

---

# Automating Metadata Documentation: A New Initiative

(OR DATA, DATA, EVERYWHERE – WE HAVE TO STOP AND THINK)

Julia Lane

New York University

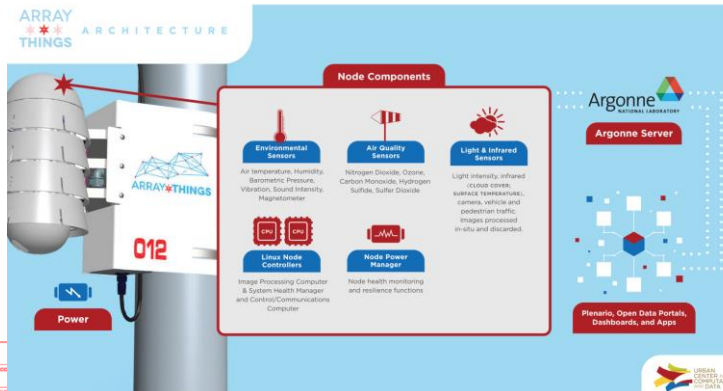
(with more contributors than I can count - but notably Frauke Kreuter, Rayid Ghani, Clayton Hunter, Brian Granger, Stefan Bender, Christian Hirsch, Hendrik Doll, Stu Feldman, Drew Gordon, Jonathon Morgan, Ekaterina Levitskaya, Daniel Castellani, Rafael Alves, Ophir Frieder, Jordan Boyd Graber, Evgeny Klochikhin, Christian Herzog, Ian Mulvaney, Alan Maloney, Daniel Hook)

# Takeaways

- Massive supply of new granular data
- New demand - local/city/regions
- Data science offers engagement and feedback tools
- Statistical skills critical
- Need to rethink data infrastructure

Data availability is burgeoning; curation, access and reproducibility is not...we need new measurement, tools and institutions

# Much, much more data



22222	Valid <input type="checkbox"/>	Employer's social security number	For Official Use Only OMB No. 1545-0048
1	Employer identification number (EIN)	2	Federal income tax withheld
3	Social security wages	4	Social security tax withheld
5	Medicare wages and tips	6	Medicare tax withheld
7	Social security type	8	Allocated tips
9	Verification code	10	Dependent care benefits
11	Long-term plans	12a	See instructions for box 12
13	Statutory rate	13b	Statutory rate
14	Other	15	State income tax
16	State income tax	17	Local income tax
18	Local income tax	19	Local income tax
20	Local income tax	21	Local income tax

**W-2 Wage and Tax Statement 2017**  
 Form W-2 For Social Security Administration — Send this entire page with Copy A For Social Security Administration; photocopies are not acceptable. Form W-3 to the Social Security Administration; photocopies are not acceptable.  
 Do Not Cut, Fold, or Staple Forms on This Page.



## 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

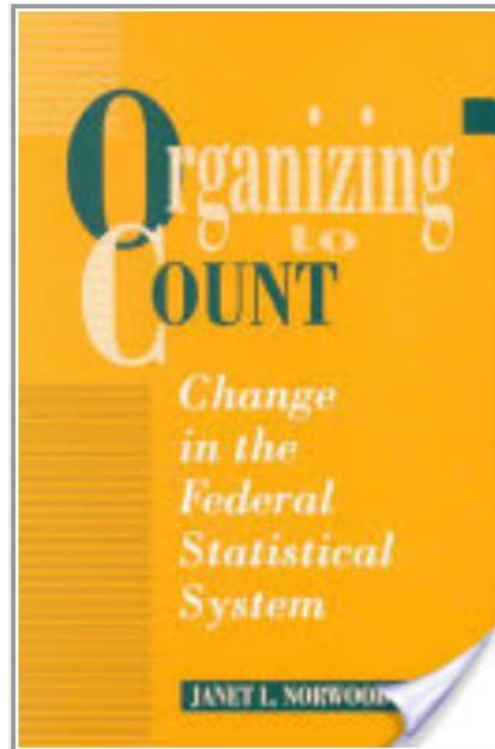
- Scaper – A Python library for soundscape synthesis and augmentation
- Audio-Annotator – A Javascript web interface for annotating audio data
- Raster Join
- Urban Pulse

# There are many core measurement issues

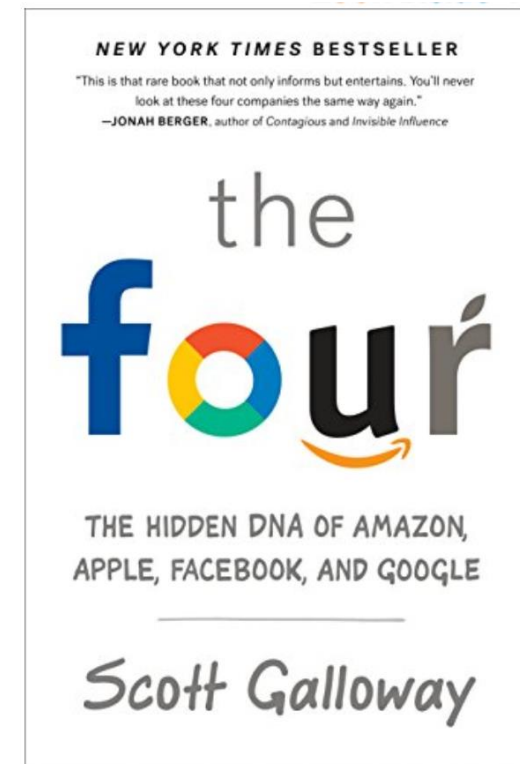
- New products
  - Services
  - Intangibles
- New people
  - Immigration
  - Globalization
- New boundaries
  - Local
  - Regional
  - Cross national
- New organizations
  - Firms
  - Platforms
  - Networks
- New tradeoffs
  - Timeliness?
  - Closeness to core measure?
  - Coverage?
  - Geographic detail
  - Longitudinal Consistency
- New privacy constraints

# Skills and institutional gaps

- The past



- The future



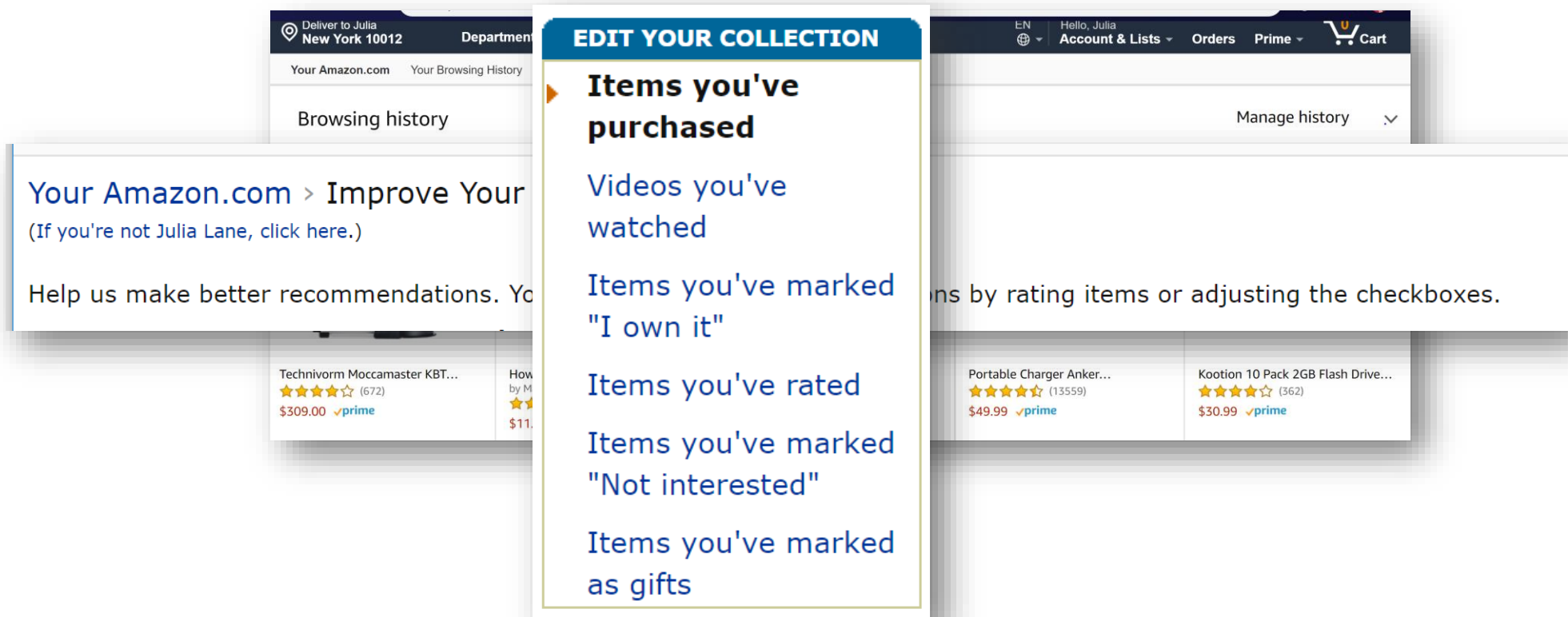


# Vision



# Private Sector businesses collect, curate and disseminate data

## Search and Discovery is automated



The image shows a screenshot of the Amazon.com 'EDIT YOUR COLLECTION' menu. The menu is overlaid on a blurred background of the Amazon website. The menu items are:

- Items you've purchased
- Videos you've watched
- Items you've marked "I own it"
- Items you've rated
- Items you've marked "Not interested"
- Items you've marked as gifts

The background shows the Amazon.com interface with a navigation bar at the top, a search bar, and a list of products. The products shown include a Technivorm Moccamaster KBT... and a Kootion 10 Pack 2GB Flash Drive... The navigation bar includes links for 'Account & Lists', 'Orders', 'Prime', and 'Cart'.

# Other sciences have done the same

## Carole Goble

From Wikipedia, the free encyclopedia

**Carole Anne Goble**, *CBE FRS* (born 10 April 1961) is a British academic who is Professor of *Computer Science* at the University of Manchester.<sup>[1][4][5]</sup> She is Principal Investigator (PI) of the *myGrid*,<sup>[1][6]</sup> *BioCatalogue*<sup>[17]</sup> and *myExperiment*<sup>[18]</sup> projects and co-leads the Information Management Group (IMG) with Norman Paton.<sup>[1][9][20]</sup>

- 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.<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>[21]</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*,<sup>[27]</sup> *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

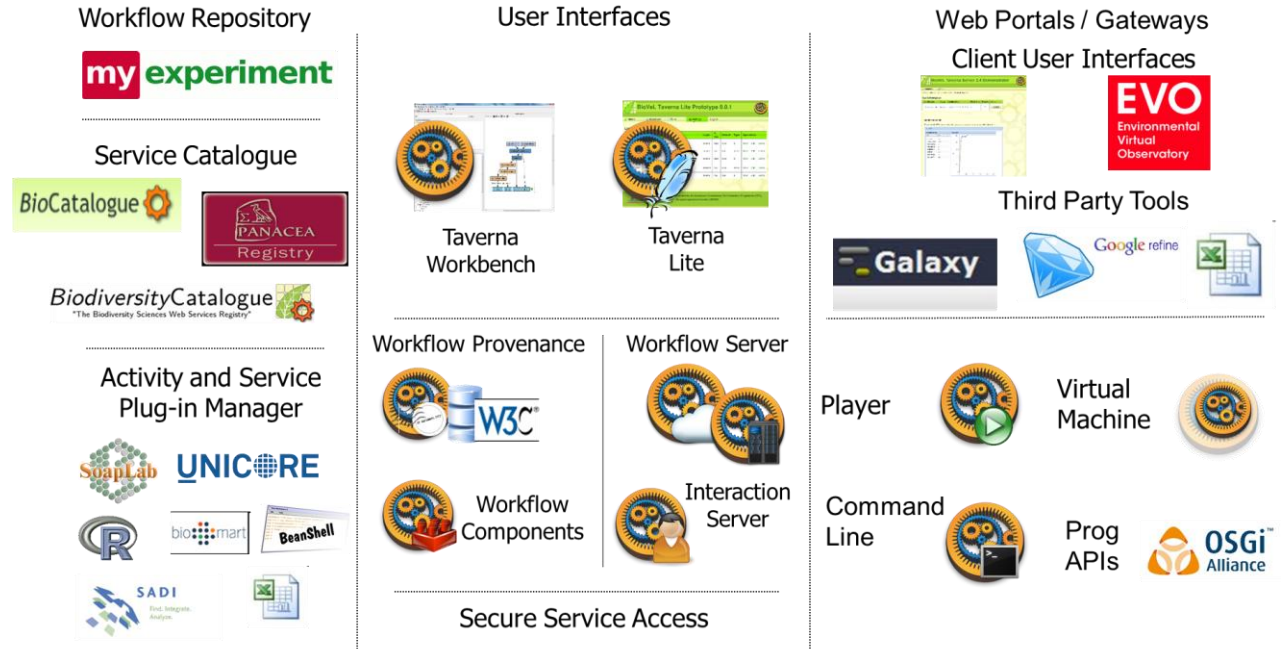
**Alma mater** University of Manchester

**Known for** myGrid  
Semantic Grid  
Open PHACTS<sup>[2]</sup>  
Taverna workbench<sup>[3][4]</sup>  
Software Sustainability Institute

**Spouse(s)** Ian Cottam (m. 2003)<sup>[5]</sup>

**Awards** The Seven Deadly Sins of Bioinformatics<sup>[6]</sup>  
Jim Gray e-Science Award (2008)

# The Taverna Suite of Tools



# Vision for Search and Discovery

User



**Isidor Nikolic**  
isidom

VSCode

Follow

Block or report user

- Microsoft
- Zurich

## TripAdvisor

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
- Free Parking
- Breakfast included
- Air conditioning
- Pool

- Non-Smoking Hotel
- Restaurant
- 5.0 Star Hotel
- All hotel details

TRAVELERS TALK ABOUT

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

OFFERS FROM MANDAPA, A RITZ-CARLTON RESERVE

- Hotel packages

**Similar hotels** [See all 199 hotels in Ubud](#)

**Four Seasons Resort Bali at Sayan**  
4.5 1,388 reviews  
#1 of 4 hotels in Sayan  
**\$518**

**Kupu Kupu Barong Villas and Tree Spa**  
4.5 1,522 reviews  
#3 of 5 hotels in Kedewatan  
**\$135**

**COMO Uma Ubud**  
4.5 1,407 reviews  
#14 of 199 hotels in Ubud  
**\$278**

Overview Rooms **Reviews** About Photos Nearby Q&A Room Tips

leonardomatheo  
London, United Kingdom  
3

5.0 Reviewed yesterday

**beautiful**

great hotel!... i'm glad I chose to stay there for 3 nights. me and my friends had a lovely time at the hotel. the staff were kind and helpful. the rooms were clean. the location was mind blowing and unbelievable. Im hoping I get the... [More](#)

[Thank leonardomatheo](#)

---

nich0lemaried  
Los Angeles, California  
5

5.0 Reviewed 4 days ago

**7 Star Luxury**

Couldn't have asked for anything more at this extraordinary place. It is truly a spiritual and healing location in the lap of the luxurious wilds of Bali. The service anticipated my every need and was most gracious about any request. I will most definitely return. [More](#)

Review collected in partnership with The Ritz-Carlton Hotel Company ©

[Thank nich0lemaried](#)

Response from Aghpt, General Manager at Mandapa, A Ritz-Carlton Reserve  
Responded today

Dear nich0lemaried, It was a great pleasure to have you staying at Mandapa, a Ritz-Carlton Reserve, in Ubud for your holiday! We are very appreciative of your wonderful comments about our facilities, the location and the service offered by our Ladies and Gentlemen. We look... [More](#)

LB American Com

← → ↺ ↻

**Elena Semenova** 9:09 PM  
HI DOC data gurus! Do you know what incarceration for lower offence class (he/she was hiding from law enforcement

**Vivek Ananda** 11:27 PM  
It mostly is bad data please email me t

**#class-3-fall17**  
☆ | 👤 97 | 🚩 0 | ➕ Add a topic

**Elena Semenova** 11:49 AM  
I asked that before and didn't get an answer between dates in fields: exit\_date, cur values between mentioned data? Also violation (work release to community

**clayton.hunter** 11:54 AM  
we may need to check with @Vivek Ananda or @Dana W suspect those are cumulative values for each individual - s

**clayton.hunter** 11:56 AM  
and ccvio\_date is a helper column that combines all ccvio properly (I believe that is the case for all columns that end

**1 reply** 6 days ago

**Drew** 12:05 PM  
@Elena Semenova sorry about this, might have gotten lost in an e-mail: Jail time is calculated on how much time in jails prior coming to prison. Thought I had circulated, but

11

Slide 8 of 65

**#class-3-fall17**  
☆ | 👤 97 | 🚩 0 | ➕ Add a topic

Wednesday, January 3rd

📞 ⓘ ⚙️ 🔍 Search @ ☆ ⋮

**Beau Anderson (CT)** 4:27 PM  
Using the idhs.hh\_indcase\_spells table, I'm able to select a sample for our project that has x unique records (based on ssn\_hash) in it that meet the criteria, and store the results into a table on our team's schema. This is well and good, but we want more info about the heads of household in our sample, such as the information in the idhs.household\_info table (sex, educational attainment, health, work experience)

**#class-3-fall17**  
☆ | 👤 97 | 🚩 0 | ➕ Add a topic


December 21st, 2017

📞 ⓘ ⚙️ 🔍 Search

1. Is there  
2. If not, I

**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 (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/1GTnuPAWxxtw3CUncX238cWwVbzx6FAdh15O1pXsuNgg/edit?usp=s>

**clayton.hunter** 11:48 AM  
shared this file:

 **Job assistance programs for welfare recipients**  
Document from Google Drive

Job assistance programs for welfare recipients

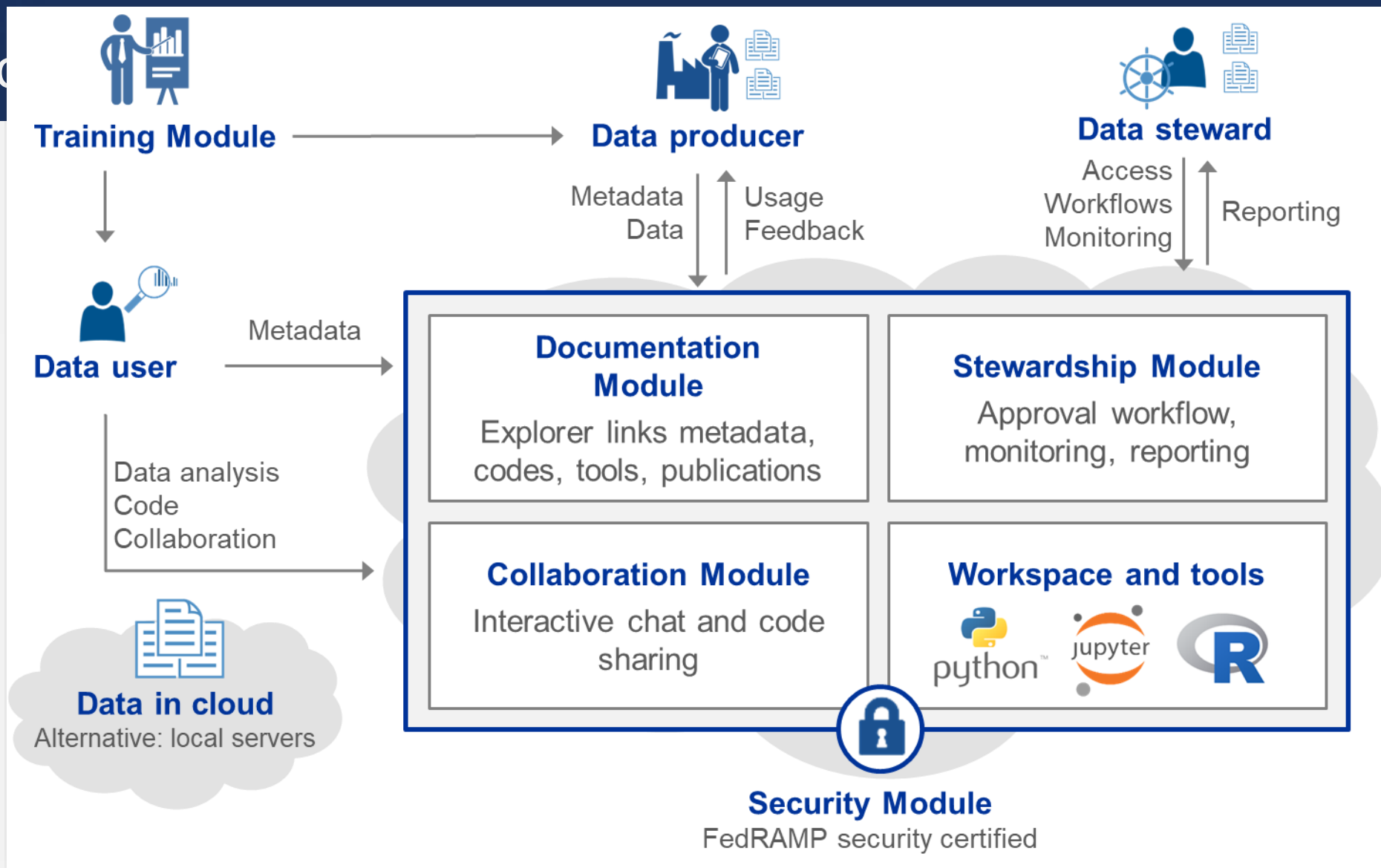
**Question posed:**  
We are trying to add some context to our project and I wondered if you had a contact person at the Illinois DHS that could help fill in some questions about programs available to TANF/benefit recipients. I looked on the DHS website and while they do have some information, there's 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

## Data Discovery

- Step 1: Create the set of corpora and metadata (computer science technology) - Competition
- Step 2; Figure out how you learn from it and automate it (machine learning techniques) - Engagement
- Step 3: Gamification – recognize and emphasize patterns (with human curation) – Rinse and repeat

# Operational





		<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 <u>feedback</u> 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



# Competition



Ident

FIND DATA ▾

START SHARING DATA ▾

MEMBERSHIP ▾

SUMMER PROGRAM ▾

TEACHING &amp; LEARNING ▾

DATA MANAGEMENT &amp; CURATION ▾



ICPSR

Membership in ICPSR

Log In/Create Account



MEMBER LIST

HOW TO JOIN

OFFICIAL REP TOOLS

PROMOTING ICPSR

NEWS &amp; PUBLICATIONS

## Obtaining ICPSR Metadata

Batch XML via [OAI-PMH](#)

Batches of ICPSR study-level metadata are available in three formats:

DDI 2.5 (with or without citations)  
Dublin Core  
MARC21

To obtain multiple files, please use our [OAI-PMH](#) service.

### Individual Records

ICPSR study-level metadata for individual studies are available in a variety of formats:

DDI 2.5 and 3.x  
Dublin Core  
MARC21 XML  
DATS 2.1

To obtain individual files, simply look for the export links near the bottom of a study home page [\(example\)](#).

## Metadata Records

- [General Information](#)
- [Conditions of Use](#)
- [Obtaining ICPSR Metadata](#)

## General Information

ICPSR has produced study descriptions, or metadata records, describing its holdings since the organization began in 1962. In the 1980s the metadata records were converted to a standardized, fielded text format. In 1999, ICPSR received an NSF Infrastructure in the Social Sciences award (SES-9977984), which enabled us to enrich the content of the records and to convert them to XML. The first XML records were produced to be compliant with the Data Documentation Initiative (DDI) Version 1.0 tag set. The records now comply with DDI versions 2.5 and 3.1. In 2006, ICPSR, in partnership with the University of Michigan Library, produced MARC records for its collection. More recently, ICPSR has added machine-readable schema.org markup in JSON-LD format to dataset landing pages.

Please note that 'metadata records' does not refer to the larger PDF documentation files associated with each study, nor to the actual data files.

Some researchers and librarians find it useful to have access to our metadata records. On this page, we describe how to obtain those files and how we maintain the metadata files over time.

## Conditions of Use

ICPSR shares its metadata records with the membership to promote wider awareness and use of ICPSR's social science data resources. In particular, we encourage members to integrate the records into local Online Public Access Catalogs (OPACs) intended primarily for the use of faculty, staff, and students at their institutions.



ICPSR metadata records are licensed under a [Creative Commons Attribution-Noncommercial 3.0 United States License](#).

ICPSR also encourages users to ensure good use of the metadata records:

# Create

inputs > ICPSR\_citations\_metadata\_enriched\_full\_combined\_coded\_status.csv

inputs > data\_sets

```
[
  {
```

```
    "date": "2016-09-20",
    "url": "http://www.jstor.org/stabl",
    "title": "ANES 1952 Time Series St",
    "name": "ANES 1952 Time Series St",
    "description": "ANES 1952 Time Series St",
    "surveys": "ANES 1952 Time Series St",
    "fields": "ANES 1952 Time Series St",
    "present": "ANES 1952 Time Series St",
    "data": "ANES 1952 Time Series St",
    "social": "ANES 1952 Time Series St",
    "and": "ANES 1952 Time Series St",
    "p": "ANES 1952 Time Series St",
    "opinions": "ANES 1952 Time Series St",
    "on": "ANES 1952 Time Series St",
    "1952": "ANES 1952 Time Series St",
    "National": "ANES 1952 Time Series St",
    "attitudes": "ANES 1952 Time Series St",
    "an": "ANES 1952 Time Series St",
    "interview": "ANES 1952 Time Series St",
    "sc": "ANES 1952 Time Series St",
    "collect": "ANES 1952 Time Series St",
    "data": "ANES 1952 Time Series St",
    "before": "ANES 1952 Time Series St",
    "and": "ANES 1952 Time Series St",
    "a": "ANES 1952 Time Series St"
```

2015-11-10	10.3886/ICPSR07355	ANES 1974 Time Series St	1978-01-01	10.1017/S00071234000C	doi	Journal Article	literature	Generational replacemen
2015-11-10	10.3886/ICPSR07381	ANES 1976 Time Series St	1978-01-01	10.1017/S00071234000C	doi	Journal Article	literature	Generational replacemen
1992-02-16	10.3886/ICPSR07233	Political Change in Britain,	1978-01-01	10.1017/S00071234000C	doi	Journal Article	literature	Generational replacemen
1992-02-16	10.3886/ICPSR07234	Political Change in Britain,	1978-01-01	10.1017/S00071234000C	doi	Journal Article	literature	Generational replacemen
1992-02-16	10.3886/ICPSR07004	Political Change in Britain,	1978-01-01	10.1017/S00071234000C	doi	Journal Article	literature	Generational replacemen
2016-09-20	10.3886/ICPSR07213	ANES 1952 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2016-09-22	10.3886/ICPSR07214	ANES 1956 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2016-09-22	10.3886/ICPSR07215	ANES 1958 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2015-11-10	10.3886/ICPSR07216	ANES 1960 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2016-12-01	10.3886/ICPSR07217	ANES 1962 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2015-11-10	10.3886/ICPSR07235	ANES 1964 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2015-11-10	10.3886/ICPSR07259	ANES 1966 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2015-11-10	10.3886/ICPSR07281	ANES 1968 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2015-11-10	10.3886/ICPSR07298	ANES 1970 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2016-09-20	10.3886/ICPSR07010	ANES 1972 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2015-11-10	10.3886/ICPSR07355	ANES 1974 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2015-11-10	10.3886/ICPSR07381	ANES 1976 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2015-11-10	10.3886/ICPSR07655	ANES 1978 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
2016-02-26	10.3886/ICPSR07763	ANES 1980 Time Series St	1982-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	The decline of electoral
1999-10-07	10.3886/ICPSR09093	American National Electio	1992-01-01	http://www.jstor.org/stabl	url	Journal Article	literature	'Sophisticated' voting in

## ARTICLE

## Reported Consumption of Low-Foods by American Children and Adolescents

*Nutritional and Health Correlates, NHANES III*

Ashima K. Kant, PhD

**Objective:** To examine the contribution of foods of modest nutritional value to the diets of American children and adolescents.

**Methods:** The data were from the third National Health and Nutrition Examination Survey, 1988 to 1994, and included 4852 children and adolescents, aged 8 to 18 years. Foods reported in the 24-hour dietary recall were grouped into the following low-nutrient-density (LND) food categories: visible fat; table sweeteners, candy, and sweetened beverages; baked and dairy desserts; salty snacks; and miscellaneous. The independent association of the number of LND foods mentioned in the recall with intake of food groups, macronutrients, micronutrients, and body mass index was examined by means of regression procedures to adjust for multiple covariates.

**Results:** The LND foods contributed more than 30% of daily energy, with sweeteners and desserts jointly

account for 30% of total energy intake. The LND foods contributed more than 30% of daily energy, with sweeteners and desserts jointly

**M**ANY AMERICAN children and adolescents consume diets that provide marginal amounts of several nutrients, including vitamins A, E, B<sub>6</sub>, folate, and the minerals calcium, magnesium, iron, and zinc.<sup>1,2</sup> These essential nutrients have well-known metabolic functions, and many have been linked to health promotion and disease prevention.<sup>3</sup> Recent survey data also suggest an alarming trend in increasing prevalence of adiposity in US children and adolescents.<sup>7,8</sup> Estimates from the Continuing Survey of Food Intakes by Individuals, 1994 to 1996, suggest a dramatic increase in consumption of

## METHODS

The NHANES III is a multistage, stratified, probability sample of the noninstitutionalized, civilian US population, aged 2 months and older.<sup>20</sup> The survey was conducted by the National Center for Health Statistics and included administration of a questionnaire at home and a full medical examination along with a battery of tests in a special mobile examination center.<sup>20</sup> Demographic and medical history information was obtained during the household interview. The examination at the mobile examination center included physical and dental examinations, dietary interview, body measurements, and collection of blood and urine samples. Body weight, height, and circumference at various body sites were measured by standardized procedures in the mobile examination center.<sup>20</sup>

## DIETARY ASSESSMENT METHOD

A 24-hour dietary recall was collected by a trained dietary interviewer in a mobile examination center interview with the use of an automated, microcomputer-based interview and coding system.<sup>20</sup> The type and amount of foods consumed were recalled with recall aids such as abstract food models, special charts, measuring cups, and rulers to help in quantifying the amounts consumed. Special probes were used to help the recall of commonly forgotten items such as condiments, accompaniments, fast foods, and alcoholic beverages.

## ANALYTIC SAMPLE

All NHANES III respondents aged 8 to 18 years with a 24-hour recall considered complete and reliable by the National Center for Health Statistics were eligible for inclusion in this study (n=4889). Recalls of pregnant (n=34) or nursing (n=3) respondents were excluded, to result in a final analytic sample size of 4852 (2383 boys and 2469 girls) representing a population estimate of 38738796 children and adolescents.

amination of the contribution of added sugar or carbonated beverages to the diets of children.<sup>11-18</sup> There is some evidence that food selection patterns favored in childhood may track through adult years<sup>19</sup>; therefore, a better understanding of the contribution of

ages in the recall, (2) amount (in grams) of LND foods and beverages, and (3) percentage of total energy from LND foods and beverages.

The NHANES III nutrient database for individual foods, which is derived from the US Department of Agriculture's Survey Nutrient Database, was used for determining energy and nutrient content of all foods.<sup>23</sup> We examined the mean intake of selected micronutrients (vitamins A, E, C, B<sub>6</sub>, and folate, and the minerals calcium, magnesium, iron, and zinc) by tertiles of the reported number of LND foods. The reported intake of each nutrient was also examined with reference to the most recent age-sex-specific standard available, the estimated average requirement.<sup>3-6</sup> As an estimate of dietary misreporting, the ratio of reported energy intake to energy expenditure for basal needs was computed. Energy expenditure for basal needs was estimated by prediction equations developed by the Dietary Reference Intakes Committee for normal and overweight or obese 3- to 18-year-olds.<sup>24</sup>

Data on serum concentrations of vitamin C, folate, and carotenoids were obtained from the National Center for Health Statistics public release compact disks.<sup>25,26</sup> The methods used for measurement of these serum analytes and their associated errors have been described.<sup>25,26</sup> Serum folate, ascorbate, and the carotenoids— $\alpha$ -carotene,  $\beta$ -carotene,  $\beta$ -cryptoxanthin, lutein-zeaxanthin, and lycopene—were chosen because the dietary intake of each nutrient is believed to be a major determinant of serum concentration of the respective nutrient and can thus serve as a biomarker of dietary exposure.<sup>4,5</sup> Serum homocysteine has been reported as an independent predictor of coronary heart disease risk and may have dietary determinants.<sup>27,28</sup>

## STATISTICAL ANALYSES

The number, amount, and proportion of daily energy contributed by all LND foods and subgroups were computed, and each LND food variable was examined as weighted tertiles. The mean daily energy, percentage of energy from macronutrients, and

# Run competition

(<https://coleridgeinitiative.org/richcontextcompetition>)

Training Computing Connecting Rich Context Resources Events About

## Rich Context Competition

### PROBLEM DESCRIPTION

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

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 research methods and 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 both the scientific methods and fields used in the analysis and the research fields.

### COMPETITION SPECIFICS

The competition has two phases (details below).

### PARTICIPANT INFORMATION

- Problem Description
- Competition Goal
- Competition Specifics
- Sponsors
- The Bigger Picture
- Competition Schedule
- How to Participate
- Remuneration
- Judges
- Program Requirements
- First Phase
- Second Phase
- Competition Terms And Conditions
- Teams

## Hold breath...

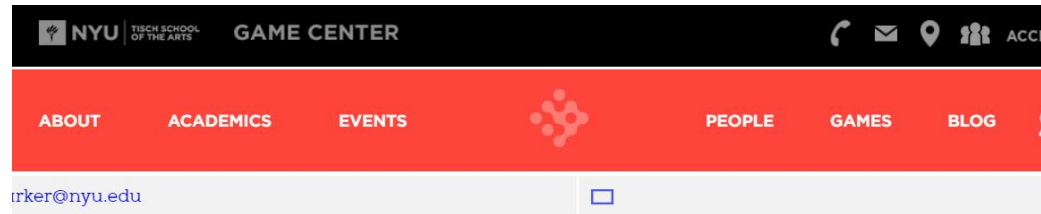
- Twenty teams
  - Half research labs
  - Half students
- Eight countries
  - US
  - Korea
  - Germany
  - Russia
  - Singapore
  - Indonesia
  - Spain
  - India



# Engagement



# NYU Game Center



## Matt Parker



MATT PARKER

### COURSES TAUGHT:

- DESIGNING FOR IMPACT: ICIVICS
- DESIGNING FOR MUSEUMS: AMNH
- CODE LAB 0
- CODE LAB 2
- CODE LAB 1
- THESIS 1
- GAMES 101

Matt Parker is a game designer, teacher, and new media artist. His work has been displayed at the American Museum of Natural History, Brooklyn Academy of Music, SIGGRAPH Asia, the NY Hall of

## Alexander King (MFA '17)



ALEXANDER KING (MFA '17)

Alexander King is an independent developer, freelance game designer and consultant. His work centers on data-driven design and simulation, and his games have been featured in festivals like ALT.CTRL.GDC and IndieCade. Alexander has an MFA in Game Design from the NYU Game Center. Before working in games, Alexander was an analytics consultant working in finance and ecommerce. Now the economies he models are largely fictional ones.

# A Game Designer's View of Engagement

## BACKGROUND

### What exists

- A large database for researchers and lawmakers to access.

### The problem:

- Data is difficult to parse and utilize effectively.
- Users are not encouraged to interact with or contribute to data.

### Project goals

- Have researchers establish context and meaning to datasets.
- Enable users to easily find related datasets and research.
- Create an effective system which encourages researchers to add metadata and organizational system to datasets.

# Jupyter

## 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



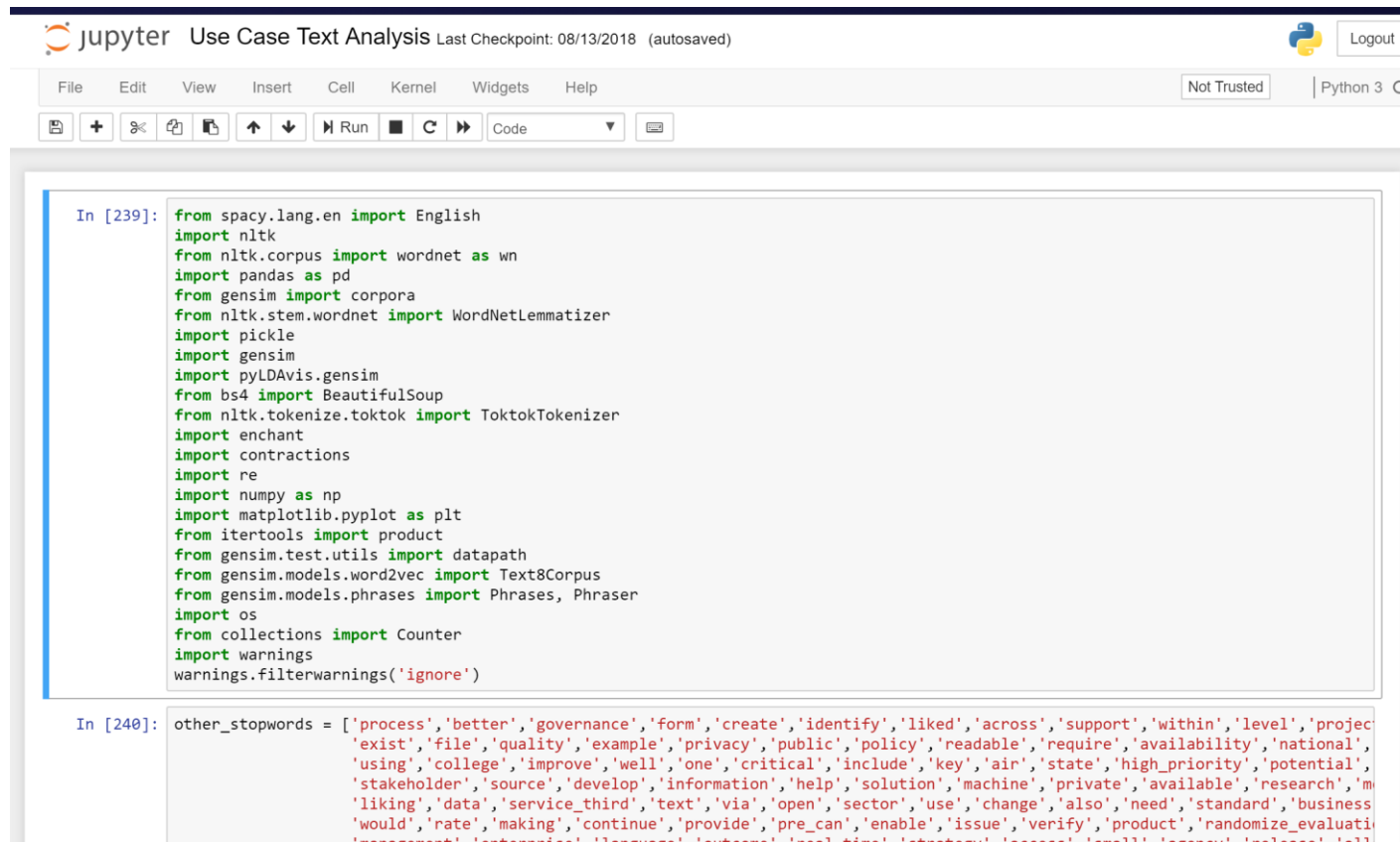
Alfred P. Sloan  
FOUNDATION

SCHMIDT **FUTURES**



Overdeck Family  
Foundation

# The role of Jupyter Notebooks



The screenshot shows a Jupyter Notebook titled "Use Case Text Analysis" with a last checkpoint of "08/13/2018 (autosaved)". The interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help), a status bar (Not Trusted, Python 3), and a toolbar with icons for file operations and execution. The notebook content consists of two code cells. The first cell, labeled "In [239]:", contains a series of import statements for various Python libraries used in text analysis, including spacy, nltk, pandas, gensim, and matplotlib. The second cell, labeled "In [240]:", shows the definition of a list of stopwords for text processing.

```
In [239]: from spacy.lang.en import English
import nltk
from nltk.corpus import wordnet as wn
import pandas as pd
from gensim import corpora
from nltk.stem.wordnet import WordNetLemmatizer
import pickle
import gensim
import pyLDavis.gensim
from bs4 import BeautifulSoup
from nltk.tokenize.toktok import ToktokTokenizer
import enchant
import contractions
import re
import numpy as np
import matplotlib.pyplot as plt
from itertools import product
from gensim.test.utils import datapath
from gensim.models.word2vec import Text8Corpus
from gensim.models.phrases import Phrases, Phraser
import os
from collections import Counter
import warnings
warnings.filterwarnings('ignore')
```

```
In [240]: other_stopwords = ['process', 'better', 'governance', 'form', 'create', 'identify', 'liked', 'across', 'support', 'within', 'level', 'projec
'exist', 'file', 'quality', 'example', 'privacy', 'public', 'policy', 'readable', 'require', 'availability', 'national',
'using', 'college', 'improve', 'well', 'one', 'critical', 'include', 'key', 'air', 'state', 'high_priority', 'potential',
'stakeholder', 'source', 'develop', 'information', 'help', 'solution', 'machine', 'private', 'available', 'research', 'm
'liking', 'data', 'service_third', 'text', 'via', 'open', 'sector', 'use', 'change', 'also', 'need', 'standard', 'business
'would', 'rate', 'making', 'continue', 'provide', 'pre_can', 'enable', 'issue', 'verify', 'product', 'randomize_evaluati
'management', 'antennae', 'language', 'outcome', 'real_time', 'staple', 'access', 'small', 'agency', 'balance', 'all
```

# A computer scientist's view of engagement

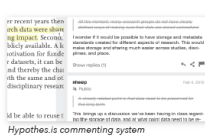
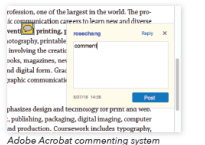
## Designing Comments

By Rose Chang, Meredith Granger, Alena Mueller, and Taka Shimokoba  
Thanks to Tim George, Brian Granger, and Ana Ruvalcaba

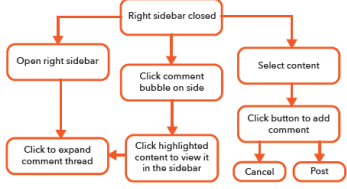
JupyterLab is adding real-time collaboration features to enable multiple users to view and edit documents at the same time. In the context of real-time collaboration, a commenting and annotation system becomes important to enable effective remote teamwork.

### 1 Background Research

The first step of any design is researching similar existing systems. We read W3C's Web Annotation Working Group documents to help us build a conceptual model and vocabulary. In particular, we examined users' feedback about existing systems to discover the advantages and issues of each.



We planned by making an **interaction map**:  
What are all the things a user might want to do?



**Why it's needed**  
Users of JupyterLab will soon be able to edit documents in real-time. In these situations, users often need to leave contextual comments and annotations for themselves and others.

**What it should be**  
Our commenting system needs to work in notebooks, datasets, text files, and other formats.

### 2 Defining Why and What

### 3 Answering Big Questions

As we planned interactions, big questions started to come up. What does it mean to resolve a comment? How is that different from deleting? What happens on a past version of the document? Here are some of the rules we came up with:



- Time and History**
  - The document changes in real-time and only comments on current content show.
  - If a previous version of the document is viewed, only the comments from that version display.
  - When content is deleted, comments on it are only visible in past versions.
  - In previous versions of the document, comments are shown in the state (resolved/open) that they were.



- Delete and Resolve**
  - Delete removes a comment from all versions.
  - Resolving a comment hides it unless viewing resolved comments.
  - Resolved comments only show if their content exists.
  - Anyone can resolve or reopen a comment thread.

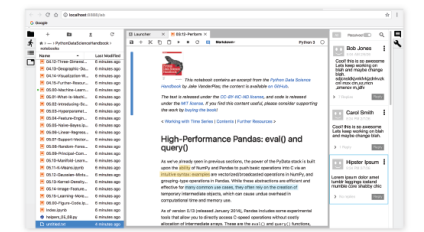


- Comment Threads**
  - If multiple comments are in the same area of content, one click opens all threads.
  - In the comment panel, users click a thread to expand it.
  - Only the first comment is shown in unfocused threads.

### 5 Evaluating Our Work

A major idea in early iterations was having three views of the comment panel. Users can view all threads, a group of threads, or an individual thread.  
**But typical users are going to have much simpler mental models.**

The cycle of evaluating our work and building prototypes that meet JupyterLab's standards for usability is going to continue for much longer. We've started brainstorming prototypes to match that.



We used Figma, a art application, to make our prototypes. Figma is great for collaboration and allowed us to link one screen to another interactively. We created a functional prototype to test users without any coding.

### 4 Building Prototypes



# The Future



# Build a contributing c



## Transform Data Use

# A Locally Based Initiative to Support People and Communities by Transformative Use of Data

SHARE

JULIA LANE, DAVID C. KENDRICK, DAVID T. ELLWOOD

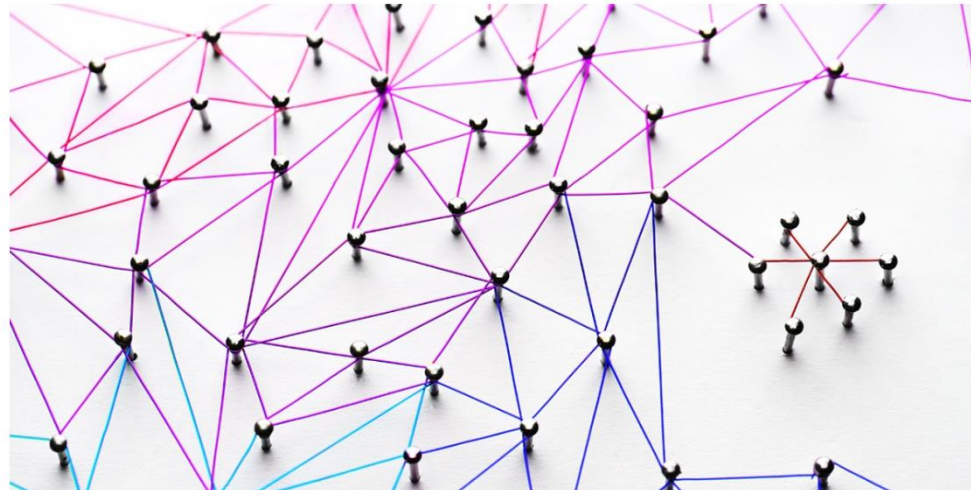
The data revolution is transforming how executives manage operations and businesses deliver goods and services. Yet when it comes to helping people escape poverty, the revolution has barely begun.

DOWNLOAD THIS PAPER

Full Idea Paper

Summary >

SIGN UP FOR OUR NEWSLETTER



BESTSELLER

forms but entertains. You'll never be the same way again." *Religious and Invisible Influence*

e

ur

OF AMAZON, MICROSOFT, AND GOOGLE

alloway

## Takeaways

- Massive supply of new granular data
- New demand - local/city/regions
- Data science offers engagement and feedback tools
- Statistical skills critical
- Need to rethink data infrastructure

Help us and be a beta tester! Email me at [julia.lane@nyu.edu](mailto:julia.lane@nyu.edu)

# Comments welcome

- [Julia.lane@nyu.edu](mailto:Julia.lane@nyu.edu)
- <https://coleridgeinitiative.org>