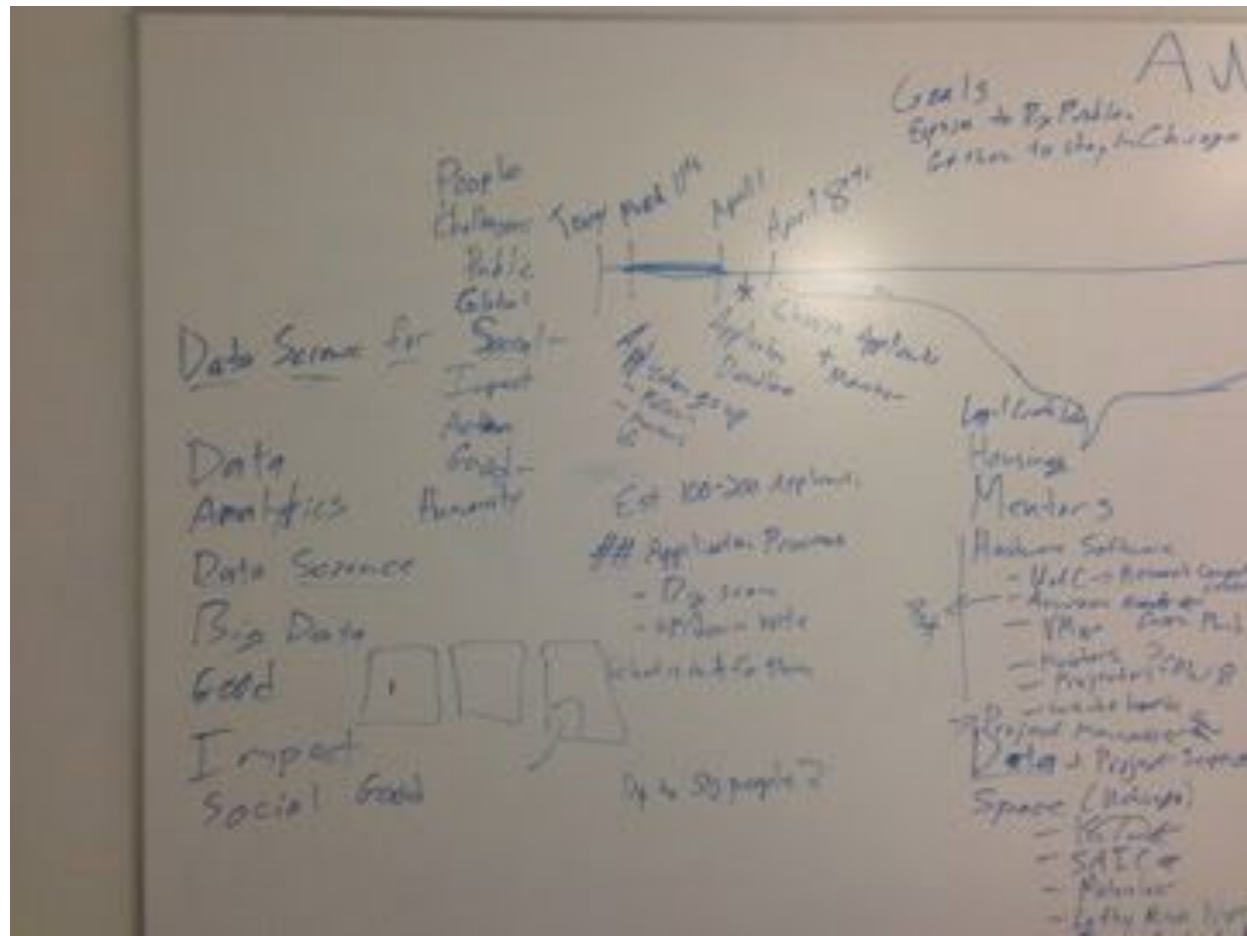# From Classroom to Clinic

Data Science for Social Good Fellowships and the Lessons Data Science Educators Can Learn from the Medical Profession

Matt Gee
National Academy of Sciences December 2018

# Data Science for….

# Goals of DSSG

| Train Fellows | Train Project Partners | Build Community |

Collaborative Projects

Lectures, Talks, Workshops

Events

Open, Collaborative, Ethical

# Step 1: Find Fellows

**600+**
applications
in 3 weeks

**120**
interviews

**36** fellows
selected

Stanford University
University of California-Irvine
Arizona State University
University of Texas-Austin
University of Wisconsin-Madison
University of Texas-Austin
Instituto Technológico Autónomo de México

University of Chicago
University of Illinois-Chicago
University of Michigan
University of Alabama
Georgia Institute of Technology
Carnegie Mellon University
Cornell University
Israel Institute of Technology

University of Maryland
City University of New York
Columbia University
Yale University
Harvard University
Massachusetts Institute of Technology
University of Cambridge
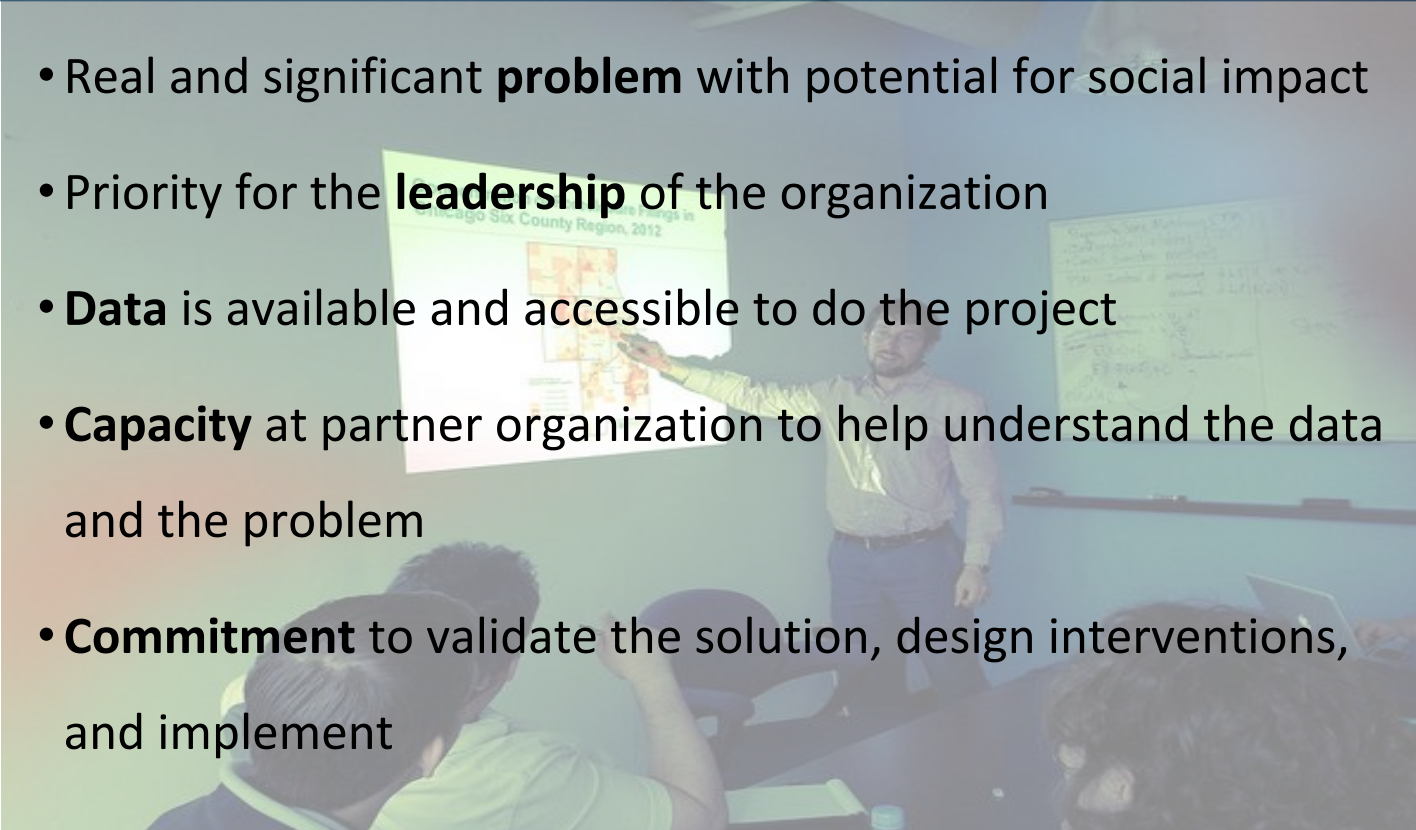


PhD/Postdoc/Professional
Masters
Undergraduate

Computer Science
Social Sciences & Public Policy
Mathematics, Statistics, Physical Sciences

# Step 2: Pick partners

**30+** potential partners

**60** discovery calls and partner interviews

**12** partners

- Real and significant **problem** with potential for social impact

- Priority for the **leadership** of the organization

- **Data** is available and accessible to do the project

- **Capacity** at partner organization to help understand the data and the problem

- **Commitment** to validate the solution, design interventions, and implement

# Step 3: Scope projects

**250 hours** of scoping calls, meetings, and negotiations
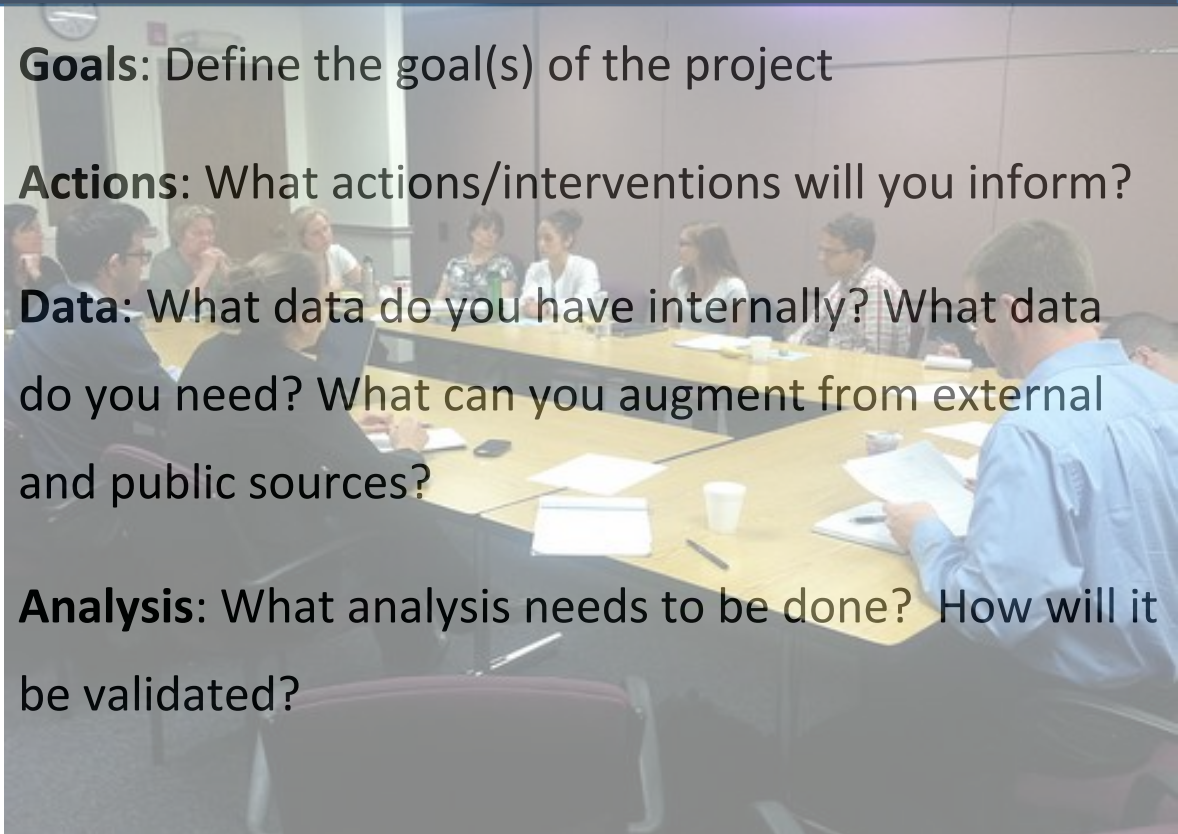
**14** finalized projects

**9** legal agreements signed by the first day of the fellowship
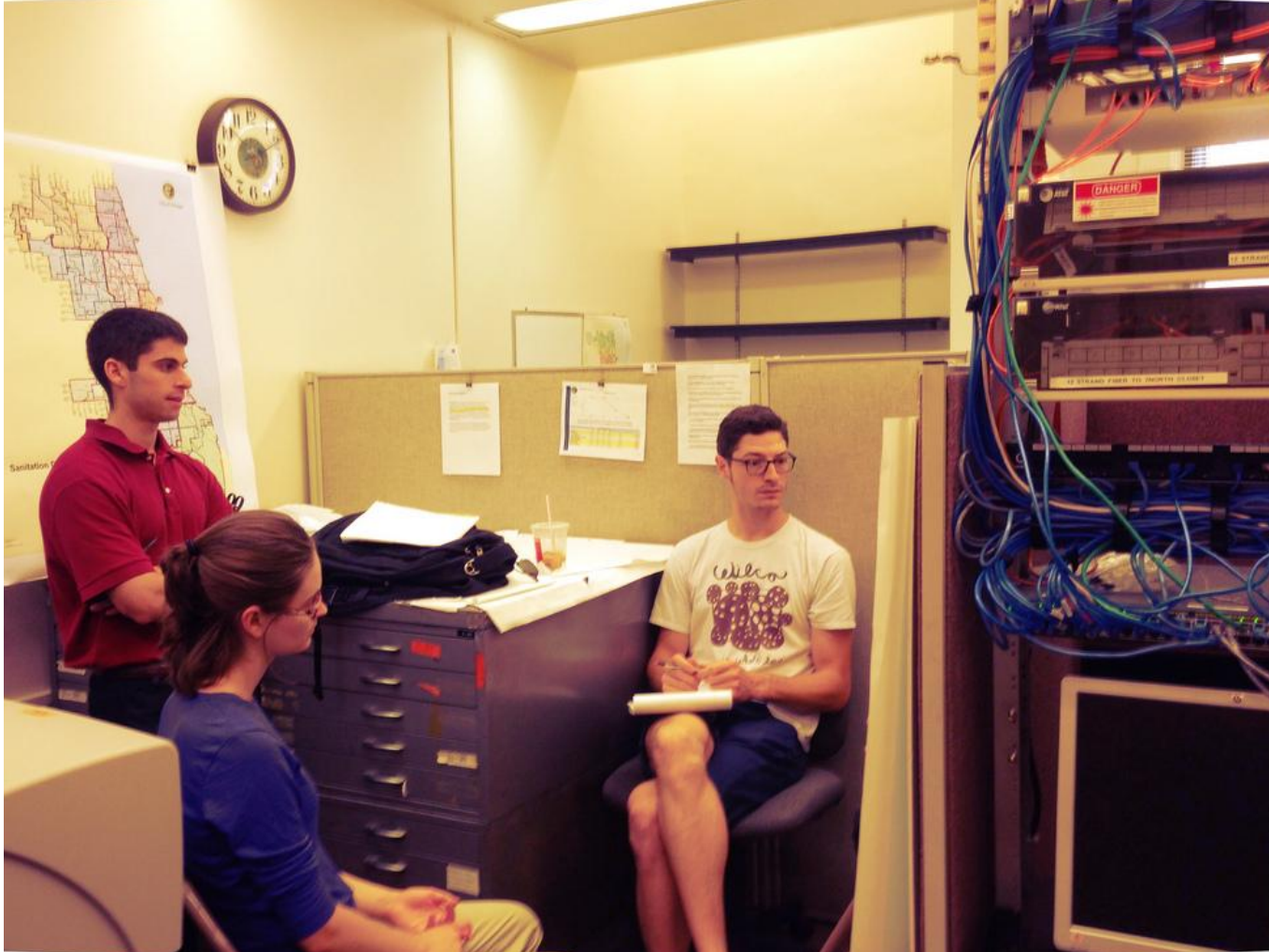
**Goals**: Define the goal(s) of the project

**Actions**: What actions/interventions will you inform?

**Data**: What data do you have internally? What data do you need? What can you augment from external and public sources?

**Analysis**: What analysis needs to be done?  How will it be validated?

**Data Science for Social Good**

"We are used to using data to justify funding decisions. Now we can use data to improve what we do."

Bill Thorland
Nurse Family Partnership

# DSSG Fellowship over the years

| 2013 | 2014 | 2015 | 2016 | 2017 | 2018 |
|------|------|------|------|------|------|
| 36 Fellows | 48 Fellows | 42 Fellows | 42 Fellows | 15 Fellows | 38 Fellows |
| 6 Mentors | 8 Mentors | 5 Mentors 3 PMs | 6 Mentors 3 PMs | 2 Mentors 2 PMs | 5 Mentors 4 PMs |
| 12 Projects | 14 Projects | 12 Projects | 12 Projects | 6 Projects | 10 Projects |
| 12 Weeks | 12 Weeks | 14 Weeks | 13 Weeks | 12 Weeks | 12 Weeks |
| THE UNIVERSITY OF CHICAGO | THE UNIVERSITY OF CHICAGO | THE UNIVERSITY OF CHICAGO | THE UNIVERSITY OF CHICAGO | NOVA NOVA SCHOOL OF BUSINESS & ECONOMICS | NOVA NOVA SCHOOL OF BUSINESS & ECONOMICS THE UNIVERSITY OF CHICAGO |

# 4000+ Applicants over the past 6 Summers



## 4000+ Applicants

## Over 70 countries and 400 universities

Computer Science, Statistics, Math, Applied Math, Sociology, Economics, Public Policy, Political Science, Physics, Chemistry, Biology, Public Health, Psychology, Engineering, BioStatistics, Geography, Business, Neuroscience, …

**224**

**Fellows**

70 Projects

Data Science For Social Good

DATA SCIENCE FOR SOCIAL GOOD
ATLANTA 2016

Data Science For Social Good
Summer Fellowship
THE UNIVERSITY OF CHICAGO

KDD 2014
This year's special theme:
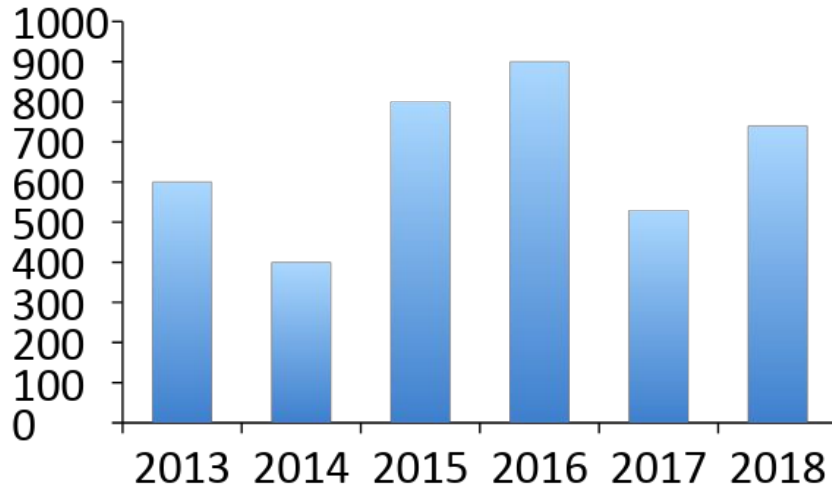Data Science for Social Good

UNIVERSITY of WASHINGTON

eScience Institute
DATA SCIENCE FOR SOCIAL GOOD

IBM Social Good Fellowship

SoGood 2016

BAYES IMPACT

Data For Good Exchange 2016

Volunteerism for the Data Generation:
Using Data Science Superpowers for Social Good
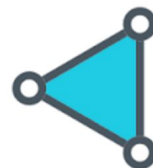
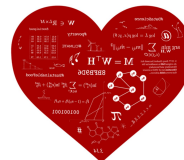LELAND STANFORD JUNIOR UNIVERSITY
DIE LUFT DER FREIHEIT WEHT
1891

Statistics for Social Good

We're a group of Stanford students, researchers, and faculty exploring the potential to promote social good through effective data analysis.

#Data4Good

NC Data4Good
"data crunch for social good"

Partnership for Social Good

Data Science Initiative
UNC CHARLOTTE

O'REILLY
Data and Social Good
Using Data Science to Improve Lives, Fight Injustice, and Support Democracy
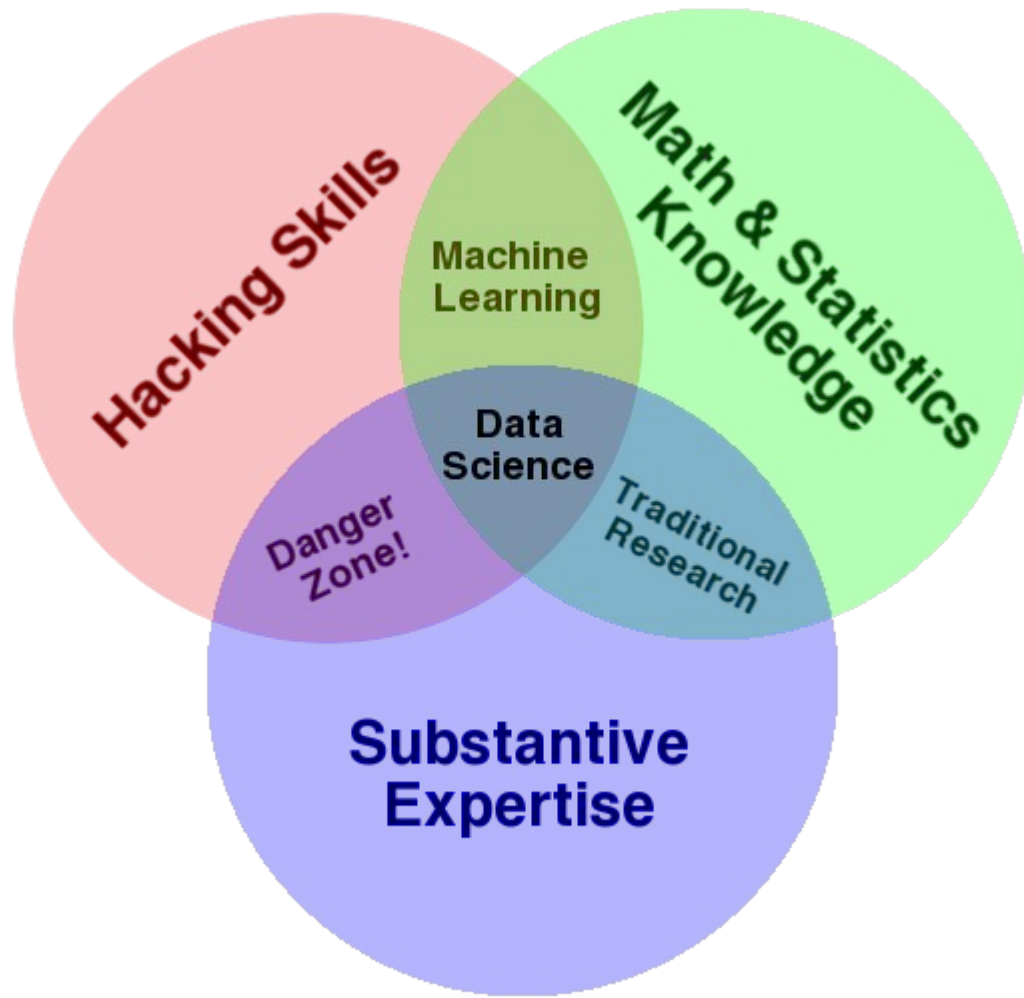Mike Barlow

DO GOOD DATA 2015

DataKind
USING DATA IN THE SERVICE OF HUMANITY

DRIVENDATA

Data Science For Social Good Europe
Summer Fellowship 2018
CASCAIS          NOVA NOVA SCHOOL OF BUSINESS & ECONOMICS          THE UNIVERSITY OF CHICAGO

# What we got wrong

*Drew Conway*

*Joel Grus*

Hacking Skills

Machine Learning

Math & Stats Knowledge

Data Science

Not *that* dangerous, in retrospect

DSSG
~~NSA~~

~~Outside~~ Committee Member

You at Your PhD Hooding Ceremony

Grad School Officemate's Thesis Advisor

CfA Brigade Captain

Substantive Expertise

The Dalai Lama
~~James Bond~~ Villain

Good
~~Evil~~
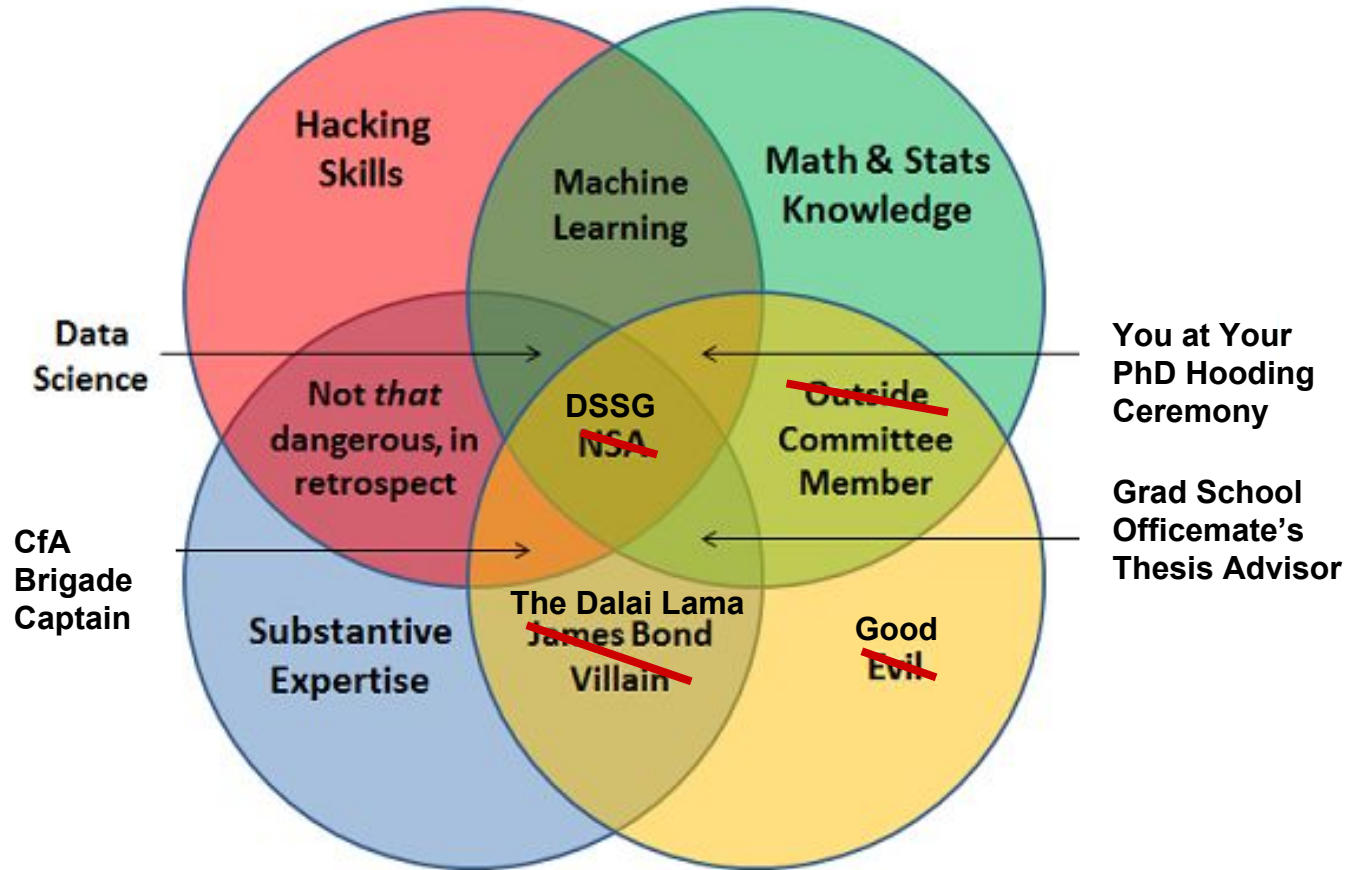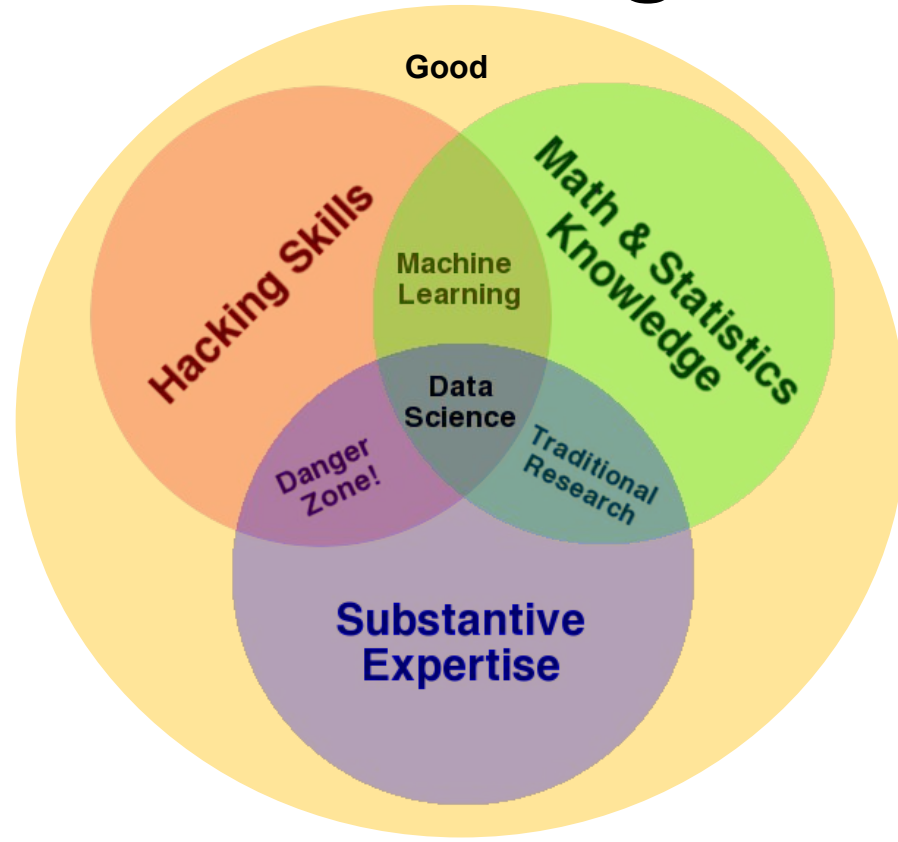
# All data science must be grounded in a sense of the good

What will it take for socially beneficial data science to go from
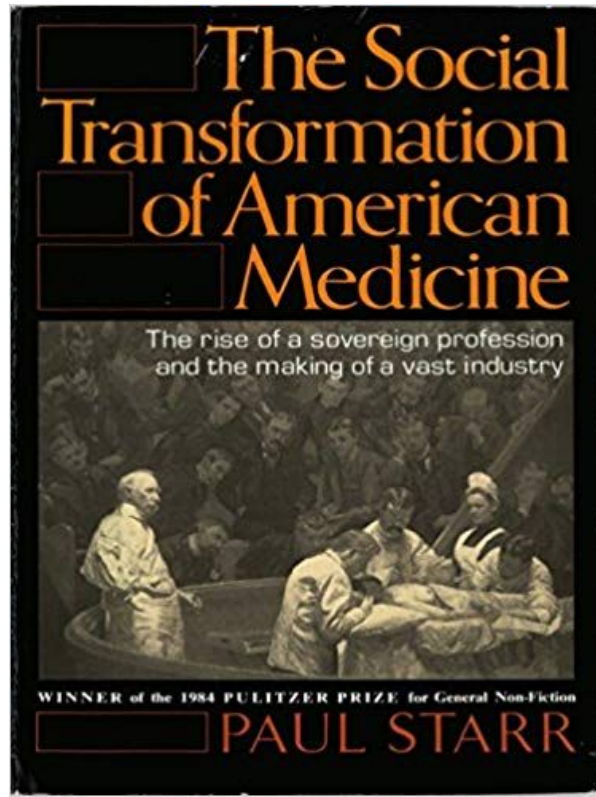**doing data science for good**
to
**doing good data science**?

# Learning from our early experience as well as that of older professions

**Data Science For Social Good**
Summer Fellowship

THE UNIVERSITY OF CHICAGO

+

The Social Transformation of American Medicine

The rise of a sovereign profession and the making of a vast industry

WINNER of the 1984 PULITZER PRIZE for General Non-Fiction
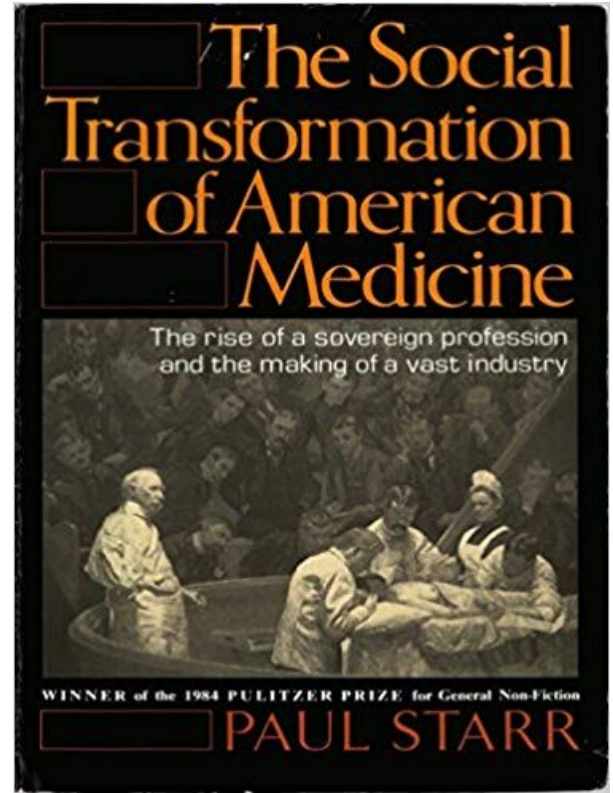
PAUL STARR

# **Three reasons why**
all data science should be taught as data science for social good

# Reason 1

# As a check on the increasing power of the data scientist
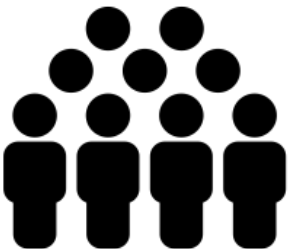
"The dream of reason never took power into account."

# We teach data science at a distance



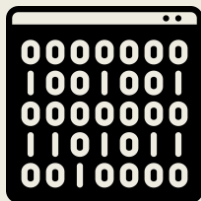**Social Context** of Individual Actions
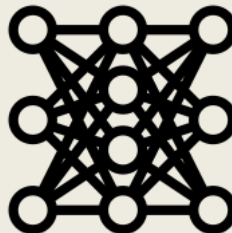
Digitally Instrumented **Individual Actions**

**The Sandbox of Data Science Education**

**Assemble Data** from Digital Traces

Build & Test **Models**

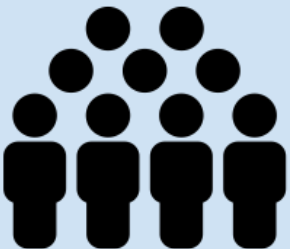Design **Outputs**

Influence **Individual Action**
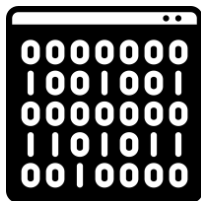
**Social Context** of Influenced Action

# Data science happens in a social context

# If you are a data scientist, your work will have social consequences

# Reason 2

# To establish professional norms

"A profession, sociologists have suggested, is an occupation that regulates itself through systematic, required training and collegial discipline;
that has a base in technical, specialized knowledge;
and that **has a service rather than profit orientation, enshrined in its code of ethics.**"

The Social Transformation of American Medicine

The rise of a sovereign profession and the making of a vast industry

WINNER of the 1984 PULITZER PRIZE for General Non-Fiction

PAUL STARR

# FORTS Framework

FAIRNESS - I make a dedicated effort to understand, mitigate and communicate the presence of bias in both data practice and consumption.

OPENNESS - I practice humility and openness. Transparent practices, community engagement, and responsible communications are an integral part of my data ethics practice.
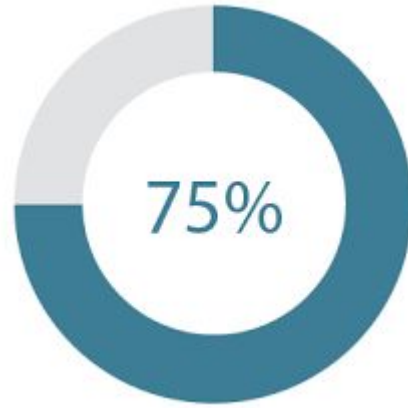
RELIABILITY - I ensure that every effort is made to glean a complete understanding of what is contained within data, where it came from, and how it was created. I also extend this effort for future users of all data and derivative data.

TRUST - I work to build public confidence in data practitioners. I make every effort to use data and algorithms in ways that maximize the informed participation of people around the world.

SOCIAL BENEFIT – I place people before data and am responsible for maximizing social benefit and minimizing harm. I consider the impact of my work on communities of people, other living beings, ecosystems and the world-at-large.

# Reason 3

To attract and keep the best and brightest minds to our profession

75%

% of adults under the age of 35 willing to take a pay cut to work at a job that has social purpose.

"You do things that you get a kick out of, and luckily, sometimes they turn out to be important."

Rainer Weiss
Nobel Laureate in Physics

*Nobel Minds, 2017*

It's time to update what it means
to do
good data science

**Three suggestions for how**
all data science can be taught as
data science for social good

# Suggestion 1

# Add **clinical practice requirements** to all data science programs

# From classroom to clinic
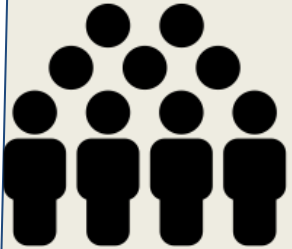# Putting data scientists into the social context of their algorithms



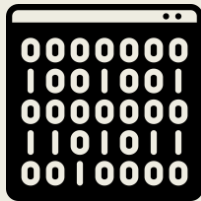**Clinical Practice in Data Science Education**
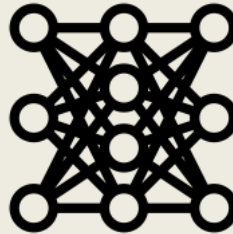
**Social Context** of Individual Actions

Digitally Instrumented **Individual Actions**

**Assemble Data** from Digital Traces

Build & Test **Models**

Design **Outputs**

Influence **Individual Action**

**Social Context** of Influenced Action

Chicago protests over police shooting of Laquan McDonald

STOP POLICE CRIMES

BREAKING NEWS

POLICE RELEASE VIDEO OF OFFICER SHOOTING TEEN

A Disturbing Number Of Chicago Cops Have 30 Or More Complaints, Data Shows

BY KATE SHEPHERD IN NEWS ON NOV 11, 2015 12:51 PM
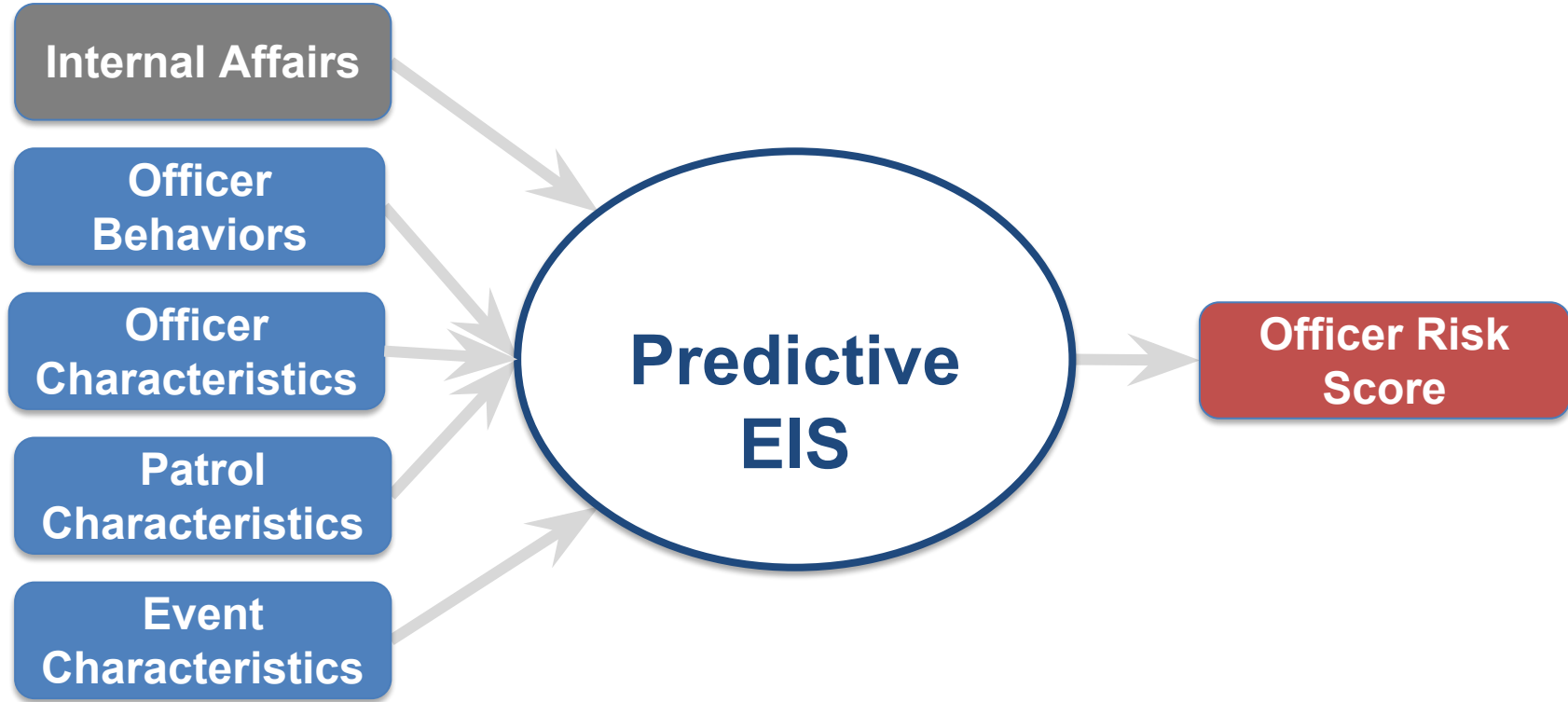
APR 10, 5:27 PM EDT

International Business Times

Police Misconduct Cases Have Cost Chicago $662 Million Since 2004: Report

BY ADAM LIDGETT ON 03/19/16 AT 4:56 PM
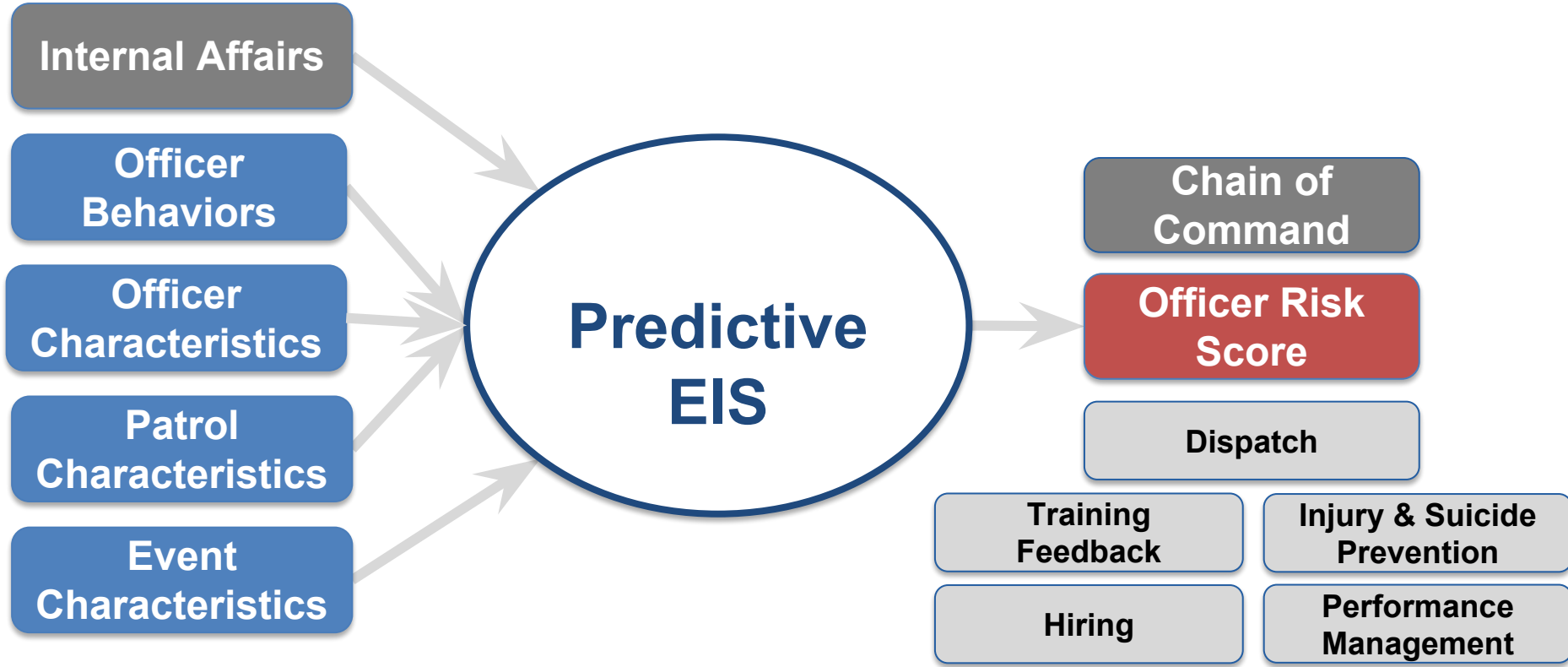
# Predicting officer's risk of adverse event

# From classroom to clinic
## Putting data scientists into the social context of their algorithms

# Toward a more holistic solution

Internal Affairs

Officer Behaviors

Officer Characteristics

Patrol Characteristics

Event Characteristics

→ Predictive EIS →

Chain of Command

Officer Risk Score

Dispatch

Training Feedback

Injury & Suicide Prevention

Hiring

Performance Management

This will be difficult for many data science programs to do.

We're here to help.

# Suggestion 2

Add written or verbal **discussion of the social and ethical implications** of models into every problem set in data science coursework

# Questions we've asked

- What biases may exist in the data you've been given? How can you find out?
- How will your choices with tuning parameters affect different populations represented in the data?
- How do you know you aren't getting the right answer to the wrong question?
- How would you justify what you'd built to someone whose welfare is made worse off by the implementation of your algorithm?
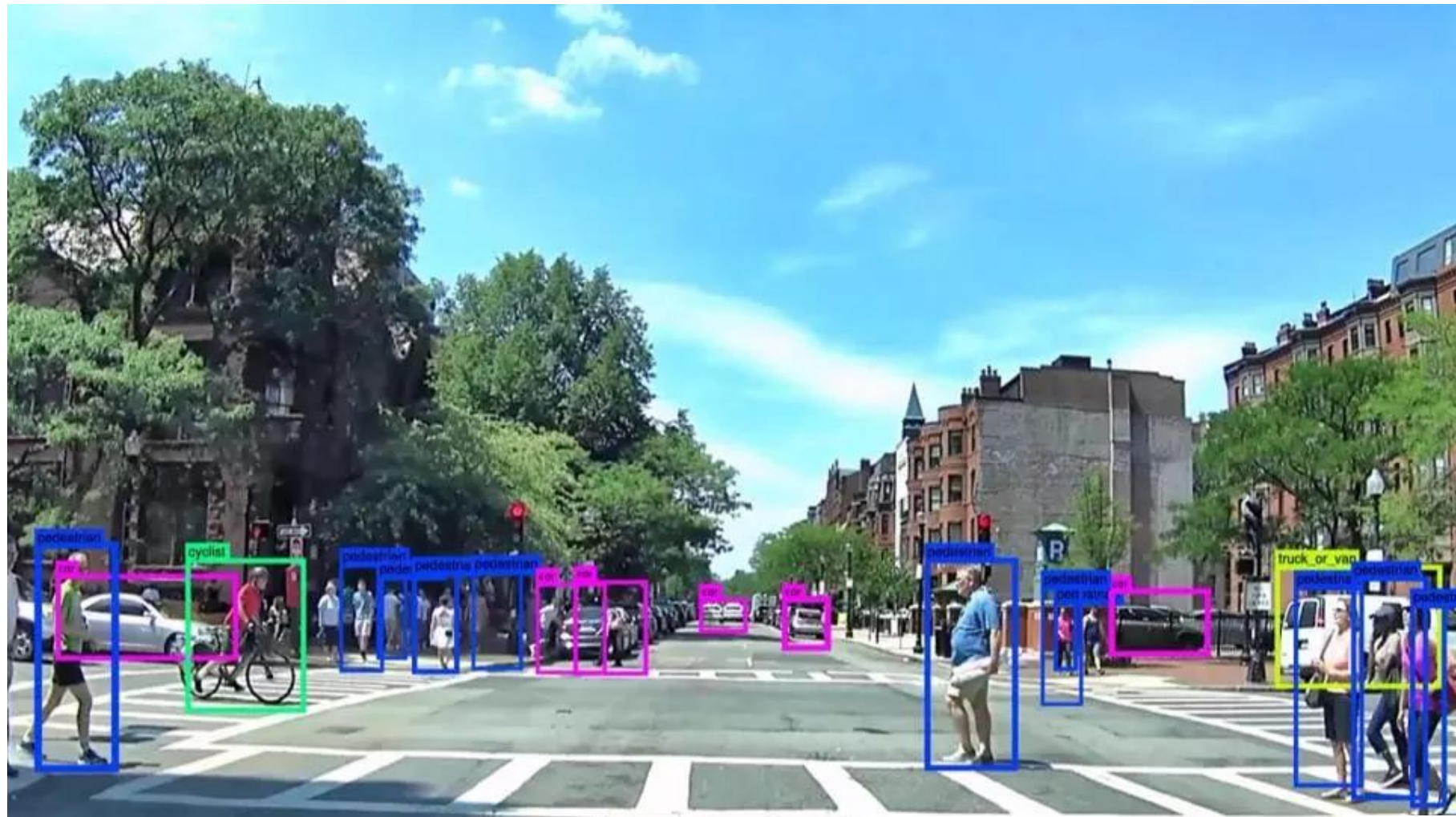
# Suggestion 3

**Provide guidance to employers** for incorporating ethics case studies into hiring, apprenticeships, and mentorship opportunities.

What could happen if we teach data scientists to always be doing data science for social good

NUERALA, 2018

Thank You!

@matthewgee
mattgee@uchicago.edu