

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
114th 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

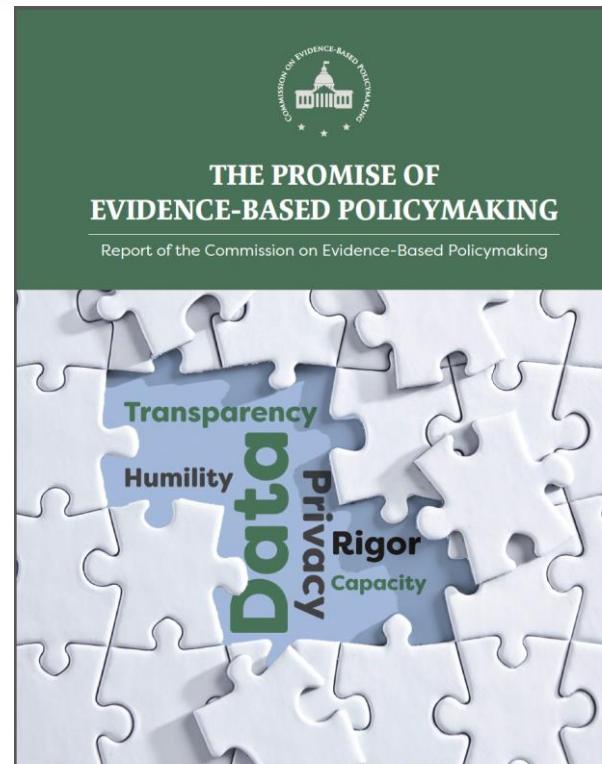
Sponsor: **Paul Ryan**
Representative for Wisconsin's 1st congressional district
Republican

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.¹ 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

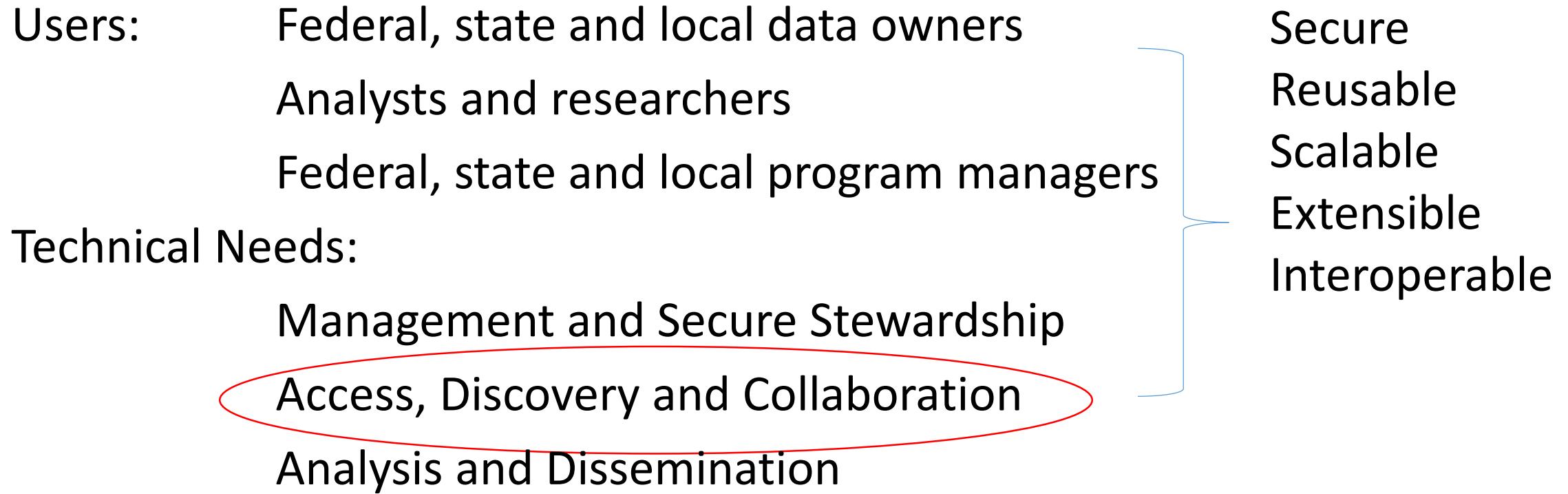
Data

- ▶ 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 (<=4s) 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

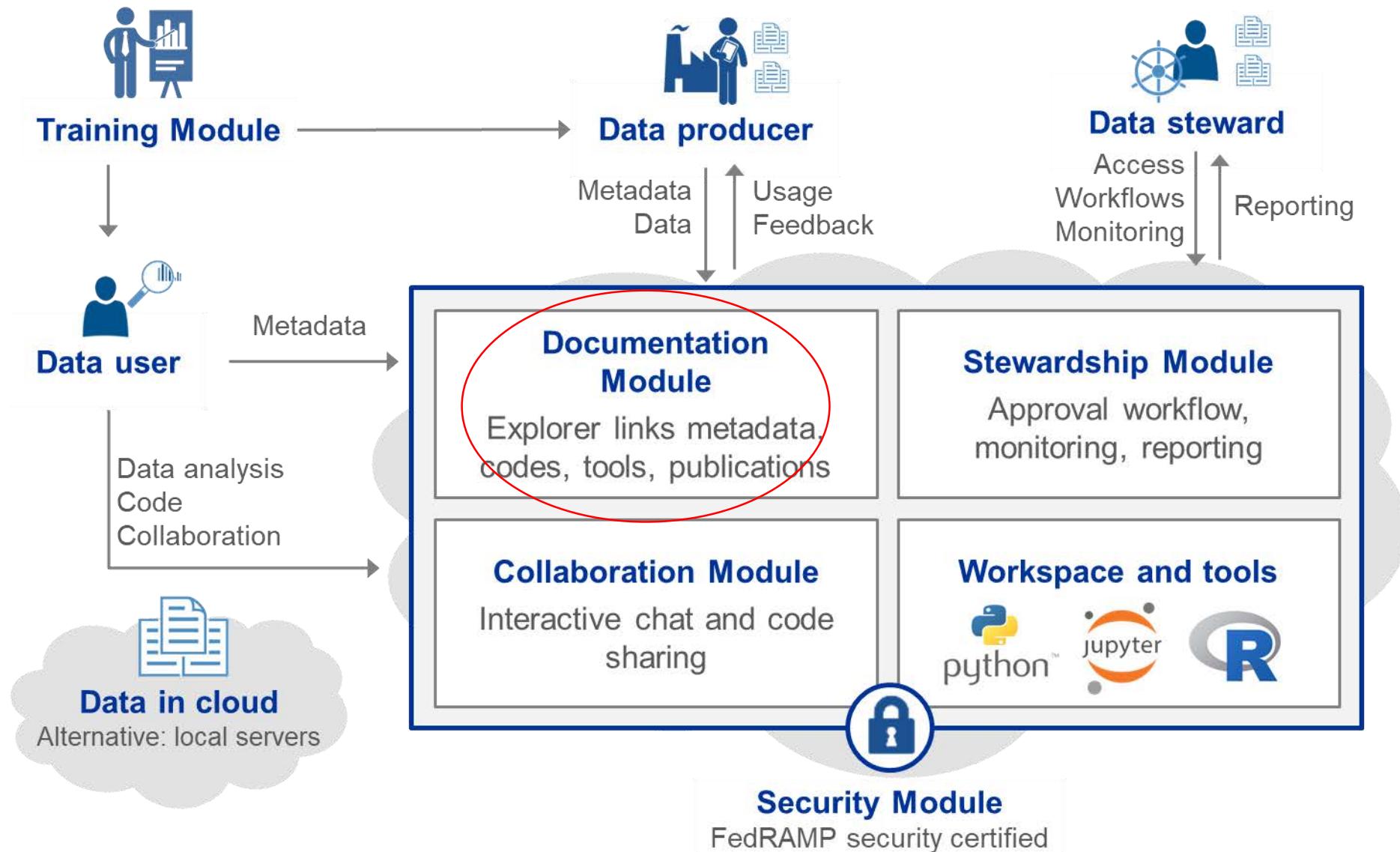
Code

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

Build technical environment



Functional characteristics



Inspiration

Carole Goble

From Wikipedia, the free encyclopedia

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

- Contents [hide]
- 1 Education
- 2 Research
- 3 Career
- 4 Awards and honours
- 5 References

Education [edit]

Goble was educated at Maidstone Grammar School for Girls.^[1] 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^[21] to 1982.

Research [edit]

Her current research interests^{[1][22]} include Grid computing, the Semantic Grid,^[23] the Semantic Web, Ontologies,^{[24][25][26]} e-Science, medical informatics,^[27] Bioinformatics, and Research Objects. She applies advances in knowledge technologies and workflow systems^[28] to solve information management problems for life scientists and other scientific disciplines^[citation needed]. 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),^[29] Biotechnology and Biological Sciences Research Council (BBSRC),^[30] 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.^[31]

Her work has been published in leading peer reviewed scientific journals including *Nucleic Acids Research*,^[32] *Bioinformatics*,^{[32][33]} *IEEE Computer*,^[10] the *Journal of Biomedical Semantics*,^[34] *Briefings in Bioinformatics*,^{[35][36][37]} *Artificial Intelligence in Medicine*,^[27] the Pacific Symposium on Biocomputing conference,^[24] the *International Journal of Cooperative Information Systems*, the *Journal of Biomedical Informatics*,^[38] *Nature Genetics*,^[39] and *Drug Discovery Today*.^{[40][41][42][43][44][45]}

Career [edit]



The Taverna Suite of Tools

Workflow Repository



Service Catalogue



BiodiversityCatalogue

The Biodiversity Sciences Web Services Registry

Activity and Service Plug-in Manager



User Interfaces



Taverna
Workbench



Taverna
Lite

Workflow Provenance



Workflow Components



Workflow Server



Interaction Server

Secure Service Access

Web Portals / Gateways Client User Interfaces



Third Party Tools



Player

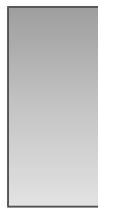


Command Line



RESEARCH

▪ Github



Isidor Ni
isidorn

VSCode



Block or report u

Microsoft



Zurich

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

LWJGL / lwjgl3

Watch 159 Star 1,430 Fork 229

Issues 32 Pull requests 0 Projects 1 Wiki Insights

New issue

Java 9 got released! #334

Open dustContributor opened this issue 25 days ago · 13 comments

dustContributor commented 25 days ago

So, how well LWJGL 3 plays with it and Jigsaw? I remember @Spasi mentioning something about moving to VarHandles instead of relying on sun.misc.Unsafe. What would this entail?

Spasi added the Type: Question label 25 days ago

Spasi commented 25 days ago

LWJGL 3 works great on Java 9. You don't have to do anything special and it runs without --illegal-access warnings.

There are no JPMS modules, but the JARs include Automatic-Module-Name entries in their manifests. This should be good enough for users that want to make their projects modular. I don't think we'll do anything more involved with modules until LWJGL 4.

The other Java 9 feature in LWJGL is multi-release JAR files. The core library is such a JAR and it includes custom code that uses the new StackWalker API to improve performance when MemoryStack debugging is enabled.

something about moving to VarHandles instead of relying on sun.misc.Unsafe

Unfortunately java.lang.invoke.VarHandle is not a full replacement for sun.misc.Unsafe. The first problem is worse performance when accessing off-heap memory. The second is that the security mechanisms in

Microsoft/vscode 53 commits
Microsoft/vscode-node-debug 2 commits

Following 0

Assignees
No one assigned

Labels
Type: Question

Projects
None yet

Milestone
No milestone

Notifications
Subscribe
You're not receiving notifications from this thread.

5 participants
dustContributor, Spasi, TM, MM, u5

Oct
More

2017
2016
2015
2014

RESEARCH



ttja
Lincoln
Unit
King
11

- Github
- Data.world
- Pinterest
- TripAdvisor



Linda
6

Overview Rooms Reviews About Photos Nearby Q&A Room Tips \$545 View Deal

Overview

5.0 732 reviews

Excellent 92% Very good 5% Average 1% Poor 1% Terrible 1%

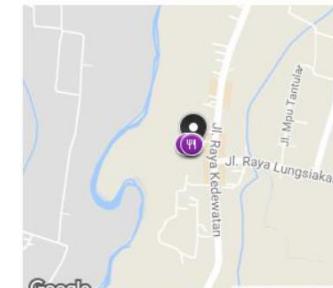
Free Wifi Non-Smoking Hotel
Free Parking Restaurant
Breakfast included 5.0 Star Hotel
Air Conditioning All hotel details
Pool

TRAVELERS TALK ABOUT

- "kubu restaurant" (52 reviews)
- "ayung river" (85 reviews)
- "kids club" (23 reviews)

OFFERS FROM MANDAPA, A RITZ-CARLTON RESERVE

Hotel packages



See all 199 hotels in Ubud

Similar hotels



Four Seasons Resort Bali at Sayan
 1,389 reviews
#1 of 4 hotels in Sayan
\$518



Kupu Kupu Barong Villas and Tree Spa
 1,529 reviews
#3 of 5 hotels in Kedewatan
\$135



COMO Uma Ubud
 1,407 reviews
#14 of 199 hotels in Ubud
\$278

TERIBBLE 0

0 POINTS 1

300 points to go

View Collection

Leadership 300 Readers

for Reviewer 5 Reviews

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

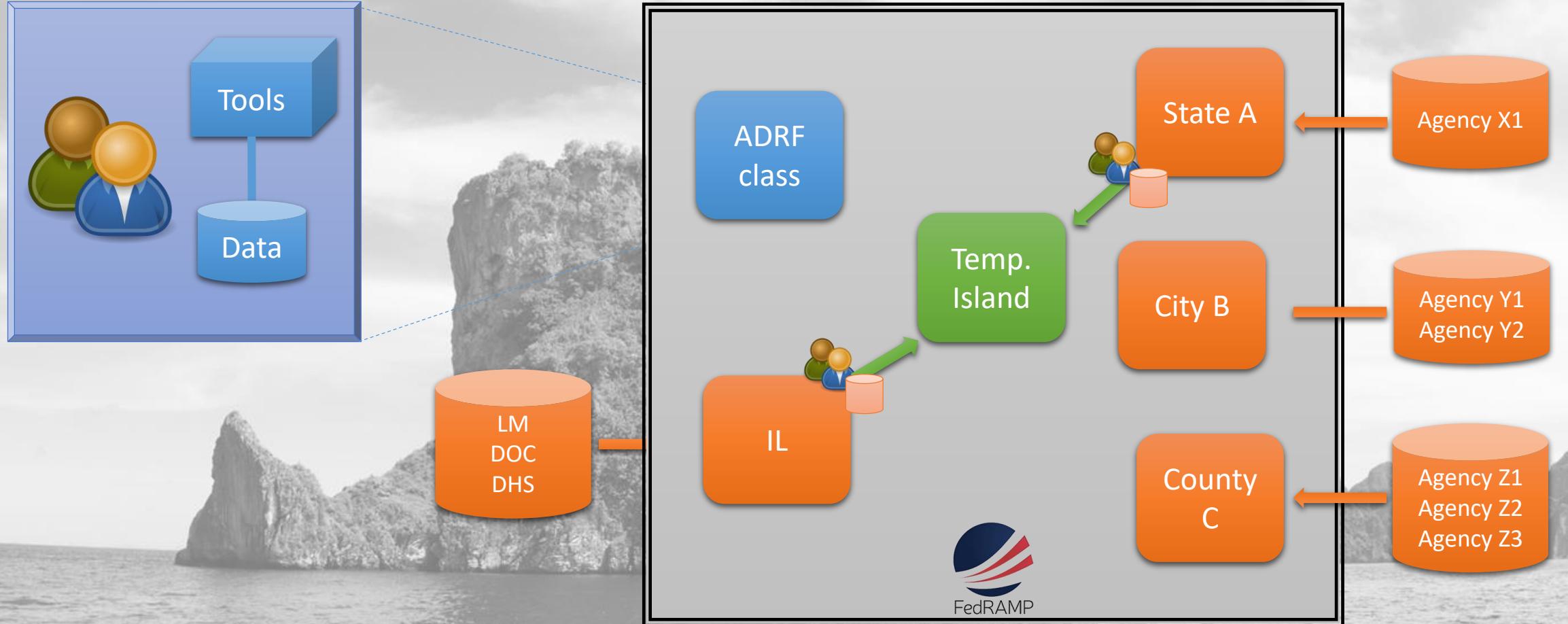


Alfred P. Sloan
FOUNDATION

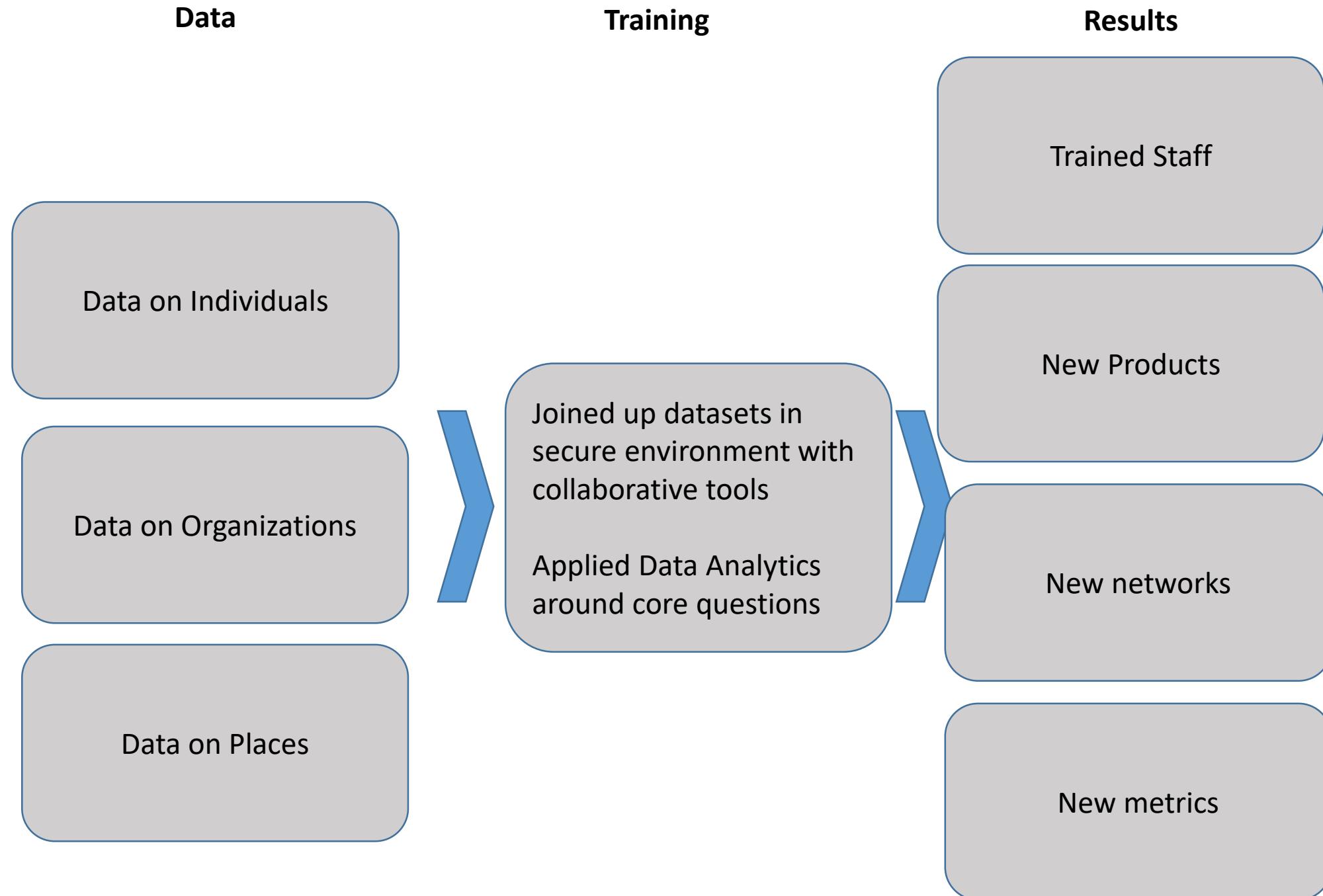
SCHMIDT FUTURES



Overdeck Family
Foundation



ADR福 SaaS



A screenshot of a Jupyter Notebook interface. At the top, there is a header bar with icons for file operations, a GitHub logo, and a Jupyter logo. Below the header, the notebook content is displayed in cells:

- Cell 1:** Python code for machine learning models (GaussianNB, DecisionTreeClassifier) and styling (sns.set_style, sns.set_context). A large red arrow points from this cell down to the "Connect to Database" section.
- Cell 2:** Python code to connect to a PostgreSQL database named "appliedda" on host "10.10.2.10".
- Cell 3:** Python code to read a table named "ides.il_wage" from the database and limit the results to 10 rows.
- Cell 4:** Python code to display the first few rows of the DataFrame.

The notebook interface includes a vertical sidebar on the left with sections for "What is JupyterHub?", "Key features", and "GitHub". On the right, there is a sidebar with links for "Blog" and "Courses".

The Machine Learning Process

[Go back to Table of Contents](#)

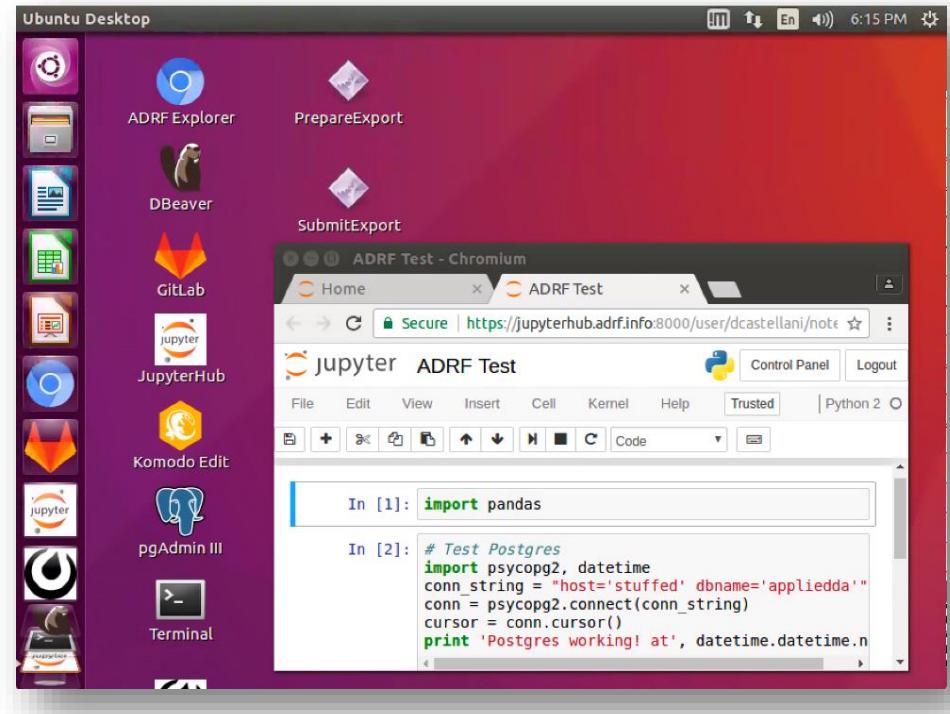
- **Understand the problem and goal.** This sounds obvious but is often nontrivial. Problems typically start as vague descriptions of a goal - improving health outcomes, increasing graduation rates, understanding the effect of a variable X on an outcome Y , etc. It is really important to work with people 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

Search and Discovery

The screenshot shows the ADRF (Advanced Data Research Framework) interface. At the top, there is a navigation bar with the ADRF logo, a search bar, and user information. Below the navigation bar, the title "Project Class1" is displayed next to a briefcase icon. The main content area features three dataset cards:

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



#class-3-f
#class-
☆ | ✎ | 897



Elena Semenova 9:09 PM

HI DOC data gurus! Do you know what the following indicates in reality? A person admitted first time in \geq 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 \geq 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?

Samip, Suren
E
L
C



Vivek Ananda 11:27 PM

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

Samip, Suren
When I LEFT
with 4x recor

1. Is there an
2. If not, how



Drew 5:29 PM

@Beau Ande
exclusively co
be able to loc
what you are
what you are

5 rep

looked on
available
towards S

Respons
Yes, we ha
a site in C
here had
data had
final resul
of employ

#class-3-fall17

When I LEFT
☆ | 897 | 0 | Add a topic



Search



Thursday, January 4th



Elena Semenova 11:49 AM

I asked that before and didn't get an answer. Does someone know how ildoc.ildoc_exit.jailtime is calculated? It doesn't equal to any interval between dates in fields: exit_date, curadm_date, cccadm_date, cccvio_date, actmsr_date. Should we consider that value at all or rely on calculated values between mentioned data? Also, ILDOC_EXTI data dictionary is missing some fields. Please confirm if cccvio_date means date of CCC violation (work release to community correctional center).

clayton.hunter 11:54 AM

we may need to check with @Vivek Ananda or @Dana Wilson for confirmation, but based on the description of jailtime in ADRF Explorer I suspect those are cumulative values for each individual - so cannot just be calculated based on that individual record

clayton.hunter 11:56 AM

and cccvio_date is a helper column that combines all cccvio* columns into a single, date formatted column so that postgres date functions work properly (I believe that is the case for all columns that end in _date)

1 reply 6 days ago

Drew 12:05 PM

@Elena Semenova sorry about this, might have gotten lost in the shuffle a while back but Vivek did provide the following information on jailtime in an e-mail: Jail time is calculated on how much time inmate spent in jail prior to coming to prison. He does get credit for time served at all jails prior coming to prison. Thought I had circulated, but maybe only updated on the metadata in the explorer (edited)

2

		July – December 2018: Design	Jan-June 2019: Make	July-Dec 2019 Measure and Analyze	Jan-June 2020 Improve
Platform	Activity	- 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
	Deliverable	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
Input Elements	Activity	-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
	Deliverable	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
Rich Context	Activity	-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
	Deliverable	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

Timeline



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
- More info at <https://coleridgeinitiative.org> and <http://jupyter.org>