apr 16, 2015**

114<sup>th</sup> Congress, 2015–2017

Status: **Enacted — Signed by the President on Mar 30, 2016**

This bill was enacted after being signed by the President on March 30, 2016.

Law: Pub.L. 114-140

Sponsor:



**Paul Ryan**

Representative for Wisconsin's 1st congressional district  
Republican

Text:

[Read Text »](#)

Last Updated: Mar 18, 2016

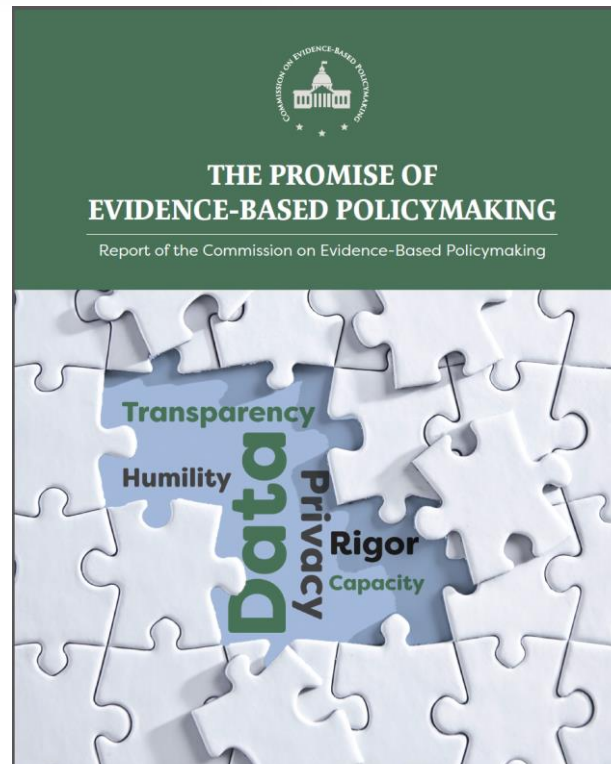
Length: 5 pages



# Context

## FY 2016 Significant Investments

- **2020 Census (\$663M):** We have the potential to save \$5 billion with the new 2020 Census design, however, we now have to build operations and systems for the 2020 Census, based on the new design.
- **CEDCaP (\$78M):** Smarter-IT Delivery Built on a Shared-Services Model.
- **American Community Survey (\$257M):** We must maintain the quality of the data while continuing our efforts to reduce respondent burden.
- **Geographic Support (\$81M):** We must make use of technology and partnerships to deliver smarter geographic solutions to our surveys and censuses.
- **Administrative Records Clearinghouse (\$10M):** Will expedite the acquisition of federal and federally sponsored administrative data sources, improve data documentation and linkage techniques, and leverage and extend existing systems for governance, privacy protection, and secure access to these data.
- **Economic & Government Censuses (\$144M):** Data products drive economic activity and are relevant to the needs businesses, policymakers, and the public. \$10.1 million increase



**Administrative Data Research Facility:** The Administrative Data Research Facility is a pilot project that enables secure access to analytical tools, data storage and discovery services, and general computing resources for users, including Federal, state, and local government analysts and academic researchers. The Census Bureau and academic partners developed the project as part of the collaborative Training Program in Applied Data Analytics sponsored by the University of Chicago, New York University, and the University of Maryland.<sup>1</sup> It is currently operating as a pilot with users accessing the Facility as part of the training program. The Facility operates as a cloud-based computing environment, with Federal security approvals, which currently hosts selected confidential data from the U.S. Department of Housing and Urban Development and the Census Bureau, as well as state, city, and county agencies, and an

# Resources

## Companion websites for publications

- ▶ [Seeing Sound: Investigating the Effects of Visualizations and Complexity on Crowdsourced Audio Annotations](#)

## Data

- ▶ [Urbansound Dataset](#) – A dataset containing 1302 labeled sound recordings. Each recording is labeled with the start and end times of sound events from 10 classes
- ▶ [Urbansound8k Dataset](#) – A dataset containing 8732 labeled sound excerpts ( $\leq 4$ s) of urban sounds from 10 classes
- ▶ [URBAN-SED Dataset](#) – A dataset of 10,000 synthesized soundscapes with sound event annotations generated using [Scaper](#)
- ▶ [Seeing Sound Dataset](#) – A dataset of 5400 crowdsourced audio annotations of 60 synthesized soundscapes

## Code

- ▶ [Scaper](#) – A Python library for soundscape synthesis and augmentation
- ▶ [Audio-Annotator](#) – A Javascript web interface for annotating audio data
- ▶ [Raster Join](#)
- ▶ [Urban Pulse](#)

te  
e Synthesis  
is in Smart  
s join  
Croix, AFP,

UI-27

MDES

QU

REF

READ AC...

TIONS BE...

THIS PAG...

TO IDEN...

EVEN IF Y...

ON AN AL...

IF ADI...

NEEDE...

RETURN...

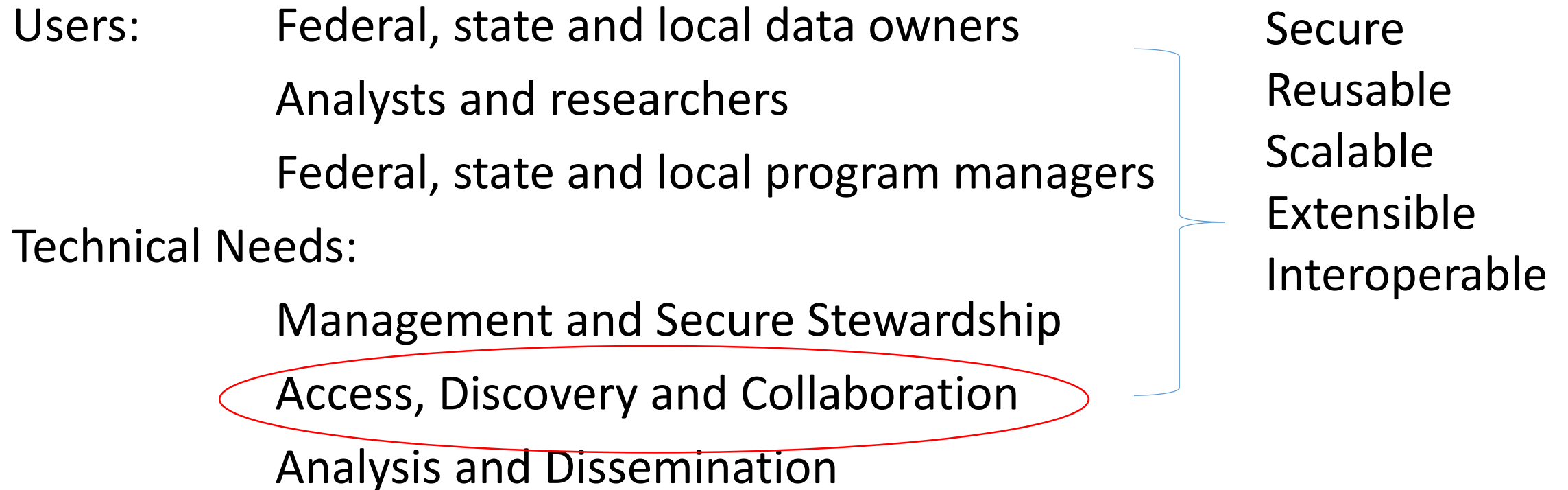
thous...

at the

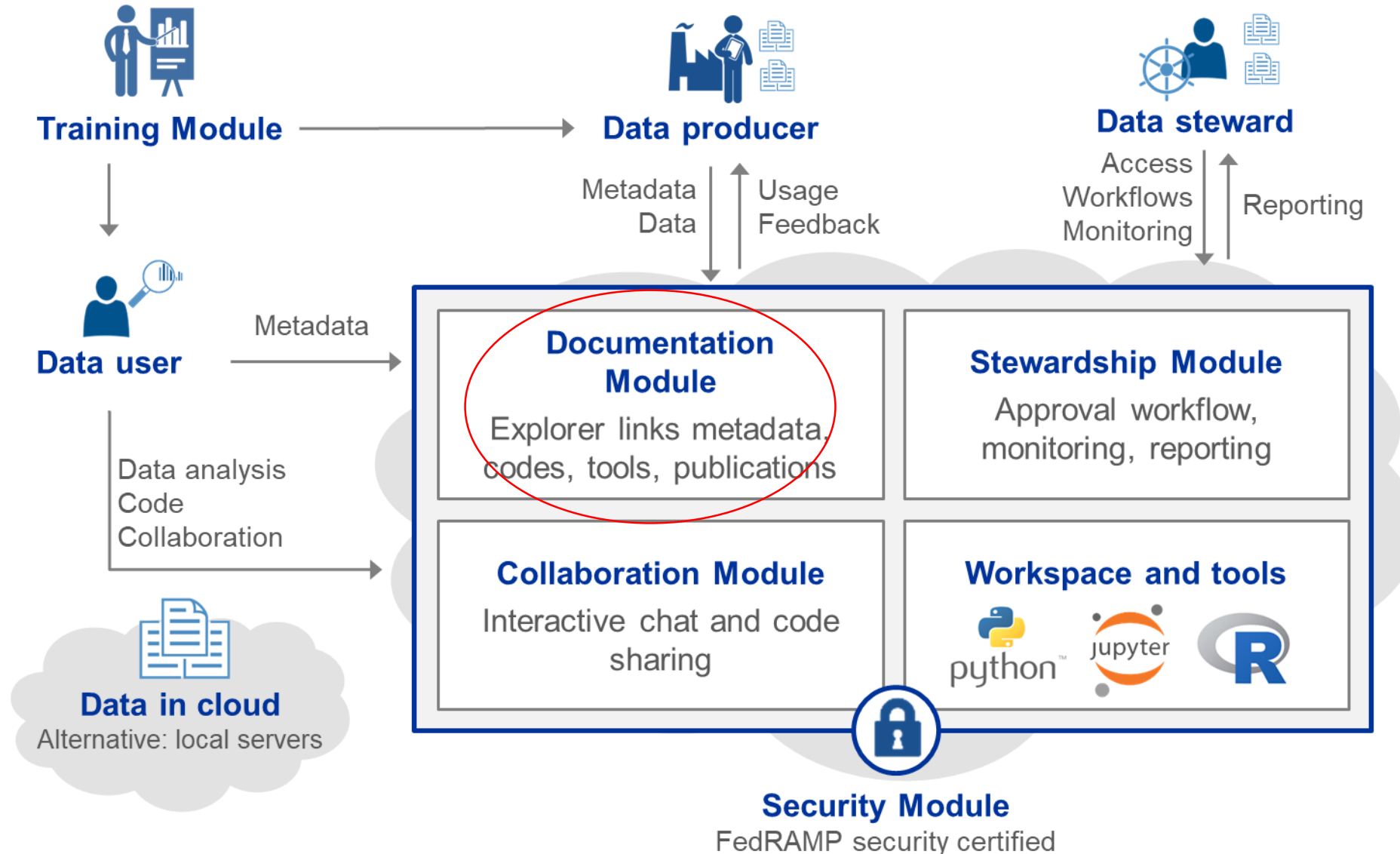
If you need help about completin 1-800-354-7271



# Build technical environment



# Functional characteristics





# Inspiration

## Carole Goble

From Wikipedia, the free encyclopedia

**Carole Anne Goble**, *CBE FRSng* (born 10 April 1961) is a British academic who is Professor of *Computer Science* at the University of Manchester.<sup>[14][15]</sup> She is Principal Investigator (PI) of the *myGrid*,<sup>[16]</sup> *BioCatalogue*<sup>[17]</sup> and *myExperiment*<sup>[18]</sup> projects and co-leads the Information Management Group (IMG) with Norman Paton.<sup>[19][20]</sup>

<b>Contents</b> <span>[hide]</span>
1 Education
2 Research
3 Career
4 Awards and honours
5 References

### Education [edit]

Goble was educated at Maidstone Grammar School for Girls.<sup>[1]</sup> Her academic career has been spent at the School of Computer Science where she gained her Bachelor of Science degree in *computing and information systems* from 1979<sup>[2]</sup> to 1982.

### Research [edit]

Her current research interests<sup>[11][22]</sup> include *Grid computing*, the *Semantic Grid*,<sup>[23]</sup> the *Semantic Web*, *Ontologies*,<sup>[24][25][26]</sup> *e-Science*, *medical informatics*,<sup>[27]</sup> *Bioinformatics*, and *Research Objects*. She applies advances in knowledge technologies and workflow systems<sup>[28]</sup> to solve information management problems for life scientists and other scientific disciplines<sup>[citation needed]</sup>. She has successfully secured funding from the European Union, the Defense Advanced Research Projects Agency (DARPA) in the US and UK funding agencies including the Engineering and Physical Sciences Research Council (EPSRC),<sup>[29]</sup> Biotechnology and Biological Sciences Research Council (BBSRC),<sup>[30]</sup> Economic and Social Research Council (ESRC), Medical Research Council (MRC), the Department of Health, The Open Middleware Infrastructure Institute and the Department of Trade and Industry.<sup>[31]</sup>

Her work has been published in leading peer reviewed scientific journals including *Nucleic Acids Research*,<sup>[3]</sup> *Bioinformatics*,<sup>[32][33]</sup> *IEEE Computer*,<sup>[10]</sup> the *Journal of Biomedical Semantics*,<sup>[34]</sup> *Briefings in Bioinformatics*,<sup>[35][36][37]</sup> *Artificial Intelligence in Medicine*,<sup>[37]</sup> the Pacific Symposium on Biocomputing conference,<sup>[24]</sup> the *International Journal of Cooperative Information Systems*, the *Journal of Biomedical Informatics*,<sup>[38]</sup> *Nature Genetics*<sup>[39]</sup> and *Drug Discovery Today*.<sup>[40][41][42][43][44][45]</sup>

### Career [edit]

**Carole Goble**



Carole Goble by Rob Whitrow

**Born** Carole Anne Goble  
10 April 1961 (age 57)<sup>[1]</sup>

**Nationality** United Kingdom

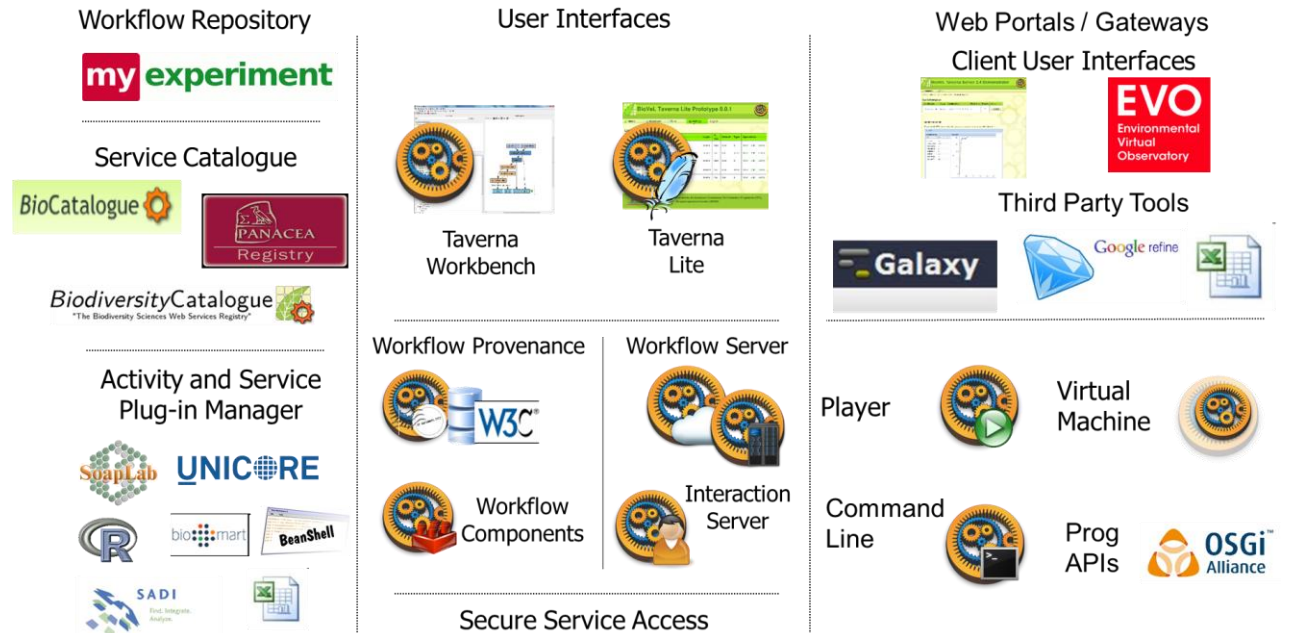
**Alma mater** University of Manchester

**Known for** myGrid  
Semantic Grid  
Open PHACTS<sup>[2]</sup>  
Taverna workbench<sup>[3][4]</sup>  
Software Sustainability Institute  
The Seven Deadly Sins of Bioinformatics<sup>[5]</sup>

**Spouse(s)** Ian Cottam (m. 2003)<sup>[6]</sup>

**Awards** Jim Gray e-Science Award (2008)

# The Taverna Suite of Tools



# USER CONTRIBUTIONS

Overview Repositories 28 Stars 3 Followers 43 Following 0

**Popular repositories**

- vscode-powershell**  
Forked from PowerShell/vscode-powershell  
Provides PowerShell language and debugging support for Visual Studio Code  
TypeScript ★ 1
- try\_git**
- gittest**  
gittest  
JavaScript
- test2**  
C++
- mvc**  
mvc  
JavaScript 1
- express**  
express starter  
JavaScript 2

1,701 contributions in the last year

Contribution activity

Jump to 2017

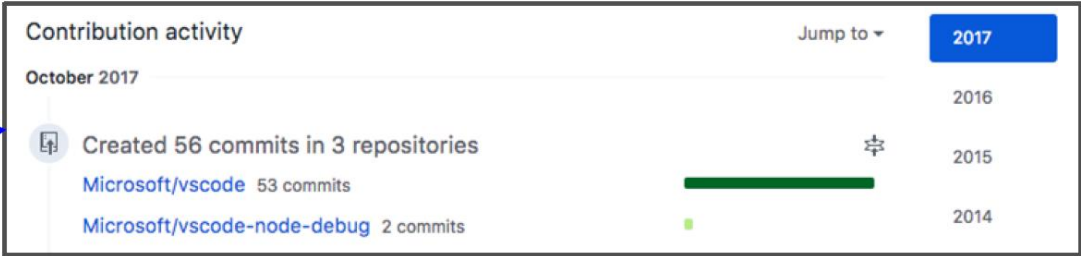
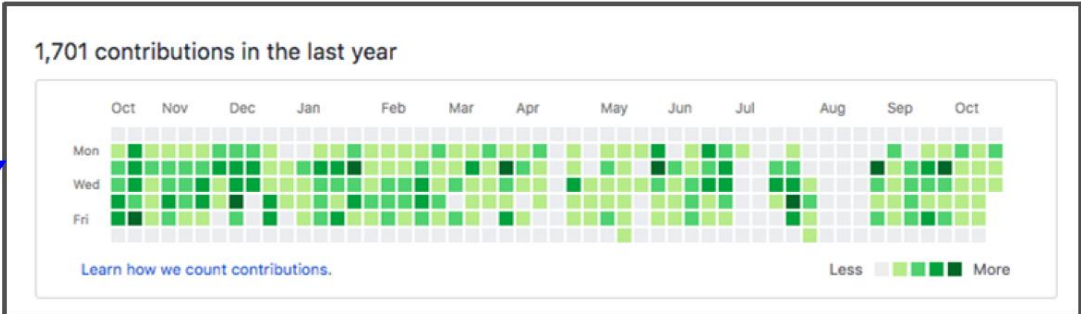
October 2017

- Created 56 commits in 3 repositories
- Microsoft/vscode 53 commits
- Microsoft/vscode-node-debug 2 commits
- Microsoft/vscode-generator-code 1 commit

Repositories 28 Stars 3 Followers 43 Following 0

**Popular repositories**

- vscode-powershell**  
Forked from PowerShell/vscode-powershell  
Provides PowerShell language and debugging support for Visual Studio Code  
TypeScript ★ 1
- try\_git**





# CONTRIBUTION TRACKING



tjac  
Lincolnshire,  
United  
Kingdom  
✍️ 11 🗳️ 5

🟢🟢🟢🟢🟢 Reviewed yesterday

## Colourful Tamassa

I cannot praise the staff at this hotel enough. They are all wonderful,

he  
in  
ar

### Lindsay H

Level **3** Contributor

TripAdvisor member since 2007

✍️ 6 Contributions

🌐 13 Cities visited

🗳️ 6 Helpful votes

📷 23 Photos

#### REVIEW DISTRIBUTION

Excellent	<div style="width: 83%;"></div>	5
Very good	<div style="width: 17%;"></div>	1
Average	<div style="width: 0%;"></div>	0
Poor	<div style="width: 0%;"></div>	0
Terrible	<div style="width: 0%;"></div>	0



Lindsay H  
✍️ 6 🗳️ 6



Lindsay H

✉️ Send Message

📅 Aug 2007  
4 year old female

6 Reviews

1 Rating

48 Photos

6 Helpful votes

### Lindsay H's TripCollective Progress

Total Points

2,051

Level **3** Contributor

### Badges (16 total)

[View Collection](#)



Passport

6 Cities



Senior Photographer

20 Photos



Readership

1,000 Readers



Helpful Reviewer

5 Votes



Senior Reviewer

5 Reviews

Total Points

0

Current Level



Next Level

1

300 points to go

# Making Computational Research with Sensitive Data Possible and Valuable

Brian E. Granger  
Associate Professor  
Cal Poly

Julia Lane  
Professor  
NYU

Fernando Perez  
Assistant Professor  
UC Berkeley

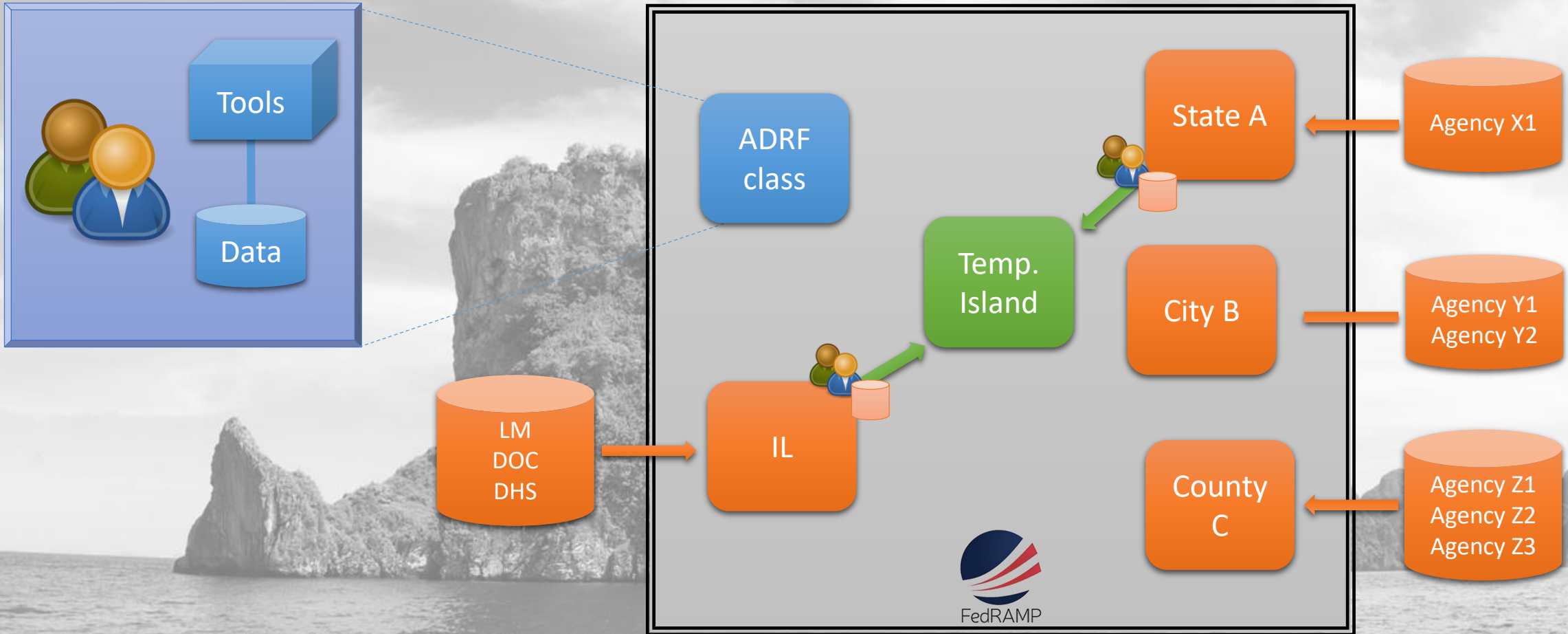


Alfred P. Sloan  
FOUNDATION

SCHMIDT **FUTURES**



Overdeck Family  
Foundation



ADRF SaaS

## Data

Data on Individuals

Data on Organizations

Data on Places

## Training

Joined up datasets in  
secure environment with  
collaborative tools

Applied Data Analytics  
around core questions

## Results


Trained Staff

New Products

New networks

New metrics

← → ↻ ↶ 🔒 Gi

 jupyter

---

What is .

JupyterHub brings data scientists - ca

JupyterHub runs in and large-scale inf

## Key feat

**Customizable** - Jupyter and more.

**Flexible** - Jupyter-

**Scalable** - Jupyterl

**Portable** - Jupyterl

```
from sklearn.naive_bayes import GaussianNB
from sklearn. import DecisionTreeClassifier
from sqlalchemy import create_engine
#import pydot
sns.set_style('white')
sns.set_context('poster', font_scale=1.25, rc={"lines.linewidth":1.25, "lines.markersize":8})
```

### Connect to Database

```
In [ ]: db_name = "appliedda"
hostname = "10.10.2.10"
conn = psycopg2.connect(database=db_name, host = hostname) #database connection
```

The database connection allows us to make queries to a database from Python.

```
In [ ]: df_tables = pd.read_sql("""SELECT * FROM ides.il_wage limit 10;""", conn)
```

```
In [ ]: df_tables.head()
```

## The Machine Learning Process

[Go back to Table of Contents](#)

- **Understand the problem and goal.** *This sounds obvious but is often nontrivial.* Problems typically start as vague descriptions of a goal - improving health outcomes, increasing graduation rates, understanding the effect of a variable  $X$  on an outcome  $Y$ , etc. It is really important to work with people who understand the domain being studied to dig deeper and define the problem more concretely. What is the analytical formulation of the metric that you are trying to optimize?
- **Formulate it as a machine learning problem.** Is it a classification problem or a regression problem? Is the goal to build a model that generates a ranked list prioritized by risk, or is it to detect anomalies as new data come in? Knowing what kinds of tasks machine learning can solve will allow you to map the problem you are working on to one or more machine learning settings and give you access to a suite of methods.
- **Data exploration and preparation.** Next, you need to carefully explore the data you have. What additional data do you need or have access to? What variable will you use to match records for integrating different data sources? What variables exist in the data set? Are they continuous or categorical? What about missing values? Can you use the variables in their original form, or do you need to alter them in some way?
- **Feature engineering.** In machine learning language, what you might know as independent variables or predictors or factors

 Blog

---

and

courses,

eract,

# Search and Discovery

The screenshot displays the ADRF interface for 'Project Class1'. The top navigation bar includes the ADRF logo, a search bar, and links for 'Project', 'Explore Data', and 'dcastel...'. Below the navigation, there are three tabs: 'Overview', 'Datasets', and 'Participants'. The 'Datasets' tab is active, showing a list of restricted datasets. The first dataset is 'ILLINOIS DEPARTMENT OF CORRECTIONS Illinois Department of Corrections (DOC) Inmate Admissions 1990-2015'. It includes a description: 'Detailed transactional data of each time a person was admitted to an Illinois Department of Corrections (DOC) facility from 1990 to 2015. Variables include demographic, charges, sentencing, conduct, security level, health and mental health status, gang affiliation. Data are collected using correc...'. It also features a 'Restricted Access' icon, a tag for 'Inmate Populations', and a calendar icon indicating a 25-year period from 1989/12/31 to 2014/12/31. The second dataset is 'US DEPARTMENT OF HOUSING AND URBAN DEVELOPMENT US Department of Housing and Urban Development Program Microdata 2004-2016 - Individuals: Illinois'. Its description is: 'Detailed transactional data consisting of tenant-level data for individuals in the US Department of Housing and Urban Development's (HUD) largest rental assistance programs: the Housing Choice Voucher Program, Public Housing, Project-based Section 8, and the Section 202/811 Programs. The dataset ...'. It also has a 'Restricted Access' icon, a tag for 'Socioeconomic Characteristics', and a calendar icon indicating a 12-year period from 2003/12/31 to 2015/12/31. The third dataset is 'ILLINOIS DEPARTMENT OF CORRECTIONS Illinois Department of Corrections (DOC) Inmate Exits - 1990-2015'. Its description is: 'Detailed transactional data of each time an inmate was released from an Illinois Department of Corrections (DOC) facility from 1990 to 2015. Variables include demographic, residence, charges, sentencing, conduct, security level, health and mental health status, gang affiliation. Data are collecte...'. It also has a 'Restricted Access' icon.

**ADRF** Search Project Explore Data dcastel...

**Project Class1**

[Overview](#) [Datasets](#) [Participants](#)

Below are the restricted datasets available for your project

**ILLINOIS DEPARTMENT OF CORRECTIONS**  
**Illinois Department of Corrections (DOC) Inmate Admissions 1990-2015**

Detailed transactional data of each time a person was admitted to an Illinois Department of Corrections (DOC) facility from 1990 to 2015. Variables include demographic, charges, sentencing, conduct, security level, health and mental health status, gang affiliation. Data are collected using correc...

Restricted Access Inmate Populations 25 years (1989/12/31 - 2014/12/31)

**US DEPARTMENT OF HOUSING AND URBAN DEVELOPMENT**  
**US Department of Housing and Urban Development Program Microdata 2004-2016 - Individuals: Illinois**

Detailed transactional data consisting of tenant-level data for individuals in the US Department of Housing and Urban Development's (HUD) largest rental assistance programs: the Housing Choice Voucher Program, Public Housing, Project-based Section 8, and the Section 202/811 Programs. The dataset ...

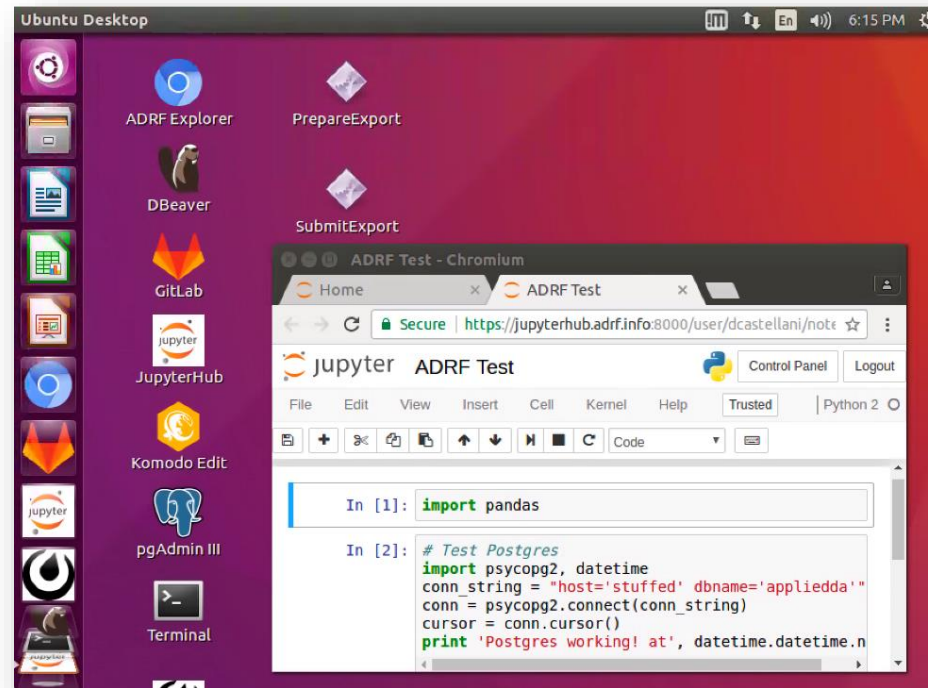
Restricted Access Socioeconomic Characteristics 12 years (2003/12/31 - 2015/12/31)

**ILLINOIS DEPARTMENT OF CORRECTIONS**  
**Illinois Department of Corrections (DOC) Inmate Exits - 1990-2015**

Detailed transactional data of each time an inmate was released from an Illinois Department of Corrections (DOC) facility from 1990 to 2015. Variables include demographic, residence, charges, sentencing, conduct, security level, health and mental health status, gang affiliation. Data are collecte...



# Collaboration





**clayton.hunter** 11:48 AM

hi folks - for anyone using IDHS data in their projects we have a bit more info on programs to help welfare recipients find stable jobs (thanks to Susan H for posing question and Rick Hendra for a great response!) - this doc will also be linked on the class website: <https://docs.google.com/document/d/1GTnuPAWxxtw3CUncX238cWwVbzx6FAdhI5O1pXsuNgg/edit?usp=sharing>



**clayton.hunter** 11:48 AM

shared this file: ▼



## Job assistance programs for welfare recipients

Document from Google Drive

### Job assistance programs for welfare recipients

#### Question posed:

We are trying to add some context to our project and I wondered if you had a contact person at the Illinois DHS that could help fill in some questions about programs available to TANF/benefit recipients. I looked on the [DHS website](#) and while they do have some information, there's not much on programs available to help recipients move to stable jobs. For instance, there's a program called [EPIC](#) directed towards SNAP recipients, but I haven't found much else.

#### Response from Richard Hendra, MDRC:

Yes, we have very specific guidance as we worked on this particular issue there. The ERA evaluation had a site in Chicago that was focused on providing TANF recipients with stable jobs. The short term report [here](#) had more detail about the program, the implementation and the interim effects. Note that the UI data had major coverage issues with the segment of the TANF caseload that we were working with. The final results are in [this](#) giant report. I'd suggest the interim (shorter) report. We used various measures of employment stability. A common measure is the extent to which individuals worked in 4 consecutive

		<b>July – December 2018: Design</b>	<b>Jan-June 2019: Make</b>	<b>July-Dec 2019 Measure and Analyze</b>	<b>Jan-June 2020 Improve</b>
<b>Platform</b>	<b>Activity</b>	- Data Model to incorporate additional metadata about datasets, users, user profiles, and user interactions (i.e., annotations, and explicit connections between datasets, people, and projects) -Telemetry Module to automatically collect structured events emitted by platform	- Deploy Data Model - Deploy Telemetry Module	- Assess Data Model Functionality -Assess Telemetry measures - Open source for community feedback	- Modify Data model with input from Rich Context - Modify Telemetry Module with input from rich context
	<b>Deliverable</b>	Data model Telemetry module	Operational Data Model Functioning Telemetry Module Functioning prototype Initial Jupyter-ADRF integration	QA report Initial prototype stabilized and productionized	Stable and complete version of the application fully integrated to the ADRF Platform. Open sourced
<b>Input Elements</b>	<b>Activity</b>	-Identify and prepare corpora (ICPSR; Bundesbank; Policy area) -Gather requirements	Generate Seed metadata generated ((ICPSR; Bundesbank; Policy area)	Review metadata developed by users Benchmark and revise	Modify and refine metadata capture and documentation
	<b>Deliverable</b>	Three corpora Set of requirements for metadata: comments and annotations on files and datasets, discussions, and contextual recommendations	Metadata for three corpora:	QA and improvement report on the quality of each element	Plan for future improvement
<b>Rich Context</b>	<b>Activity</b>	-Design gamification strategy - Design Pre/Post Survey design - Develop Telemetry measures - Research UX for the collaborative user interfaces i) an interface to help users to ingest Datasets, ii) an interface to help users to create comments and code snippets for Datasets, and iii) an interface to help users to search for Datasets -Design learning approach	Deploy interface Administer Pre survey Capture logging information Test gamification strategy Test learning approach	Review interface Administer post survey Review logging information Review feed back to platform Revise learning approach	Modify and refine interfaces, surveys and learning model
	<b>Deliverable</b>	Survey Telemetry measures Wireframes for the interfaces Learning model	Survey results Log results Gamification results Learning results	Survey results and pre/post analysis Revised UX, feedback loop Revised learning model	Functioning rich context module incorporating human and automated elements with continuous feedback loops to platform



# Rich Context Competition

## PROBLEM DESCRIPTION

Researchers and analysts who want to use data for evidence and policy can't easily find out **who** else worked with the data, on **what topics** and with **what results**. As a result, good research is underutilized, great data go undiscovered and are undervalued, and time and resources are wasted redoing empirical research.

We want you to help us develop and identify the best text analysis and machine learning techniques to discover relationships between data sets, researchers, publications, research methods and fields. We will use the results to create a rich context for empirical research – and build new metrics to describe data use.

This challenge is the first step in that discovery process.

## COMPETITION GOAL

The goal of this competition is to automate the discovery of research datasets and the associated methods and research topic fields in social science research publications. Participants should use any combination of machine learning and data analysis methods to identify the datasets used in a corpus of social science publications and infer the scientific methods used in the analysis and the research fields.

## COMPETITION SPECIFICS

PARTICIPANT INFORMATION
Problem Description
Competition Goal
Competition Specifics
Sponsors
The Bigger Picture
Competition Schedule
How to Participate
Remuneration
Judges
Program Requirements
Phase 1
Phase 2
Competition Terms And Conditions
Teams

# Key challenges to be solved with metadata – particularly for federal statistical system

- Limited internal capacity
- Security
- Legal mandates surrounding access and use
- Data sharing issues
  - cost
  - burden
  - data quality
  - data documentation
  - risk of bad analysis



# Comments and questions?

- If interested in contributing – contact me at
- [Julia.lane@NYU.EDU](mailto:Julia.lane@NYU.EDU)
- More info at <https://coleridgeinitiative.org> and <http://jupyter.org>