

# Positioning Libraries to Support Data Science

**Tim Dennis, Director, UCLA Library Data Archive**

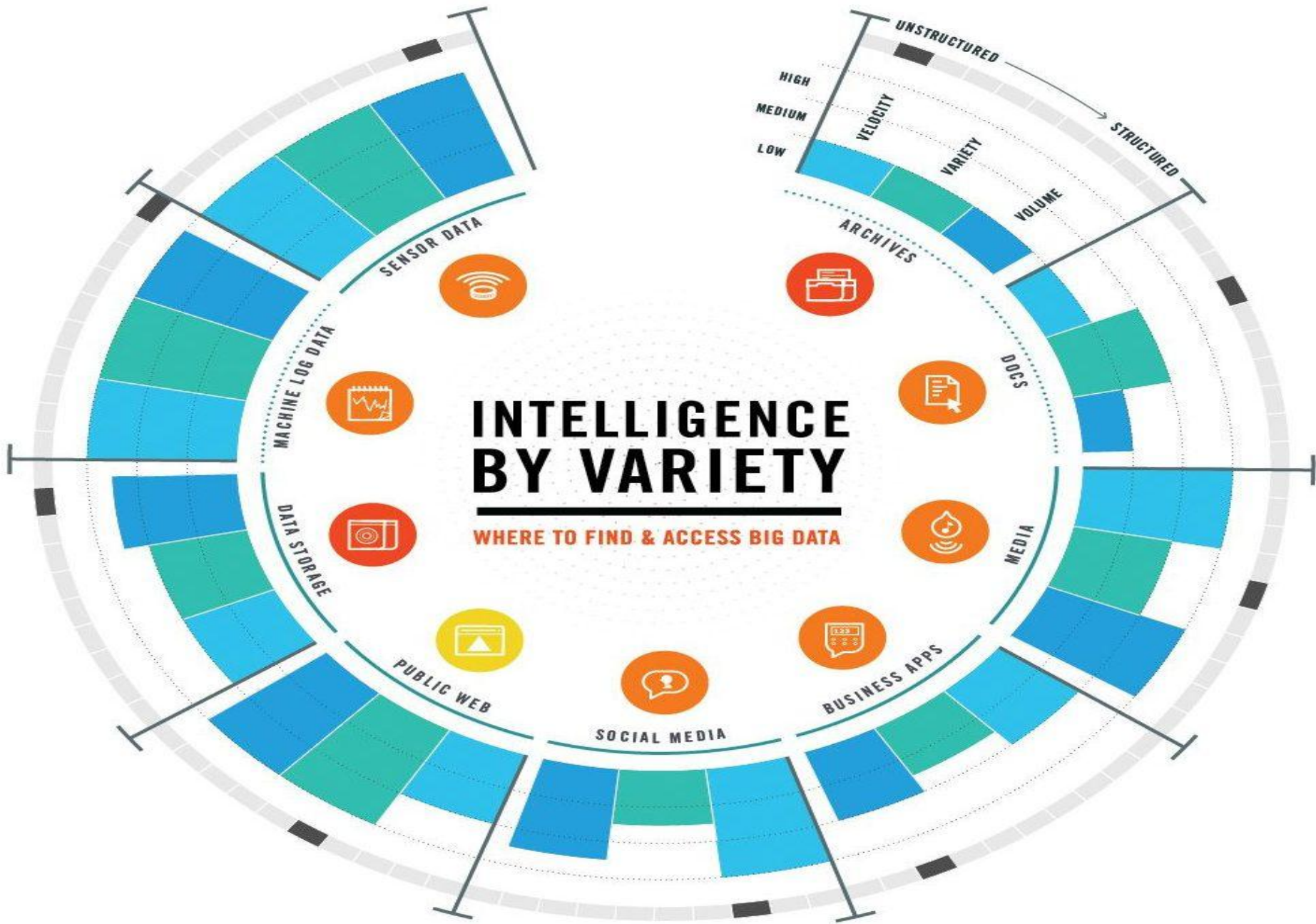
LiASA 2017 - 18th Annual Conference  
October 5, 2017

# What is Data Science?

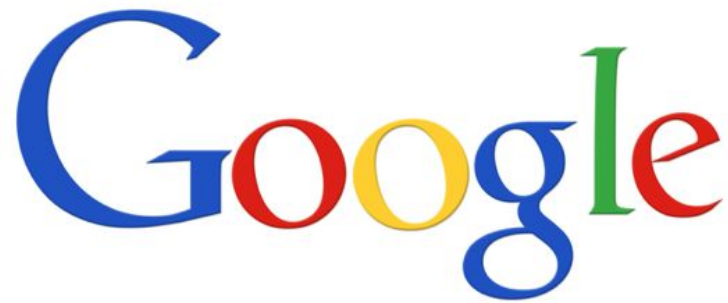


“[Data Science is] the ability to extract knowledge and insights from large and complex data sets”

-- Former US Chief Data Scientist - Dr. DJ Patil <https://goo.gl/3Y6es5>



These companies collect and profit from data



FRESHLY MINED!  
YOUR PERSONAL DATA!

Only  
2/5¢

# YOU!

MORE INVASIVE THAN EVER! YOUR SPENDING HABITS!  
YOUR INCOME! YOUR FACEBOOK PROFILE! YOUR NEIGHBORS!  
YOUR SOCIAL SECURITY NUMBER! YOUR MORTGAGE!  
YOUR BELIEFS! AND SO MUCH MORE!

NET WT. 1.4 LBS.

# Industry Context of Use

- Need to make sense of huge swaths of increasingly dynamic and complex data in order to monetize
- Need people with skills to:
  - Munge, clean and organize data
  - Create algorithms to make data actionable
  - Store data efficiently and intelligently

“

Data is the new oil”

Clive Humby



What does that mean for  
academia?



# Impact on Academia

- Many academics with higher order math/CS skills are poached (or start start-ups)
- Market pressures put on Universities to produce more grads with data skills
- Research methods in many fields transformed by computationally enabled research and access to “Big Data”

## DISCIPLINES

<< All Disciplines

Computer Science & IT (1480)

Business Information Systems (116)

Computer Sciences (701)

**Data Science & Big Data (82)**

Geographical Information Systems (GIS)  
(78)

Health Informatics (90)

Human Computer Interaction (20)

Informatics & Information Sciences (352)

IT Security (100)

Video Games & Multimedia (30)

Web Technologies & Cloud Computing (47)

## LOCATION

United Kingdom (75)

United States (82)

Show **All**

Filters: United States × Data Science & Big Data × Clear All

## Analytics, M.Sc.

American University Washington DC Kogod School of Busin... | Washington, D. C., Washington, D.C., U

 44,369 EUR / year ⓘ  1 year  On campus  English (Take IELTS Test)



The Master of Science in Analytics program (MSAn) offered by the Kogod School of Business at American University Washington DC prepares students to be experts in data analysis and to make informed organizational decisions, and to solve dynamic business problems. The program

Add to Comparison

## Graduate Pathway in Data Analytics Engineering, Pre-Master

INTO George Mason University | Fairfax, Virginia, United States

 29,759 EUR / year ⓘ  9 months  On campus  English (Take IELTS Test)

# Undergrad Data Science Education

In spring 2016, UC Berkeley's first Foundations of Data Science course attracted around 300 students. This semester (2017), nearly 1,000 have enrolled...Across the UC system, campuses are quickly adding data science programs in response to soaring workplace demand.

--By Isha Salian, [www.sfchronicle.com](http://www.sfchronicle.com)



**Cathryn Carson**

@CathrynCarson

Following



First day of class in Foundations of  
[#Datascience](#) for [@BerkeleyDataSci](#) - all  
materials at [data8.org](#)



10:23 AM - 23 Aug 2017

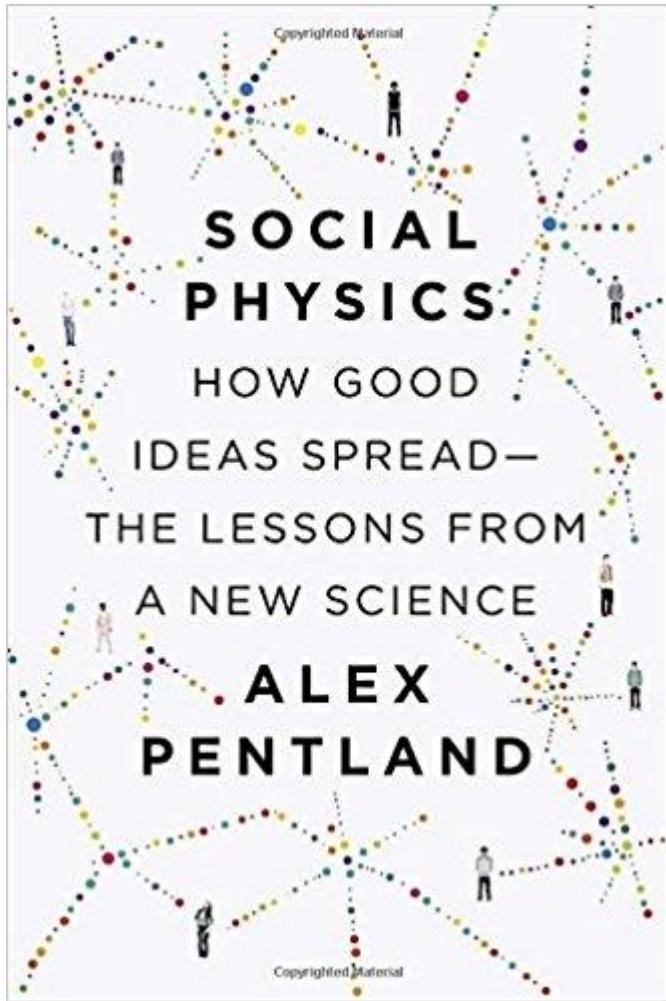
14 Retweets 47 Likes



Data Science enables new research questions through new data sources:

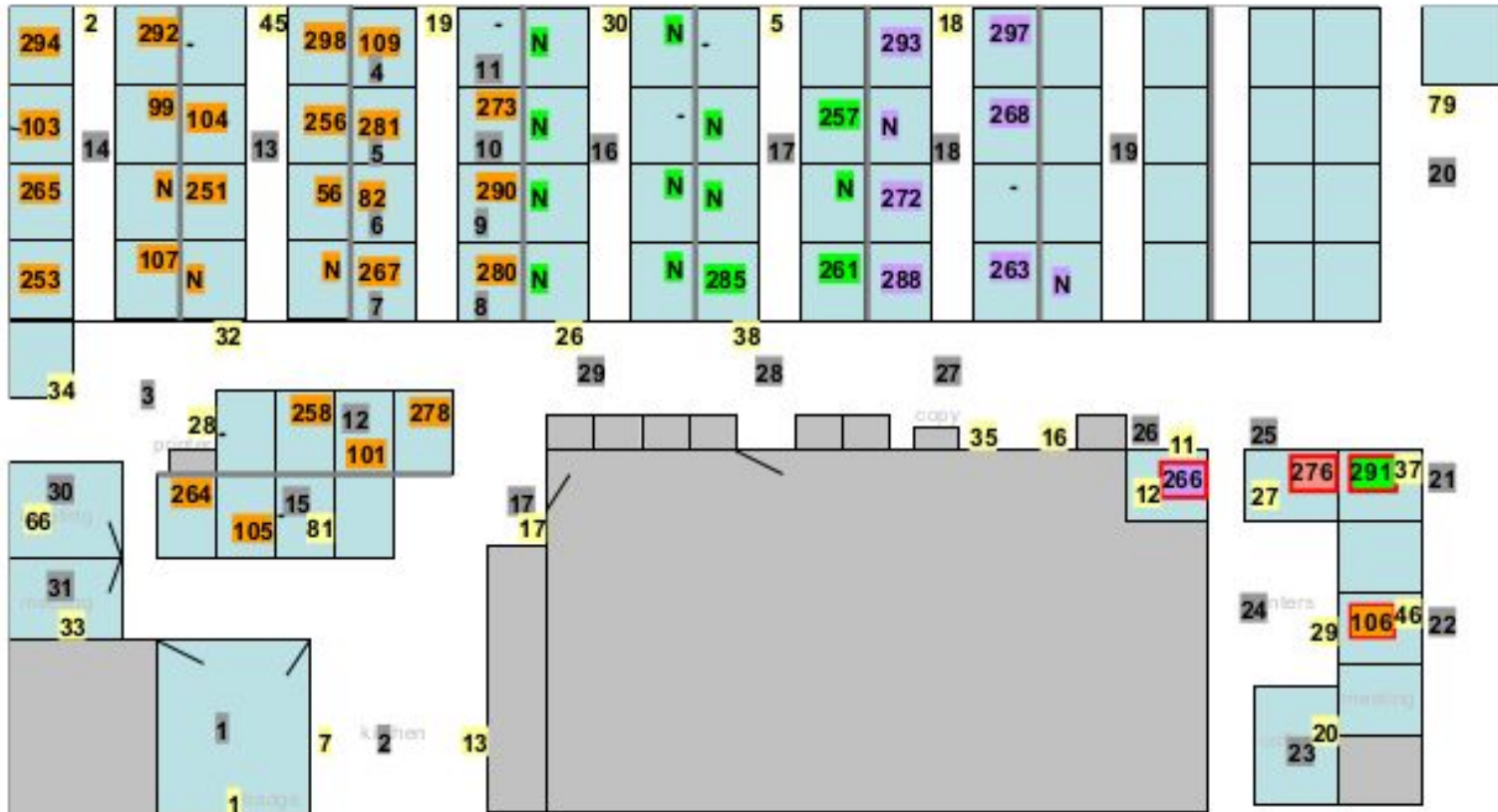
- Cell phone data
- Social media data
- Sensor data
- Textual data

# Example: Sensor Data



- Modelling humans like we do ants
- Used sensors in a business office
- Sensor badges had mic, accelerometer, and geolocation





All except workers' name

- Worker
- Configuration
- Coordinator
- Pricing
- Manager
- Base station
- RSSI

# Example: Sensor Data

## Findings:

- Individual's style of speaking (rate, emotional tone, mirroring, turn-taking predicts 30% of the variance in performance (raises)
- Can determine influence networks and nodes
- Exposure to peer behavior influences individual behavior



# Example: Sensor Data

## Data Collection Sites:

- Dorm network
- Grad school network
- Thin data on cell phone data from Ivory Coast

Pentland asserts “we are scratching the surface” on this type of research

## How Censorship in China Allows Government Criticism but Silences Collective Expression

GARY KING *Harvard University*

JENNIFER PAN *Harvard University*

MARGARET E. ROBERTS *Harvard University*

**W**e offer the first large scale, multiple source analysis of the outcome of what may be the most extensive effort to selectively censor human expression ever implemented. To do this, we have devised a system to locate, download, and analyze the content of millions of social media posts originating from nearly 1,400 different social media services all over China before the Chinese government is able to find, evaluate, and censor (i.e., remove from the Internet) the subset they deem objectionable. Using modern computer-assisted text analytic methods that we adapt to and validate in the Chinese language, we compare the substantive content of posts censored to those not censored over time in each of 85 topic areas. Contrary to previous understandings, posts with negative, even vitriolic, criticism of the state, its leaders, and its policies are not more likely to be censored. Instead, we show that the censorship program is aimed at curtailing collective action by silencing comments that represent, reinforce, or spur social mobilization, regardless of content. Censorship is oriented toward attempting to forestall collective activities that are occurring now or may occur in the future—and, as such, seem to clearly expose government intent.

# Example: Big Text Revolution

- Digital text encodes much of human activity, including the past
- Often relatively cheap to purchase
- Def. “Big Data”
- How to parse it and make sense from it is a growing research field

How can Libraries engage?

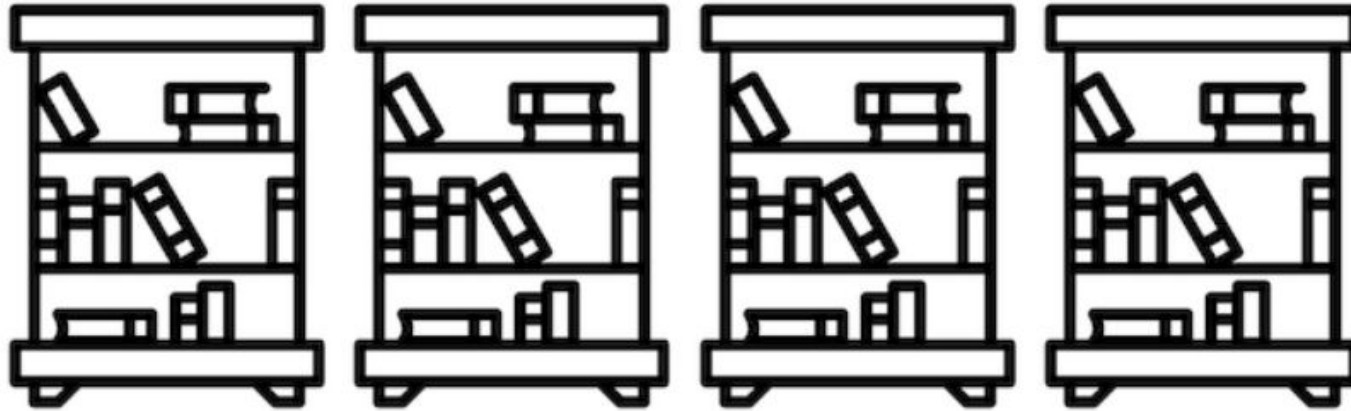
# Two areas where libraries can jump in

1. Collections
2. Data services

# Collections

1. Purchase licensed data sources  
(Corpora of text, social media data, etc.)
2. Treat purchased data as a first class citizen of your collection
3. Collections as data

# Collections as Data



**Collections as Data** is an Institute of Museum and Library Services supported effort that aims to foster a strategic approach to developing, describing, providing access to, and encouraging reuse of collections that support computationally-driven research and teaching in areas including but not limited to Digital Humanities, Public History, Digital History, data driven Journalism, Digital Social Science, and Digital Art History.

# Example: Collections as Data



- Free Speech Movement collection - artifacts from FSM digitized
- The Library & Research IT created API
- Held competition (Hackathon) for students to create novel interfaces for it



# FREE SPEECH MOVEMENT DIGITAL ARCHIVE

Enter your keyword

Search the Archive!

Text  Image  Audio



September 14, 1964

## Free Speech Movement

A short timeline of the events of the FSM



SEPTEMBER  
14, 1964

Dean Tucker  
writes to student  
political groups

Jack Weinberg is  
arrested

President Clark Kerr  
addresses faculty and  
student body

Free Speech Movement

Free Speech Movement  
is formed

March to Regent's  
Building

Dean Tucker's letter to  
student political groups

First Meeting of the Campus  
Committee on Political  
Activity (CCPA)

Shike - Occupied  
Sprout Hall

SEPT.

OCT.

NOV.

DEC.

# Data Services

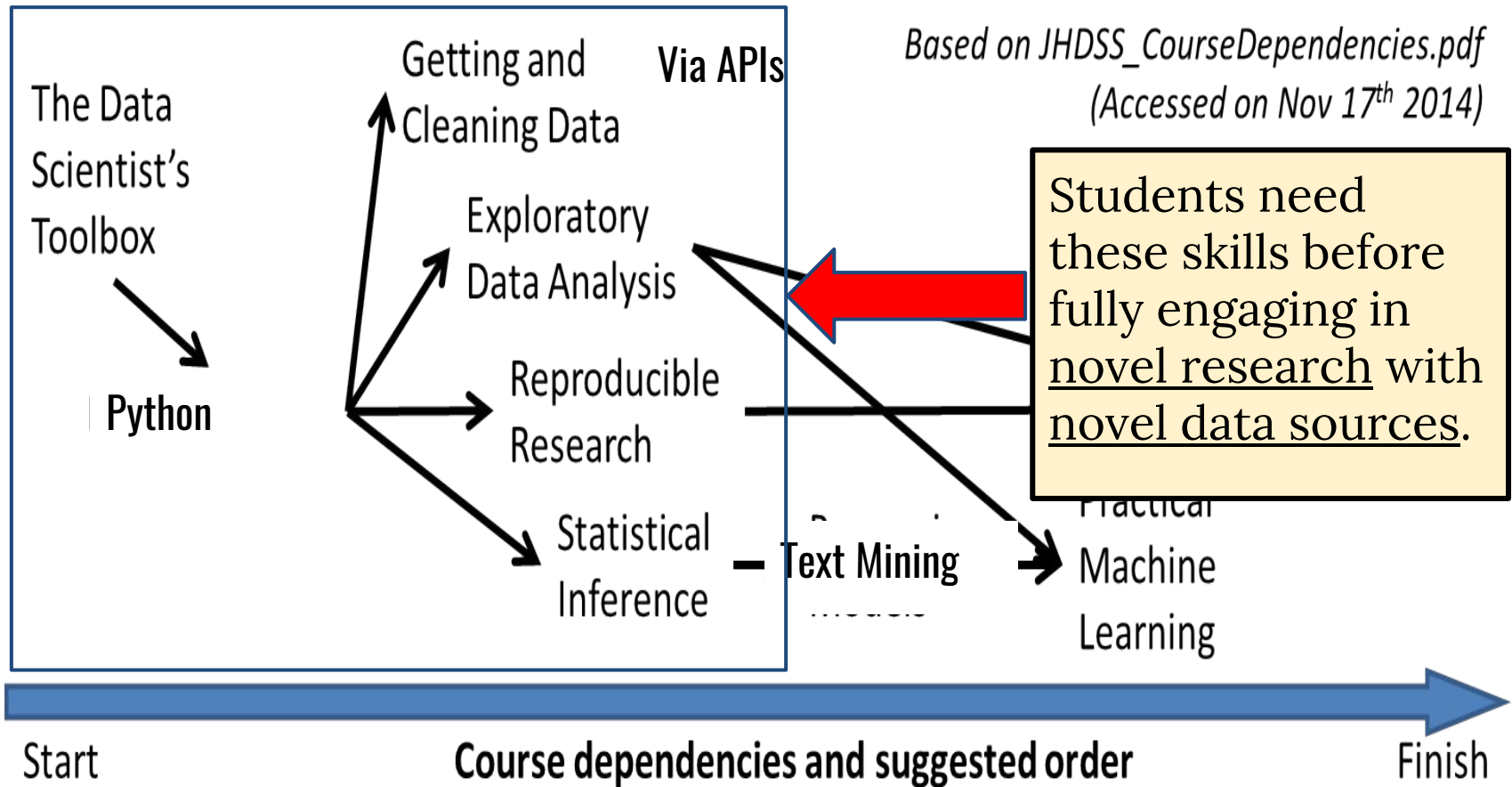
1. Finding data
2. Help with research data needs:
  - management, backup and storage
  - identify data repositories
  - educating and preparing data for sharing and publishing
  - metadata and standards support
3. Provide consultation service on using data
4. Data Instruction

# A note about metadata

**Metadata** is crucial to your Big Data project's overall success and to your enterprise data architecture organization.

--Data Science Central  
<https://goo.gl/u2ocJA>

# New Research Depends on Basic CS Skills



# Training

- If data training doesn't exist on your campus, consider starting a training program for graduate students.
  - Use Data and Software Carpentry's model to start with high quality workshops and build a local network
- If a data training program exists, partner with whoever is running it and offer to help or support
- Start meetup style groups (HackyHour, Python and Pizza, Data and Donuts)



# Software Carpentry

Software development best practices

Domain agnostic

- Command line
- Version control with git
- Programming in Python or R



# Data Carpentry

Working effectively with data  
includes domain-specific content

- Data organization
- Data cleaning
- Data analysis and visualization in R or Python

# Ok, why else should I learn this stuff?

- So you aren't replaced by a robot
- Learning data and software skills is very marketable (ask the students at your school)
- It's fun!
- It will help your researchers





# Library Carpentry

Working effectively with data using  
best practices

- Basic computation skills
- Versioning and collaboration through Github
- Data analysis and cleanup through OpenRefine
- If interested, join the community call Oct. 10, 4pm. Details: <https://goo.gl/PxnQCG>

# Data Science

Give[s] us new ways of grasping patterns, collectivities, and systematic effects that remain invisible to us without statistical and computational tools; of understanding the linkages from data to knowledge to decision-making under conditions of uncertainty; of exploiting domain-specific computational possibilities fluidly and reliably and seeking cross-fertilization across them; and of **critically engaging** the constructive and creative possibilities opened up by data collection and computation, as well as their **challenging ethical and social entanglements**.

-- Data Sciences @ Berkeley, The Undergraduate Experience

# Data Literacy for Society

- Libraries (Public and Academic) should be offering advice and education on cyber-security so people can protect their personal data
- As an equity issue they should be introducing users to free and open source replacements for expensive proprietary systems
- They should also be fostering the teaching of code because this is where the jobs of the future will exist
- Providing data literacy training so people can navigate every day society (news, advertisements, etc.)

# Questions and Discussion

Thanks so much!

Questions?

Tim Dennis <[timdennis@ucla.edu](mailto:timdennis@ucla.edu)>

Twitter: @jt14den