

# Data Science with Python

Summer 2016

# Today

- Hello
- Data Science: what is it?
- Course Overview
- Python: check-in on codeacademy
- Data Ingestion
- Jupyter / ipython
- Git / VirtualEnv
- Check-in

# What is Data Science?

- What words come to mind when you think of Data Science?
- What experience do you have with Data Science?
- Why are you taking an Introduction to Data Science Class?

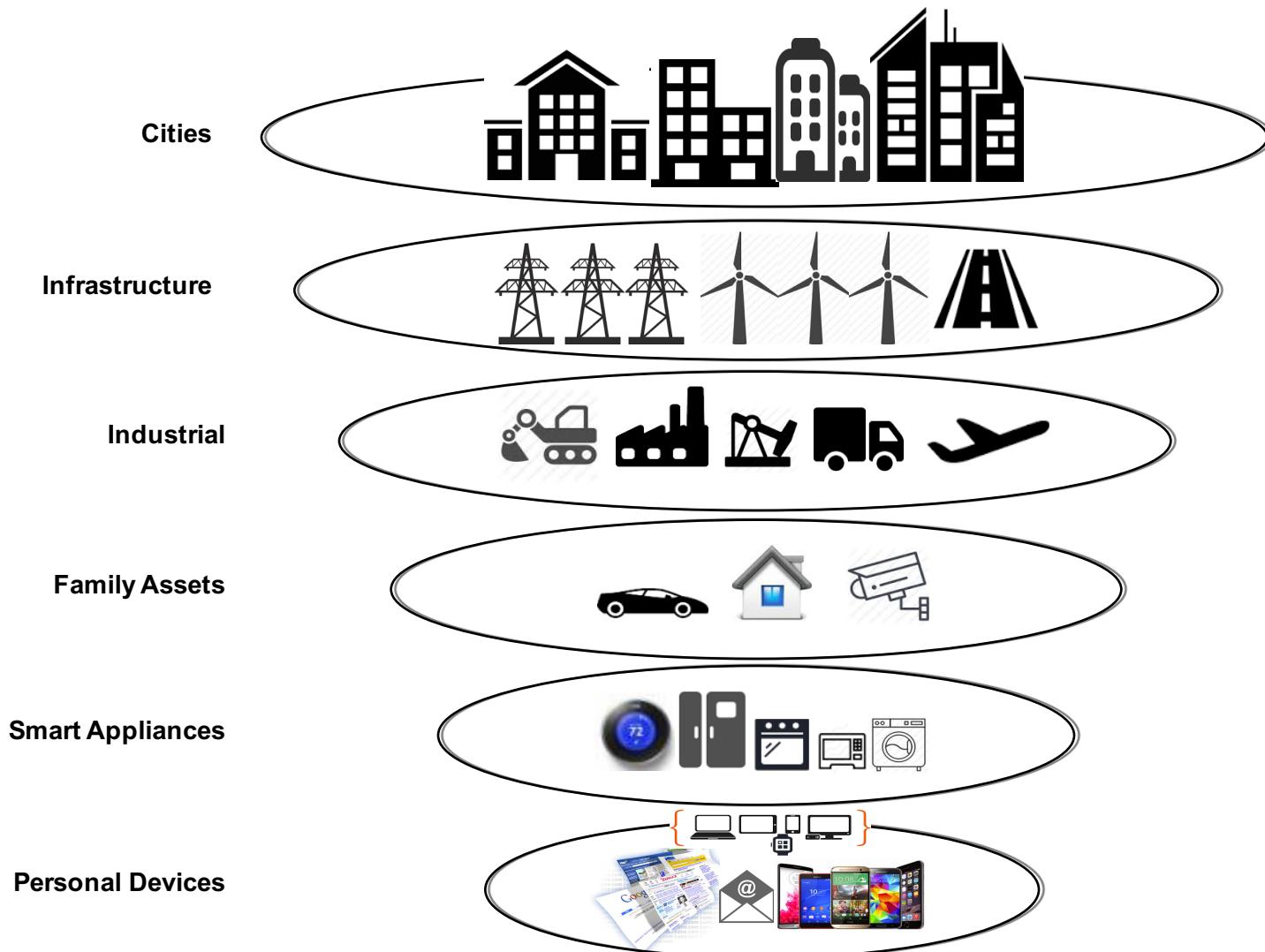
“Data”











HBR.ORG

# Harvard Business Review

OCTOBER 2012  
REPRINT R1210C

SPOTLIGHT ON BIG DATA

## Big Data: The Management Revolution

Exploiting vast new flows of information  
can radically improve your company's  
performance. But first you'll have to  
change your decision-making culture.  
*by Andrew McAfee and Erik Brynjolfsson*

# Data Scientist:

*The Sexiest Job of the 21st Century*



REVO<sup>LUTION</sup>  
ANALYTICS



ClearStory DATA  
GoodData  
+ tableau



webtrends



Google Analytics



Localytics

Spark



hadoop

cloudera

Greenplum

NETEZZA

VERTICA



ORACLE

IBM

TERADATA

eBureau®  
AnalyticsIQ Inc.

POLK

IXI™ SERVICES

causa  
nielsen

Semcasting

# What is Data Science?

- “How Companies Learn Your Secrets” *NYT*, by Charles Duhigg, February 16, 2012
- <http://bit.ly/targetdata>



<http://bit.ly/targetdata>

# What did Target Do?

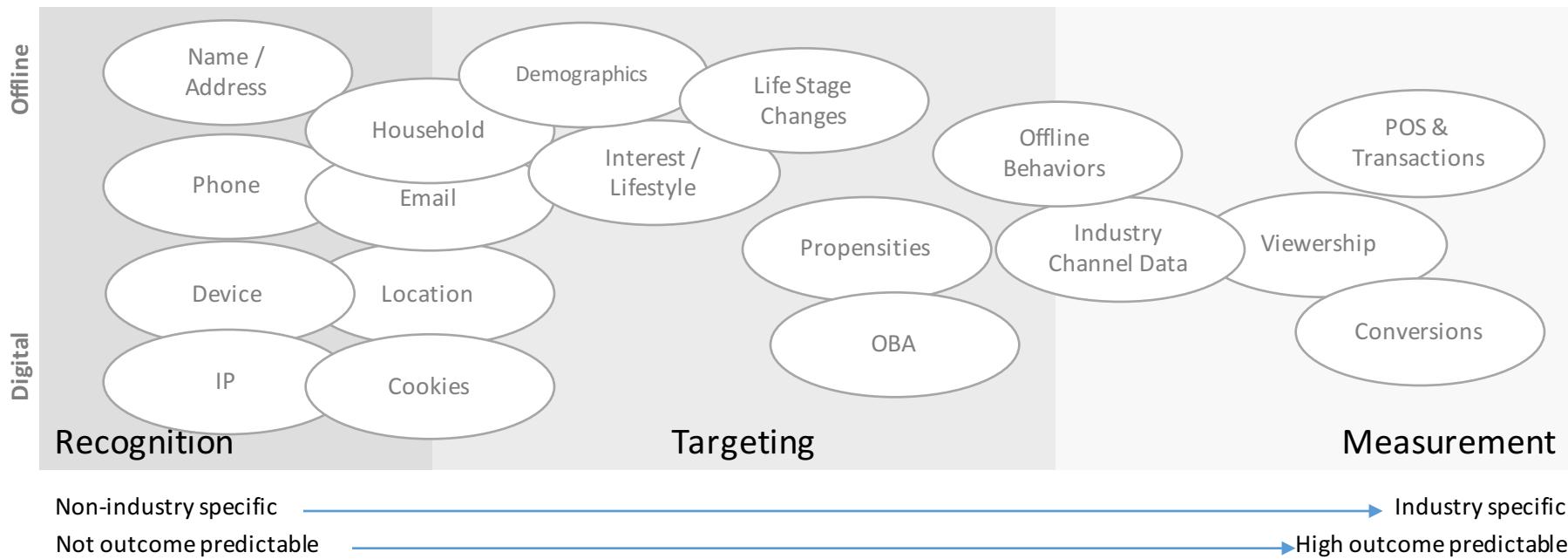
- Mining of data on shopping patterns
  - Specific products purchased
  - Combination of products purchased
  - Combined with demographic and other data
- Psychology and neuroscience
  - Habits:
    - Cue-routine-reward
    - When are habits open to change?

# Lessons from Target

- Yes, Data Science is about mining data
- There are deeper theoretical issues involved in understanding what you find
- Left out of that long article are most of the critical steps that precede the analysis
- In short, Data Science > data mining

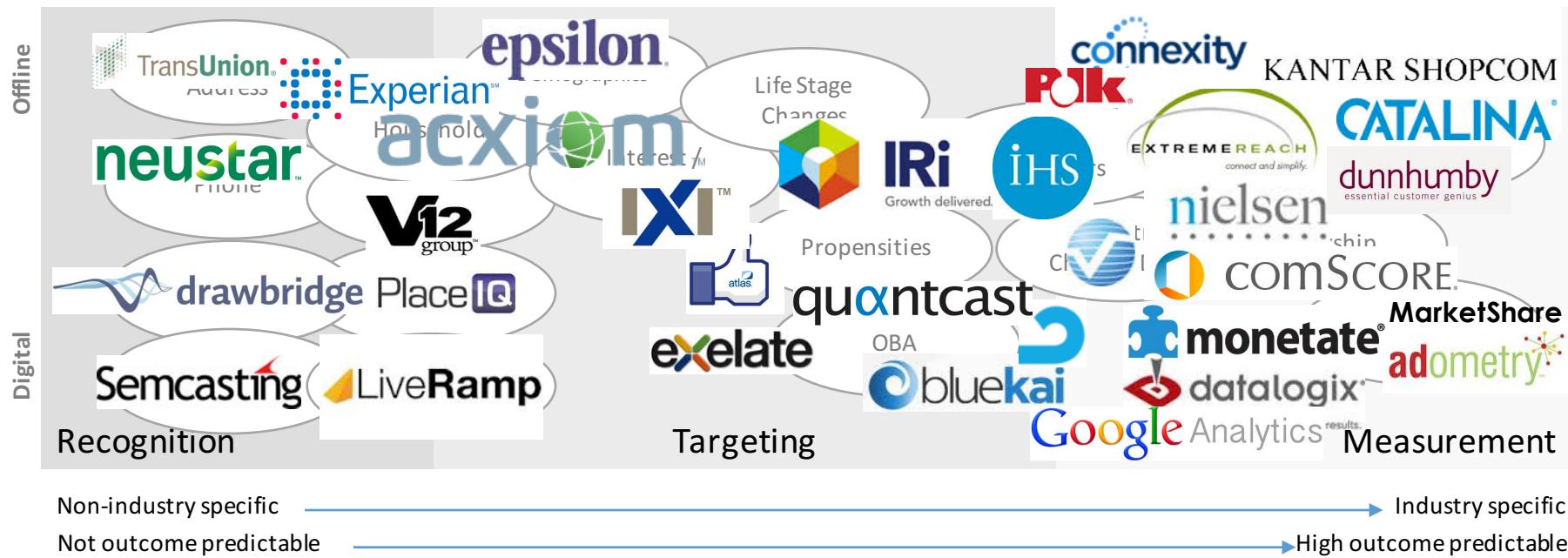
# Consumer data landscape

A framework for understanding consumer data



# Consumer data landscape

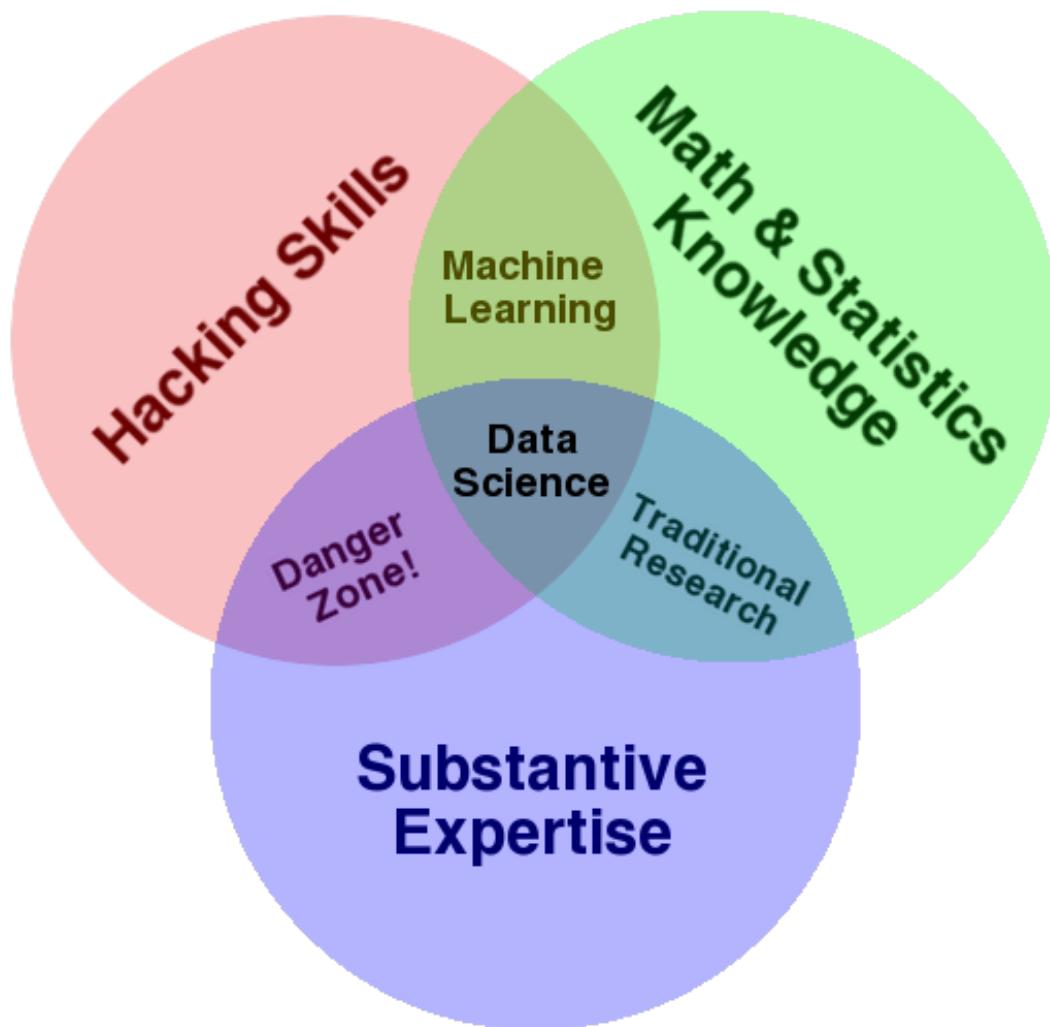
A framework for understanding consumer data



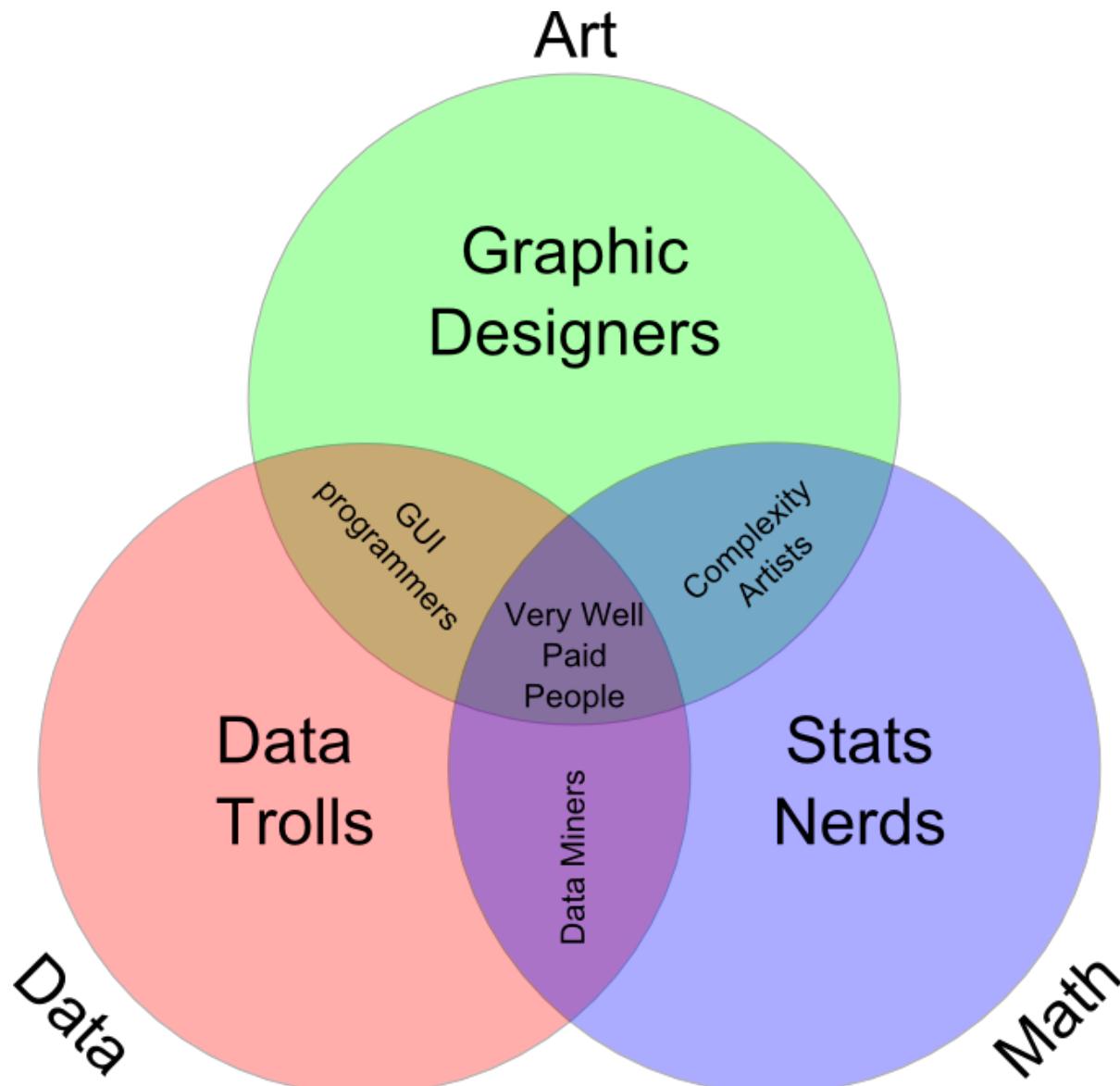
# Definition of Data Science

- There are many, but most say data science is:
  - Broad – broader than any one existing discipline
  - Interdisciplinary: Computer Science, Statistics, Information Science, databases, mathematics
  - Applied focus on extracting knowledge from data to inform decision making.
  - Focuses on the skills needed to collect, manage, store, distribute, analyze, visualize, and reuse data.
- There are many visual representations of Data Science

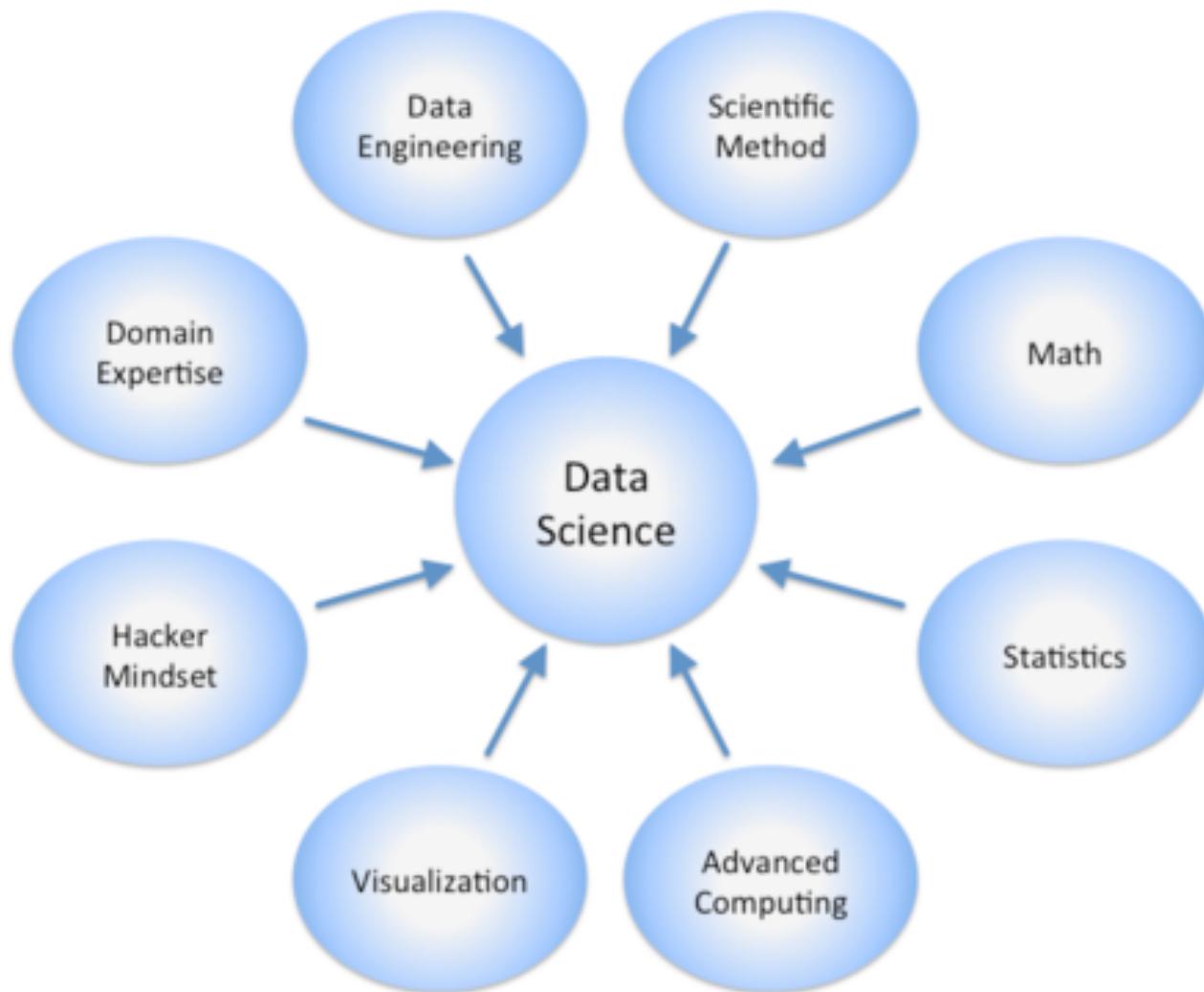
Some definitions link computational, statistical, and substantive expertise.



Other definitions focus more on technical skills alone.



Still other definitions are so broad as to include nearly everything.

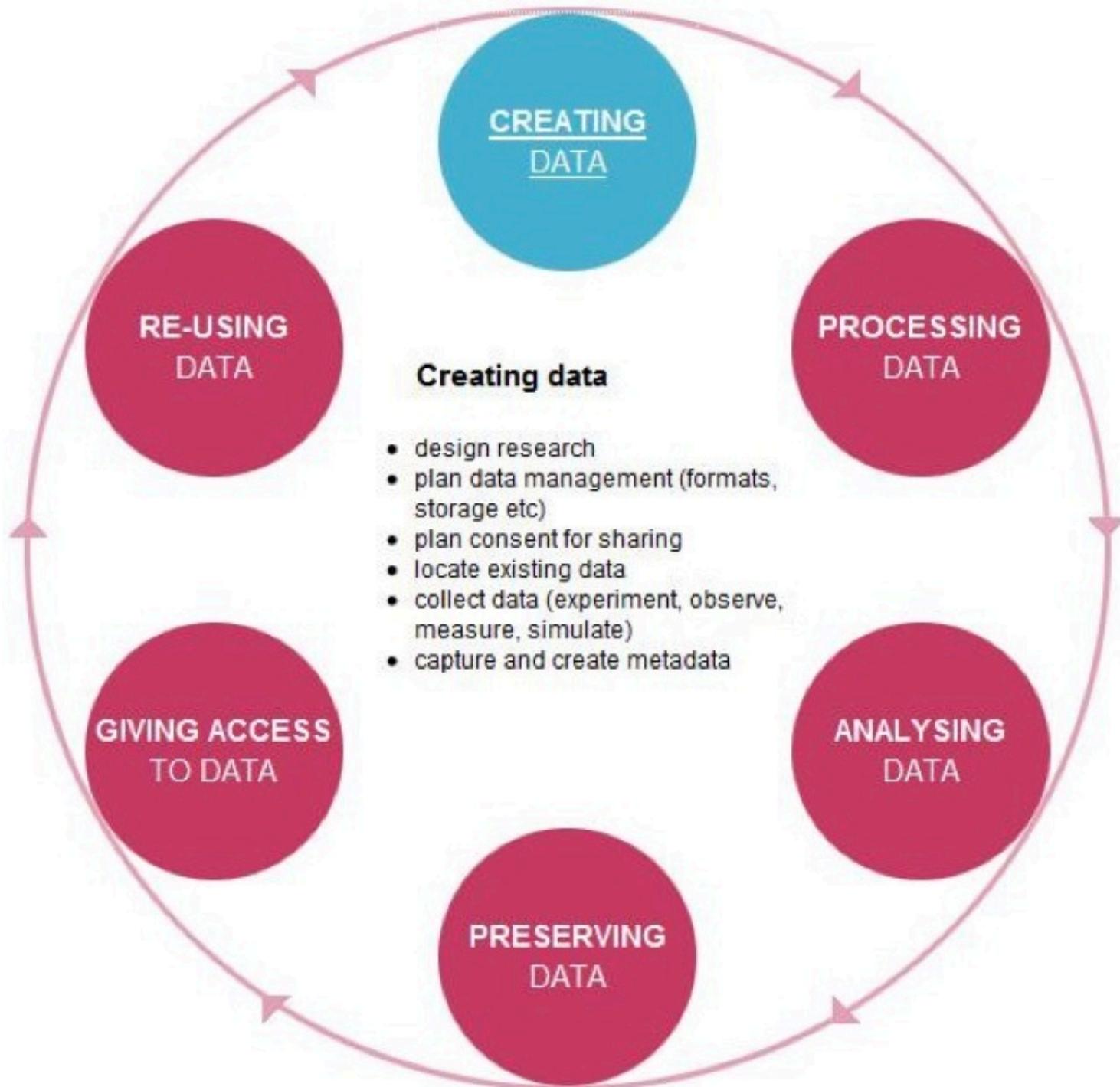


There are many “Word Cloud” representations of Data Science as well.



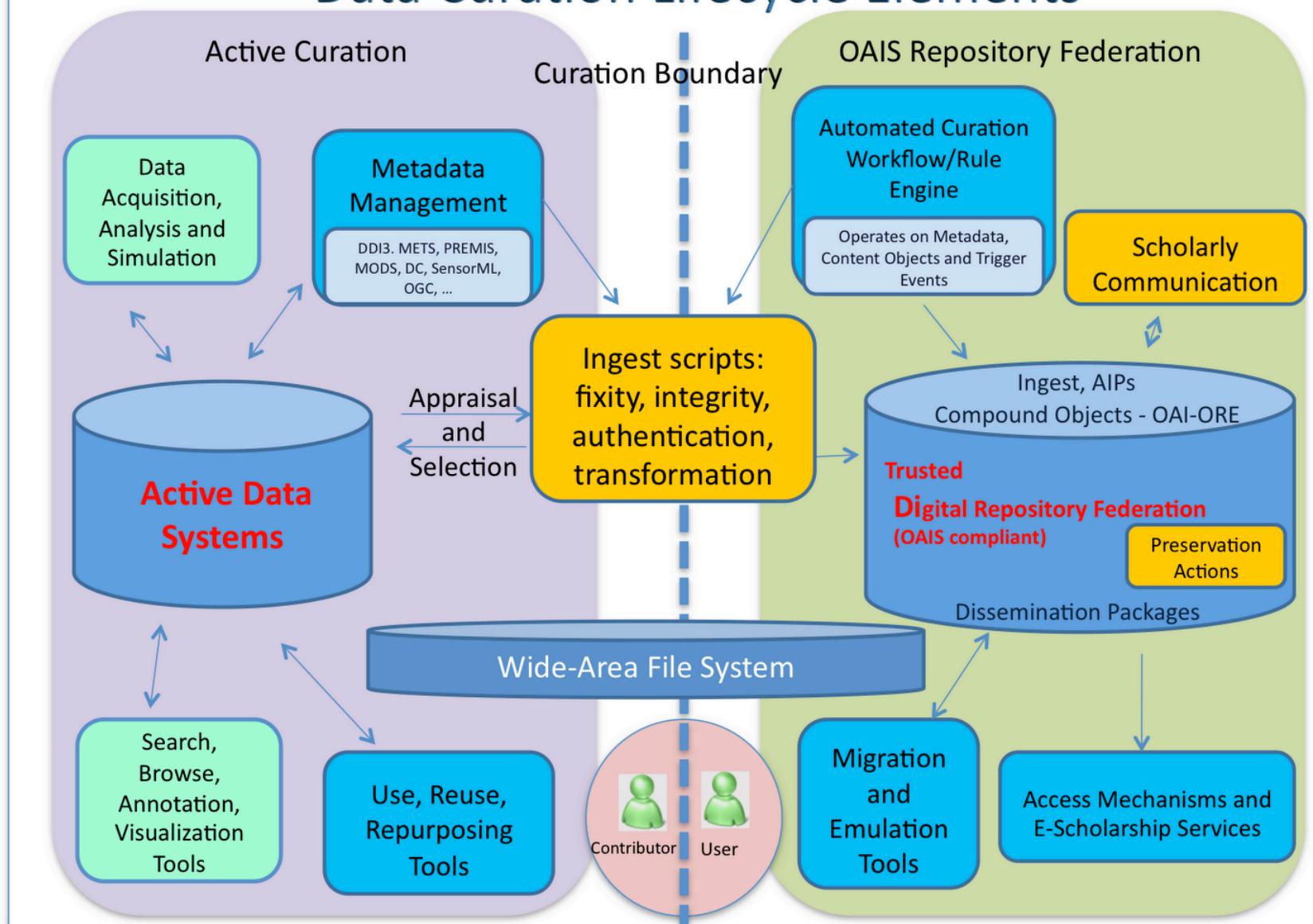








# Data Curation Lifecycle Elements



# What is Missing?

- Most definitions of data science underplay or leave out discussions of:
  - Substantive theory
  - Metadata
  - Privacy and Ethics

# Today

- Hello
- Data Science: what is it?
- **Course Overview**
- Python: check-in on codeacademy
- Data Ingestion
- Jupyter / ipython
- Git / VirtualEnv
- Check-in

Course github site:

<https://github.com/philmui/datascience>

git clone https://github.com/philmui/datascience

# Getting git

Mac:

- brew install git

PC:

- <https://git-scm.com/downloads>

# APPENDIX