

Social Data Science with Text

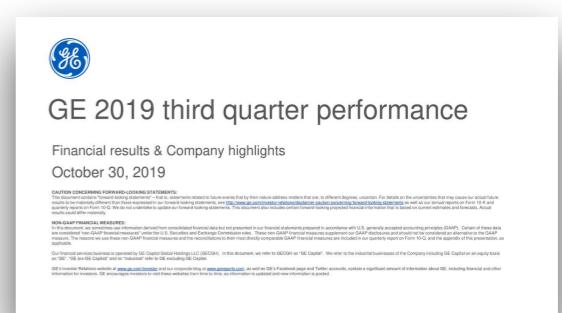
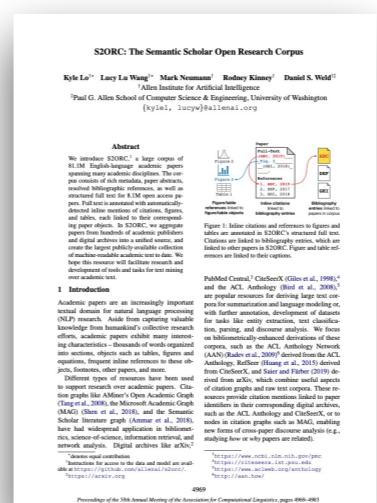
Katherine Keith

PhD Candidate

College of Information and Computer Sciences
University of Massachusetts Amherst



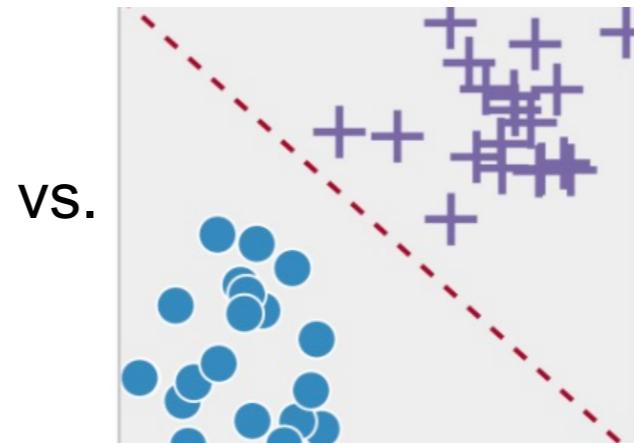
Social Data Science with Text





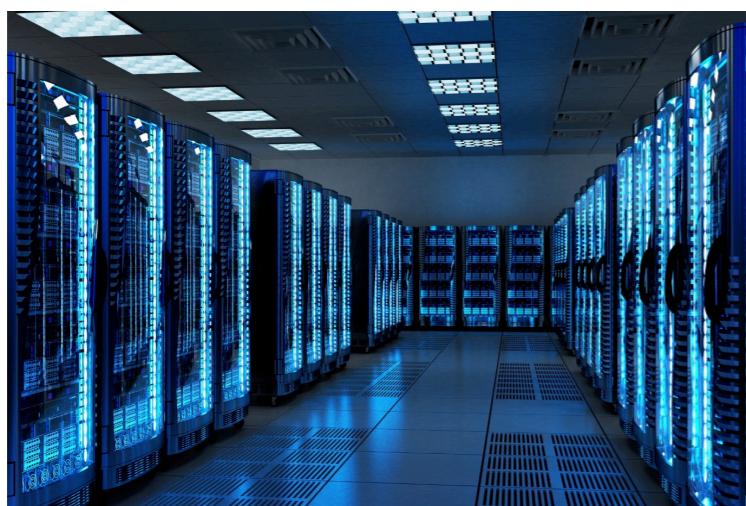
Social Data Science with Text





Machine Learning
Natural Language
Processing

Social Data Science with Text

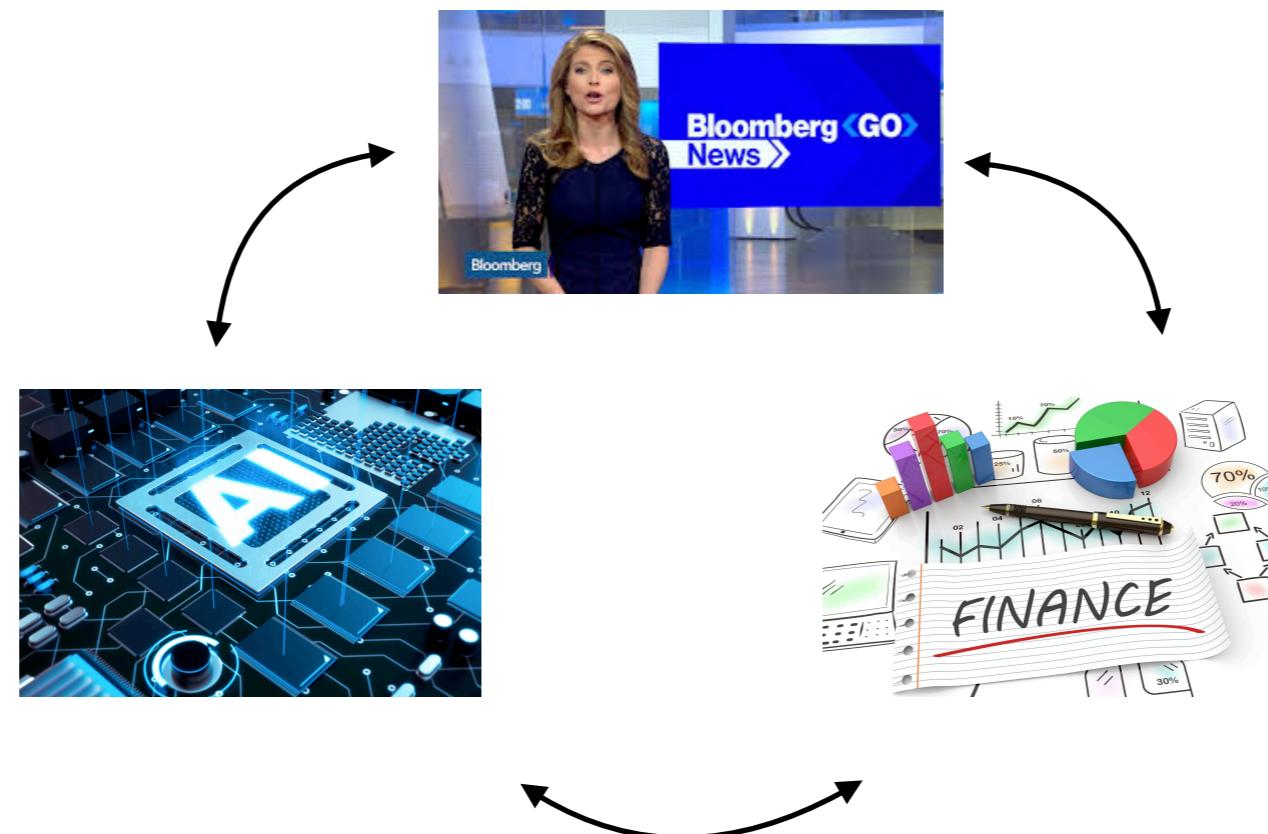


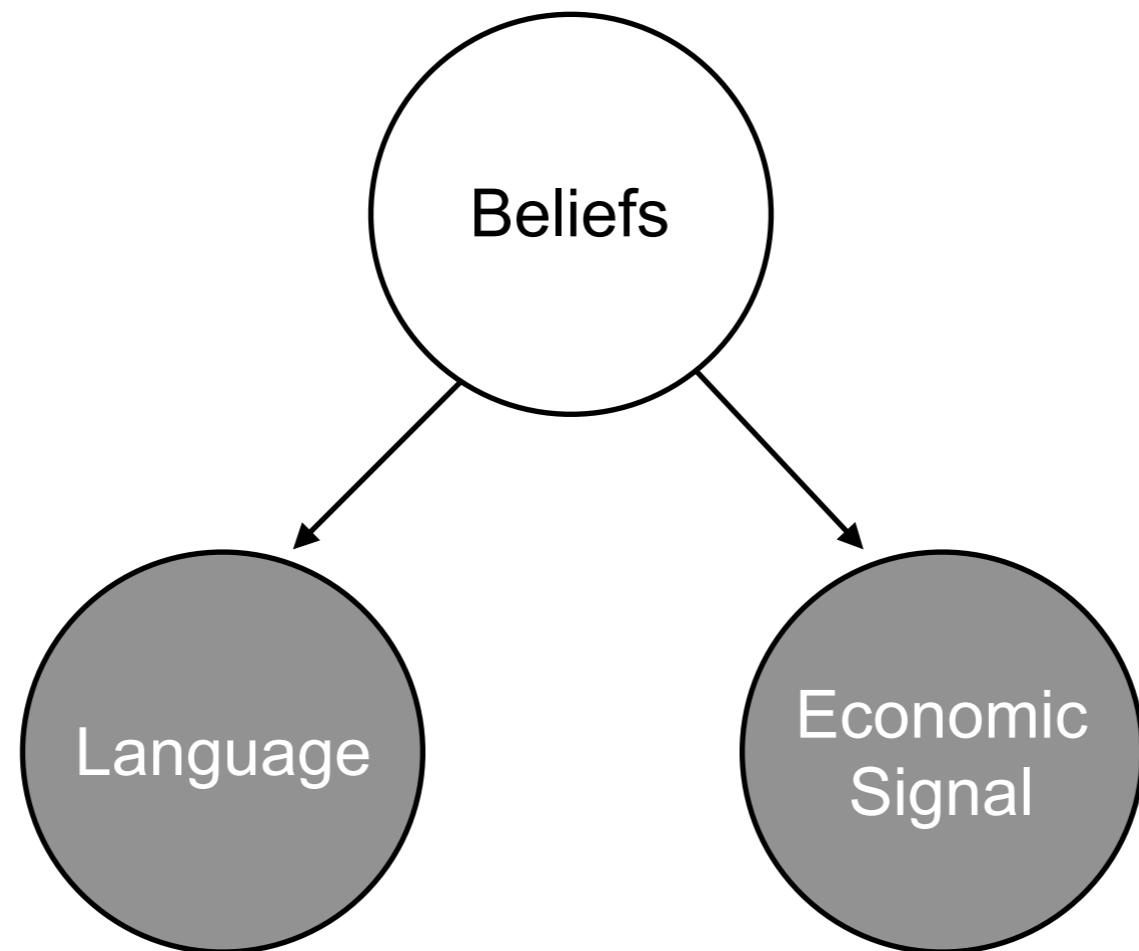
Interdisciplinary collaboration



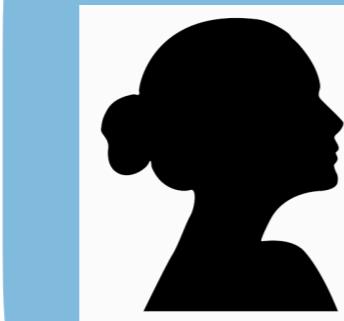
Bloomberg Data Science Fellow 2019-2021
Bloomberg Research Intern 2018, 2020

Interdisciplinary collaboration





Analyst's beliefs about a firm



Beliefs

Earnings call transcripts



GE 2019 third quarter performance

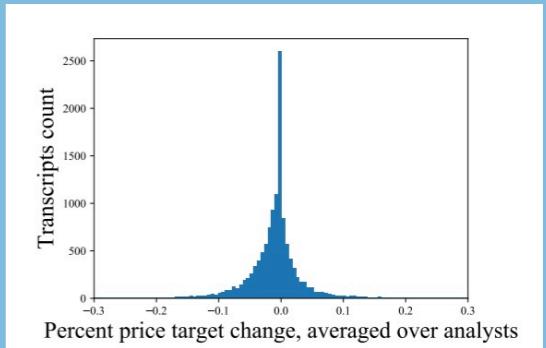
Financial results & Company highlights
October 30, 2019

CAUTION CONCERNING FORWARD-LOOKING STATEMENTS:
This document contains forward-looking statements – that is, statements related to future
results to be materially different than those expressed or implied by forward-looking statements, or
statements concerning our intentions, plans or beliefs. We do not undertake to update our forward-looking statement
results could differ materially.
NON-GAAP FINANCIAL MEASURE:
In this document, non-GAAP financial information derived from consolidated financial data is
considered "non-GAAP financial measures" under the U.S. Securities and Exchange Commission
regulations. Non-GAAP financial measures are not prepared in accordance with generally accepted
accounting principles. Our financial services business is operated by GE Capital Global Holdings, LLC (GCGH).
GE Capital is a registered service mark of General Electric Company, Inc. GE Capital is also a registered trademark of GE Capital.
GE Capital is a registered service mark of General Electric Company, Inc. GE Capital is also a registered trademark of GE Capital.
GE's Investor Relations website at www.ge.com/investor and our corporate blog at www.ge.
GE encourages investors to visit these websites from time to time for information for investors.



Language

Analysts' price target before and after call



(**Keith** and Stent, “Modeling financial analysts’ decision making via the pragmatics and semantics of earnings calls.” ACL, 2019)

Belief that policy is
driving economic uncertainty

Beliefs

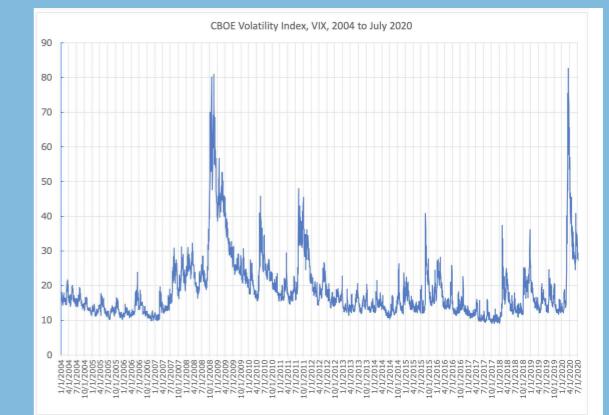
News reports



Language

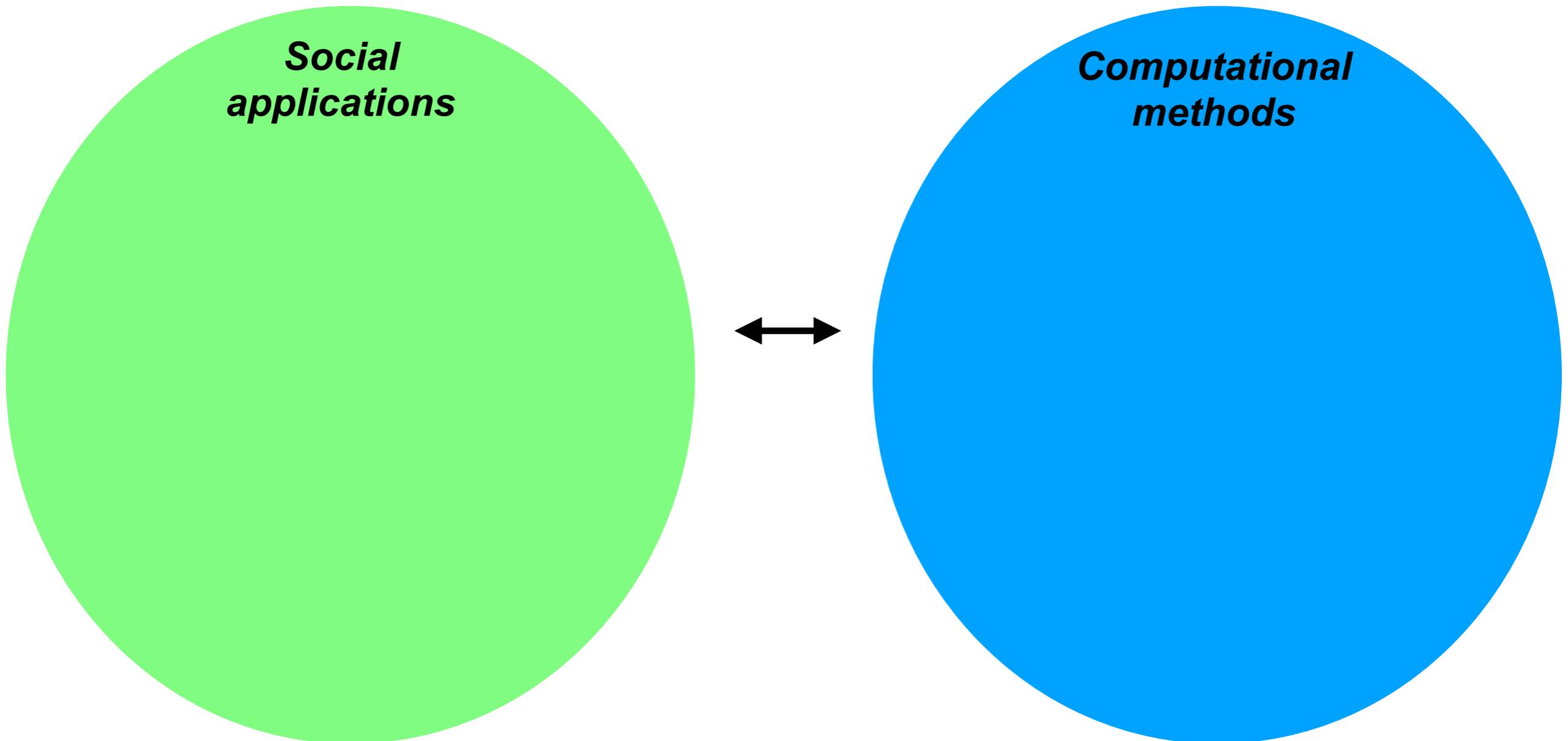
Economic
Signal

Stock volatility index

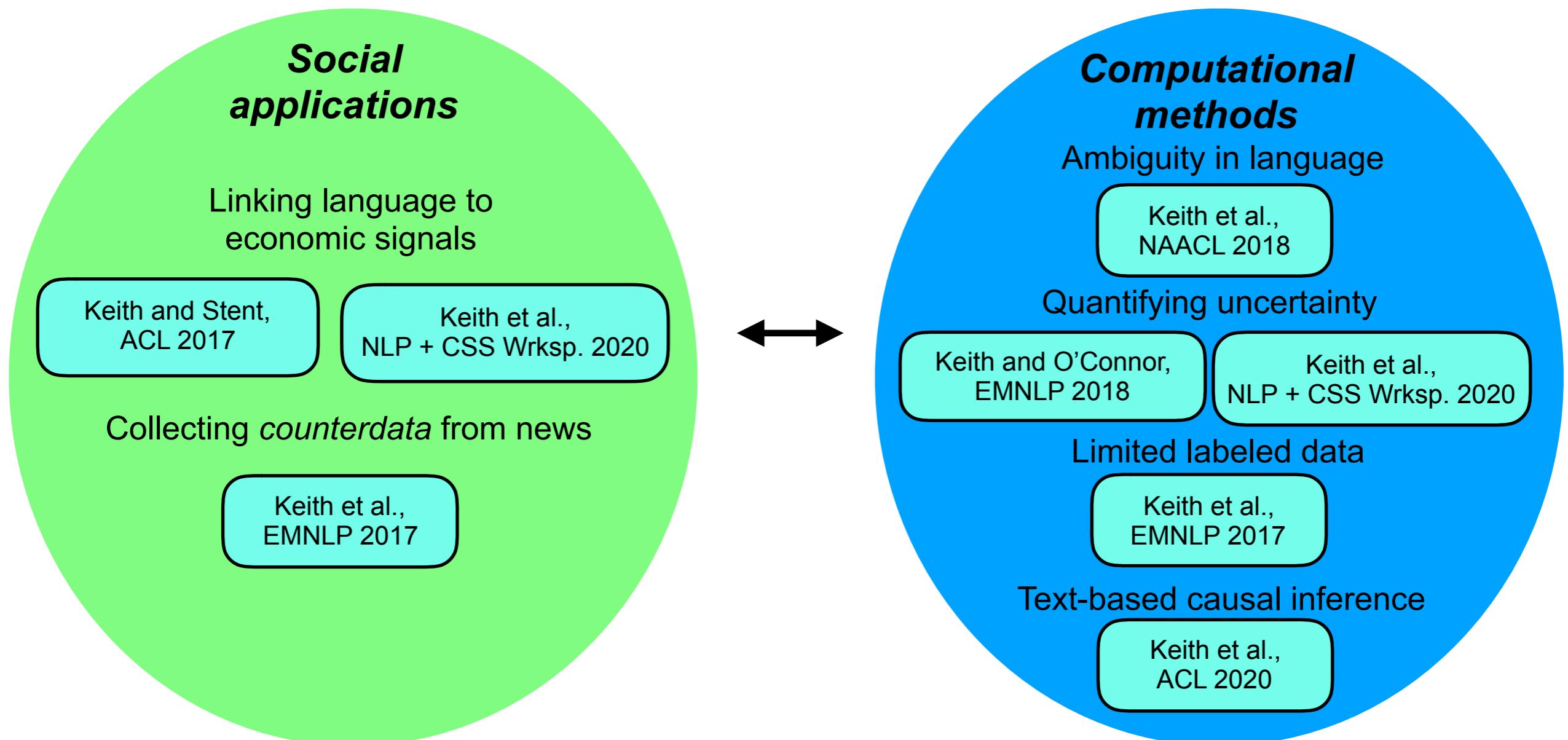


(**Keith et al.**, “Uncertainty over Uncertainty: Investigating the Assumptions, Annotations, and Text Measurements of Economic Policy Uncertainty.” NLP+CSS Workshop, 2020)

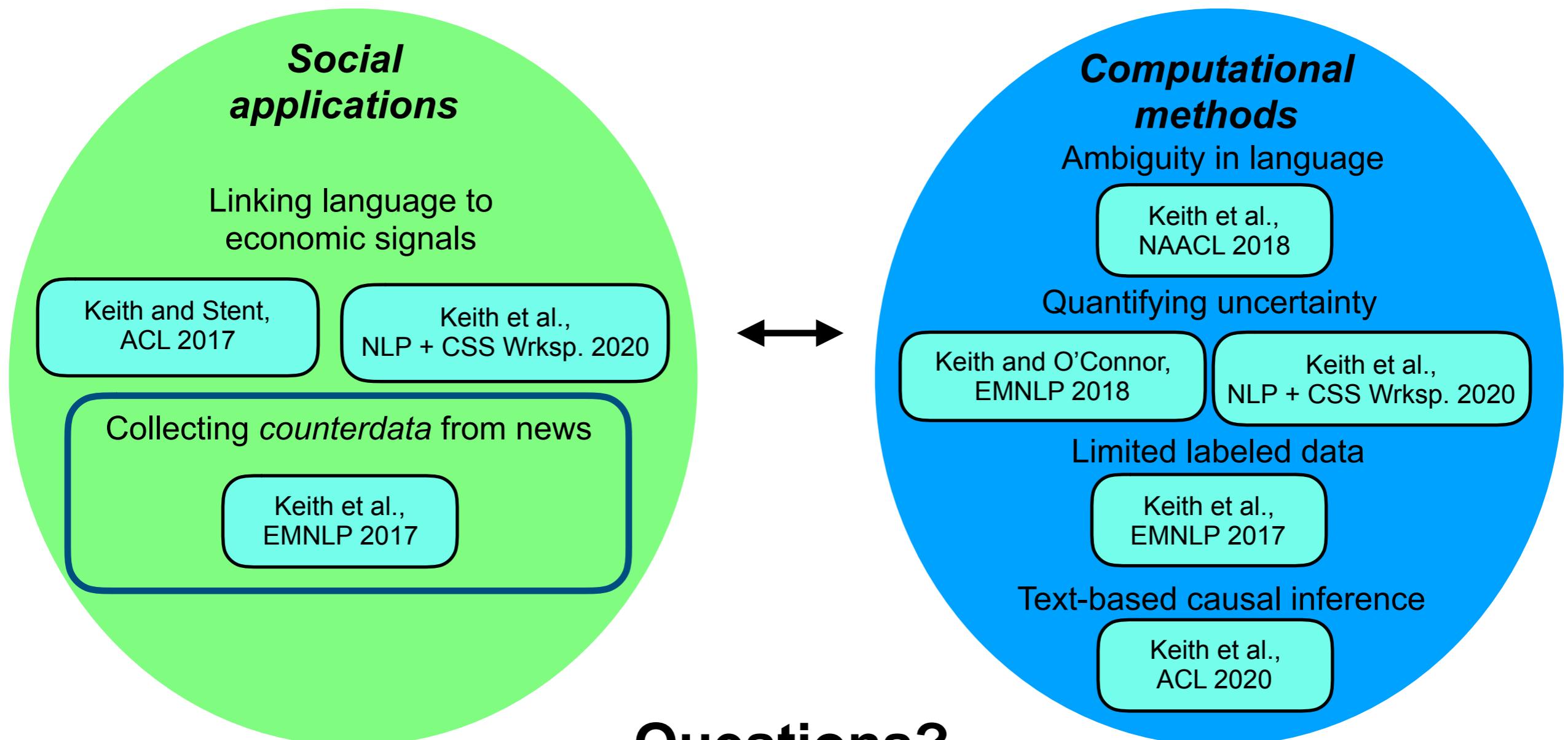
My research philosophy: *Social data science with text* requires a rich symbiosis between *domain applications* and *computational methods*.



My Research Agenda



My Research Agenda



How can **data science** contribute to **social impact**?

How can **data science** contribute to **social impact?**

How can we improve outcomes in the world?

How can **data science** contribute to **social impact?**

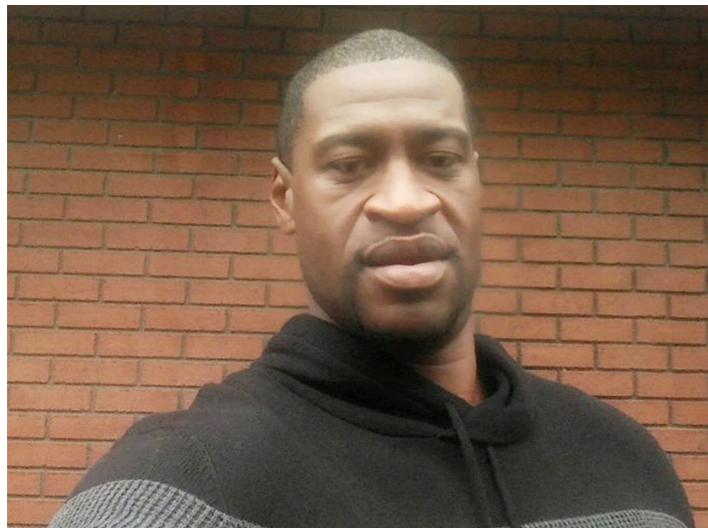
How can we *improve outcomes* in the world?

values +
measurement

How can we *improve outcomes* in the world?

value:
police should
not kill civilians

How can we *improve outcomes* in the world?

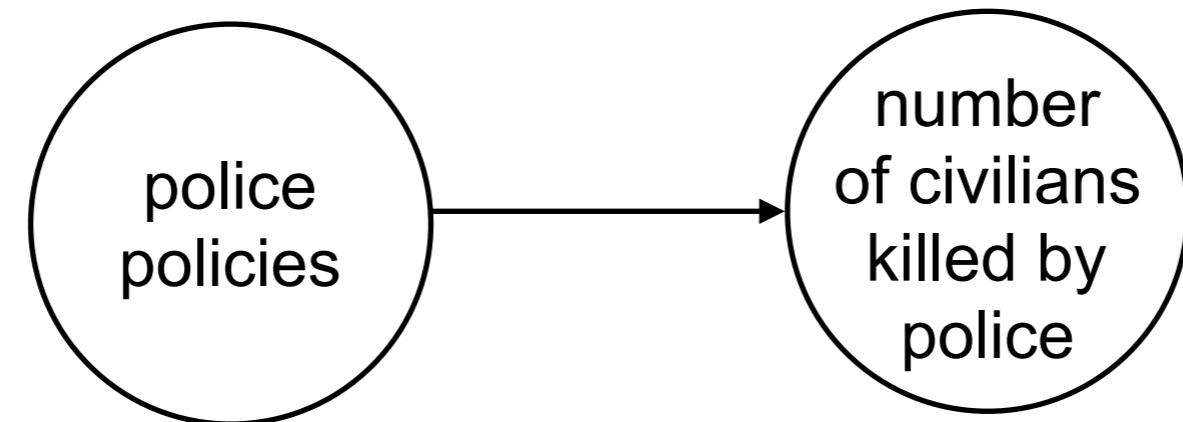


Source: Times Magazine

value:
police should
not kill civilians

How can we *improve outcomes* in the world?

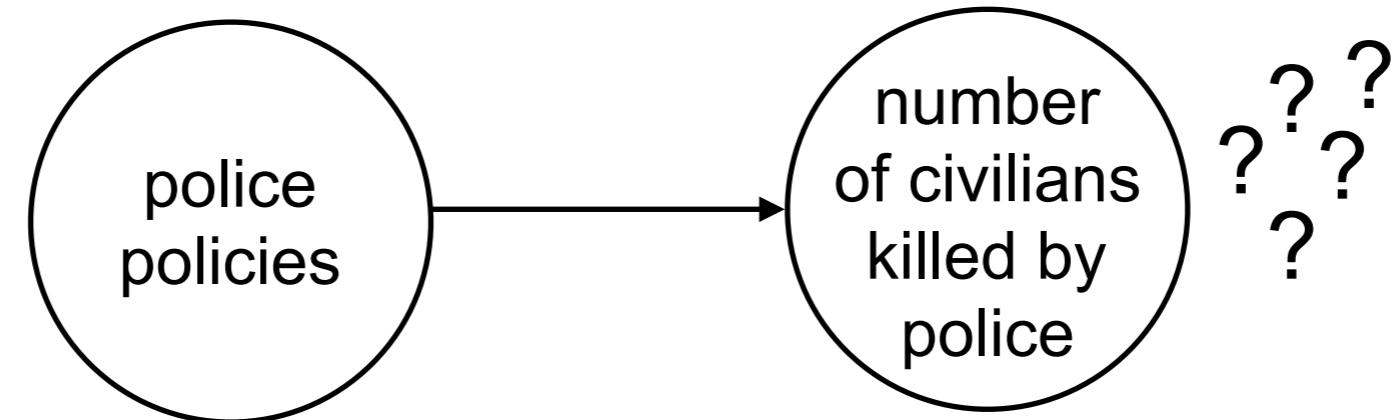
measurement:



value:
police should
not kill civilians

How can we *improve outcomes* in the world?

measurement:



U.S. federal government systematically undercounts or fails to count police fatalities

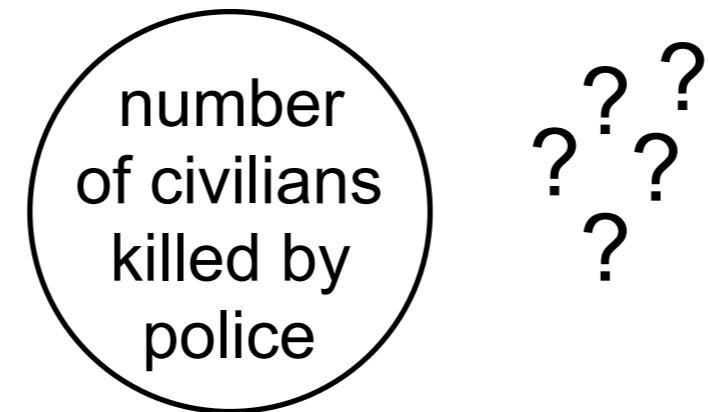
- **2013:** Obama signs *Death in Custody Reporting Act (DCRA)*
 - Requires police departments to report every time a citizen dies in custody
- **2019:** FBI begins *National Use of Force Data Collection*
 - Local law enforcement agencies are not required to participate and the data is not yet public

Counterdata: grassroots collection of missing datasets

(D'Ignazio and Klein, Data Feminism, 2020)

Counterdata: police fatalities from news reports

measurement:



Counterdata: police fatalities from news reports

measurement:

number
of civilians
killed by
police



news
reports

Approach 1: Manual

FatalEncounters.org,
KilledByPolice.net,
The Guardian,
Washington Post



Approach 2: Automated



**(Keith et al.,
EMNLP 2017)**

Issue: Cost of
human time and
emotional strain

Why is automatically detecting police fatality events hard?

Police killed PERSON.

Police killed PERSON.

Police officers spotted the butt of a handgun in Alton Sterling's front pocket and saw him reach for the weapon before opening fire, according to a Baton Rouge Police Department search warrant filed Monday that offers the first police account of the events leading up to his fatal shooting.

Police killed PERSON.

long-range dependencies

Police officers spotted the butt of a handgun in Alton Sterling's front pocket and saw him reach for the weapon before opening fire, according to a Baton Rouge Police Department search warrant filed Monday that offers the first police account of the events leading up to his fatal shooting.

Police killed PERSON.

long-range dependencies

Police officers spotted the butt of a handgun in Alton Sterling's front pocket and saw him reach for the weapon before opening fire, according to a Baton Rouge Police Department search warrant filed Monday that offers the first police account of the events leading up to his fatal shooting.

coreference

Police killed PERSON.

long-range dependencies

Police officers spotted the butt of a handgun in Alton Sterling's front pocket and saw him reach for the weapon before opening fire, according to a Baton Rouge Police Department search warrant filed Monday that offers the first police account of the events leading up to his fatal shooting.

coreference

*event
coreference*

Automatically detecting police fatality events

Domain knowledge vs. Machine Learning

Automatically detecting police fatality events

Domain knowledge vs. Machine Learning

Domain Knowledge: Keyword Matching

Input: sentences

*PERSON was **fatally shot** by **police**.*

Keyword matching

officers, police, cops,
troopers, deputy, ...

*Officers reported PERSON
was **killed** in a car accident.*

kill, killing,
shoot, shooting,
murder, homicide ...

Classification

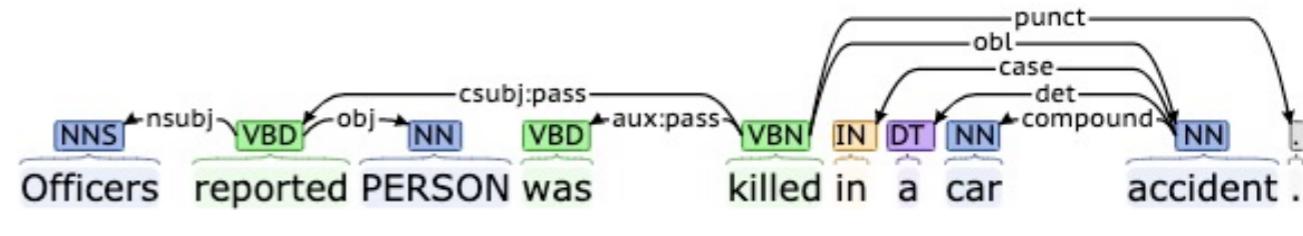
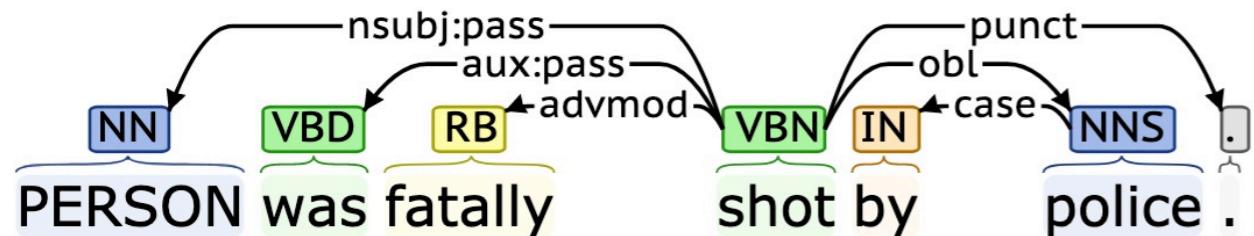
Yes

Yes

Issue: many
false positives (*low
precision*)

Domain Knowledge: Syntactic Dependency Parsing

Input: automatically infer dependency parse trees over sentences



Rules over dependency paths

PERSON <-nsubj:pass <-

kill, killing,
shoot, shooting,
murder, homicide ...

->obl ->

officers, police,
cops, troopers,
deputy, ...

Classification

Yes

No

Issue:
Difficult for a
domain expert to list
all possible rules
(*low recall*)

(e.g. Chen and Manning, EMNLP, 2014; Nivre et al. LREC, 2016; Keith et. al, NAACL, 2018)

Automatically detecting police fatality events

Domain knowledge vs. Machine Learning

Supervised Machine Learning

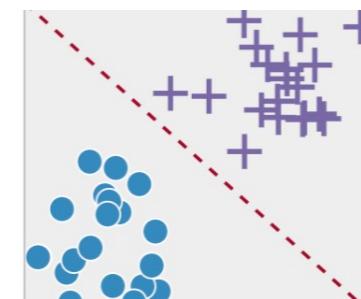
1. Gather input data

Police killed PERSON.

2. Label input data

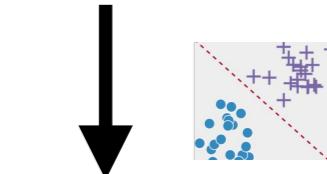
- *logistic regression* with bag of words features
- *convolutional neural networks* initialized with pre-trained word embeddings

3. Train model: statistical pattern matching between inputs and labels



4. Inference: (generalization) apply trained model on unseen inputs

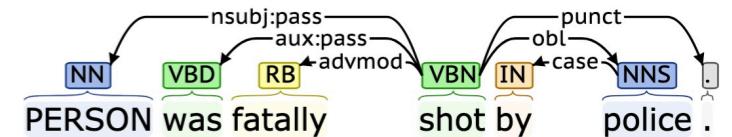
PERSON died in a police homicide.



Yes

kill, killing,
shoot, shooting,
murder, homicide ...

Need to **evaluate tradeoffs** for
methods of event extraction



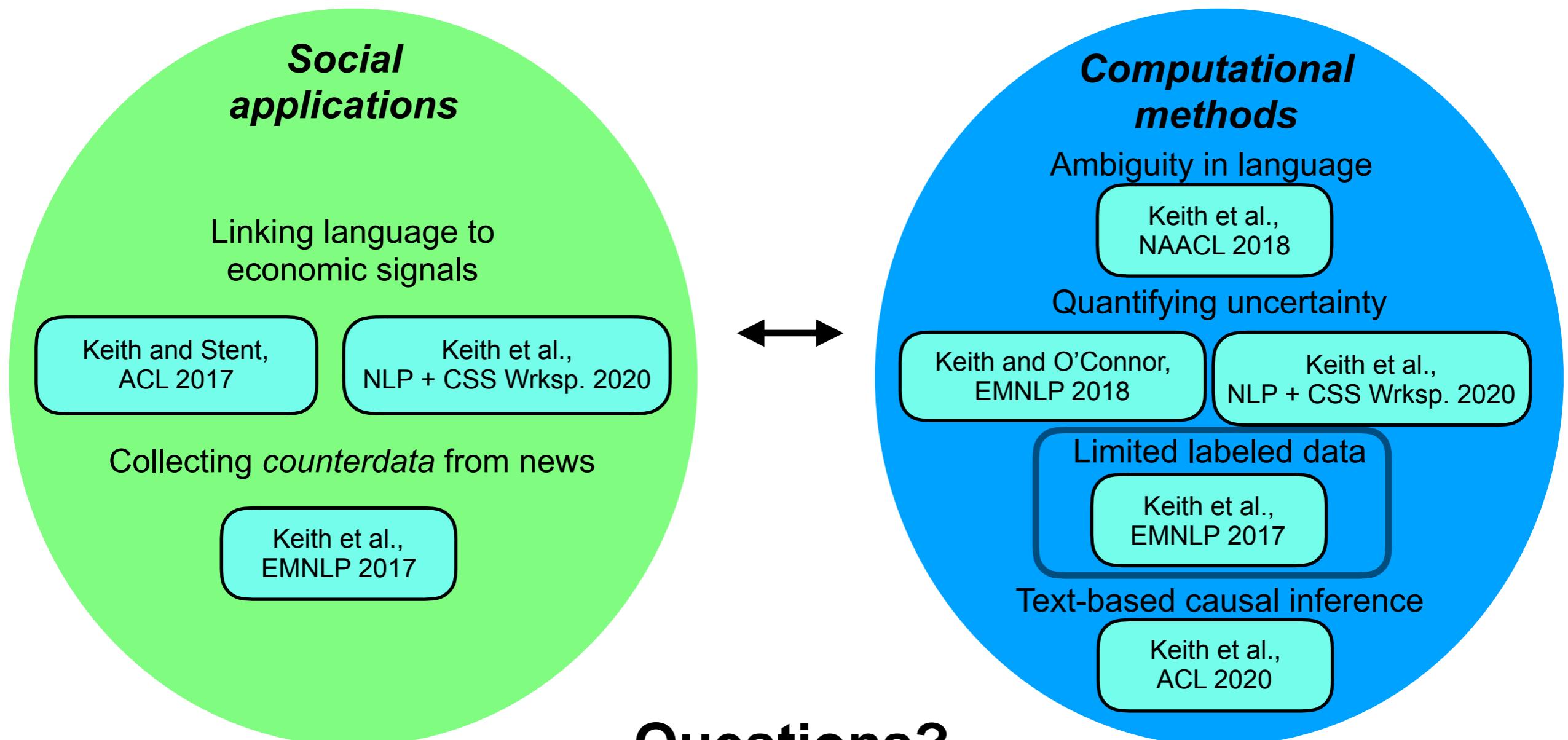
Event extraction methods can
be used to collect **counterdata**

number
of civilians
killed by
police

Without these **measurements**,
some questions can be nearly
impossible to answer.

How can we **improve**
outcomes in the world?

My Research Agenda



Supervised Machine Learning

1. Gather input data

Police killed PERSON.

2. Label input data

Yes/No

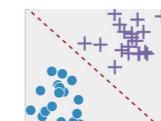
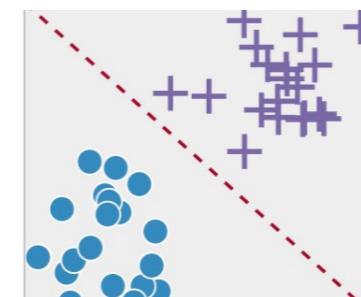
3. Train model
matching between inputs and labels

How do we design and train models with few to no labeled examples?

4. Inference: (generalization) apply trained model on unseen inputs

PERSON died in a police homicide.

Yes



Distant supervision

(Craven and Kumlien, 1999; Mintz et al., 2009)



(e.g. Fatal Encounters)

Requirement: knowledge in an external database that is not currently aligned with text

Distant supervision

(Craven and Kumlien, 1999; Mintz et al., 2009)

1. Impute positive labels from an **external database**



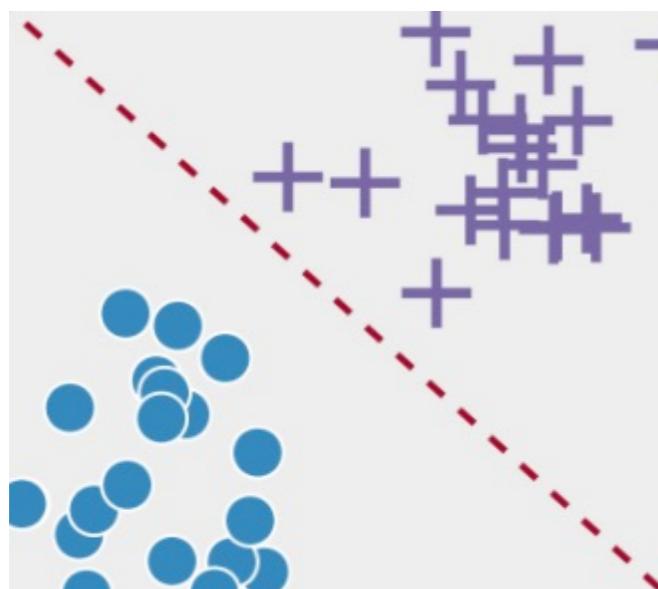
(e.g. Fatal Encounters)



John Doe was killed by police.

Police fatally shot **Jane Public**.

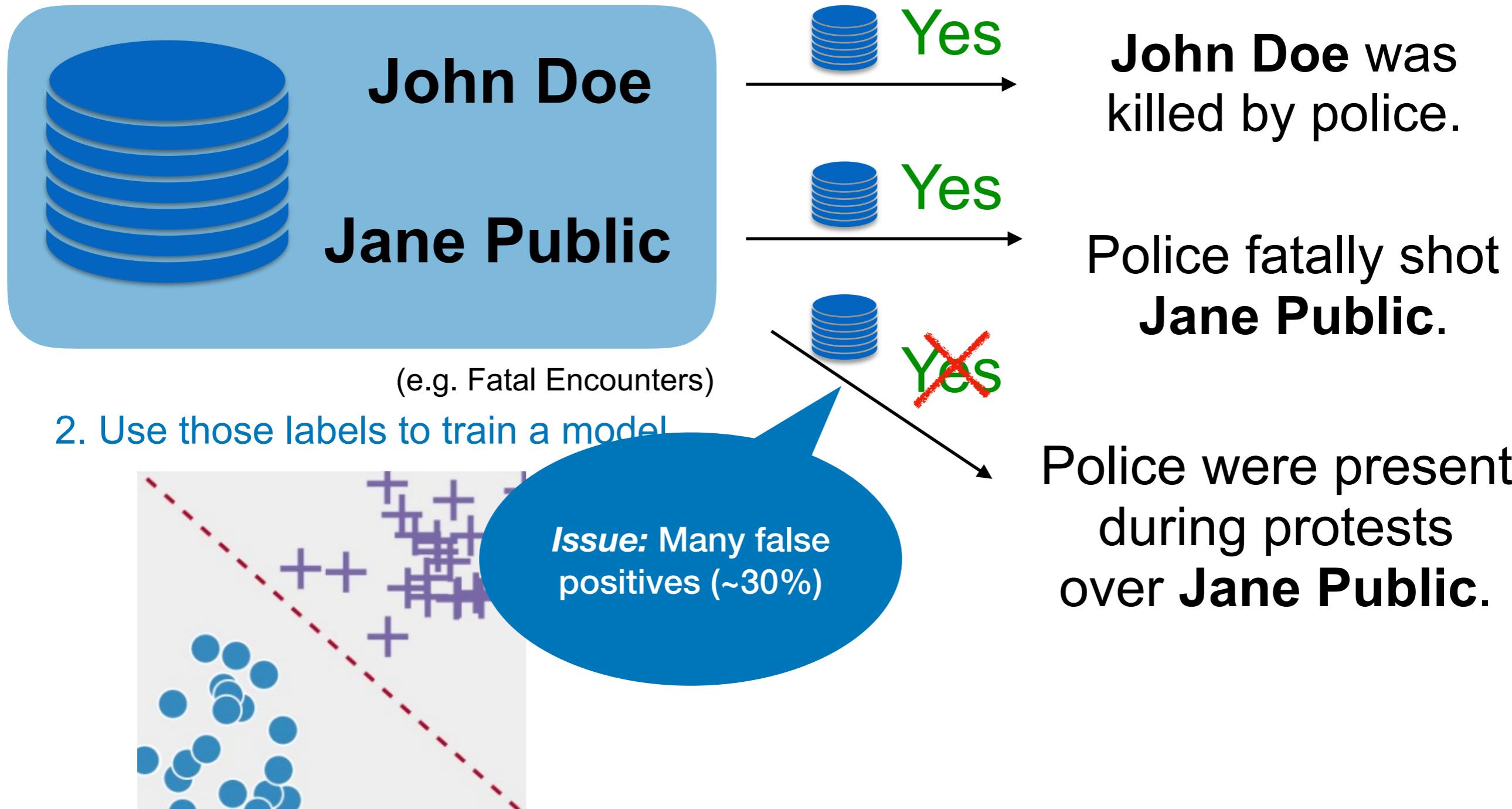
2. Use those labels to train a model



Distant supervision

(Craven and Kumlien, 1999; Mintz et al., 2009)

1. Impute positive labels from an **external database**



Latent disjunction model

(Keith et al., EMNLP 2017)

For each sentence i

$$e_i$$

Entity (person's name)

$$y_{e_i} \in \{0, 1\}$$

Entity label

$$x_{\mathcal{M}(e_i)}$$

Set of all sentences
with entity

$$z_i \in \{0, 1\}$$

Label of whether the
sentence indicates a
police fatality event

Latent disjunction model

(Keith et al., EMNLP 2017)

For each sentence i

e_i

Entity (person's name)

$y_{e_i} \in \{0, 1\}$

Entity label

$x_{\mathcal{M}(e_i)}$

Set of sentences

$z_i \in \{0, 1\}$

Label for sentence in police fatality event

Expectation

Algorithm

Builds in assumption that *at least one sentence* must be a positive per entity.

E-Step:

$$q(z_i) := P(z_i | x_{\mathcal{M}(e_i)}, y_{e_i})$$

M-step:

$$\max_{\theta} \sum_i \sum_{z \in \{0,1\}} q(z_i = z) \log P_{\theta}(z_i = z | x_i)$$

Parameters for sentence classifier

Iterate until convergence

Empirical evaluation

Police fatality data

Google News

	Train	Test
Document dates	Jan-Aug 2016	Sept-Dec 2016
Total Docs.	793,010	317,345
Total Entities	49,203	24,550

Data publicly available: <http://slanglab.cs.umass.edu/PoliceKillingsExtraction/>

Distant supervision vs. Latent disjunction model

Cheaper than standard supervision

Empirical results: Improves entity-level F1 by 18% on the test set

Reduces distantly-labeled *false positives*

Error Analysis

Logan Clarke was **shot by a campus police officer** after waving kitchen knives at fellow students outside the cafeteria at Hug High School in Reno, Nevada, on December 7.

Model prediction: Yes

True value: No

Error Analysis

Logan Clarke was **shot** by a **campus police officer** after waving kitchen knives at fellow students outside the cafeteria at Hug High School in Reno, Nevada, on December 7.

Model prediction: Yes
True value: No

Model has not learned to distinguish *killed* vs. *shot*

Only **official police** are in the database so campus police or security guards count as errors

User Interface Prototype

The screenshot shows a web-based application for filtering police fatalities. At the top, there's a header bar with navigation icons (back, forward, search, etc.) and a title "Police Fatalities". Below the header is a "Filter" section on the left and a "Results" section on the right.

Filter Section:

- Name:** Daniel Gills
- In Fatal Encounters?**: Both
- Published:** Start Date → End Date
- Imported:** 10/14/2017 → 10/20/2017

Results Section:

Collapsible headers: [Collapse all](#) [Uncollapse all](#)

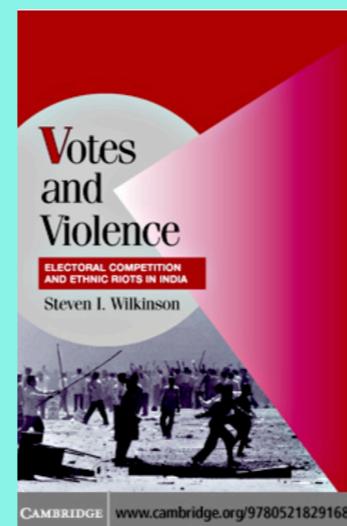
Name (359 capped at 500)	Confidence	Number of Sentences	In Fatal Encounters? (17 342)
Gilbert Flores	1.30	4	In FE
J.C. Hawkins	0.489	1	Not in FE
old J.C. Hawkins Jr. was shot and killed by police on Friday after a sexual assault and robbery at a home on Riverside Avenue.			
Published: 2017-10-14 Imported: 2017-10-14			
http://www.newsplex.com/content/news/Officers-placed-on-paid-administrative-leave-following-shooting-450911743.html			
Tamir Rice	0.212	1	In FE
David Armstrong	0.111	1	Not in FE
Steve Kemmlein	0.0562	1	Not in FE

- Ongoing work: Fatal Encounters used our monitoring system for weekly updates
- Dozens of cases and updates found

My ongoing and future work on supporting *counterdata* collection

My ongoing and future work on supporting *counterdata* collection

Using news articles to automatically detect police actions during communal violence in India



kaggle Open Data Research Grant



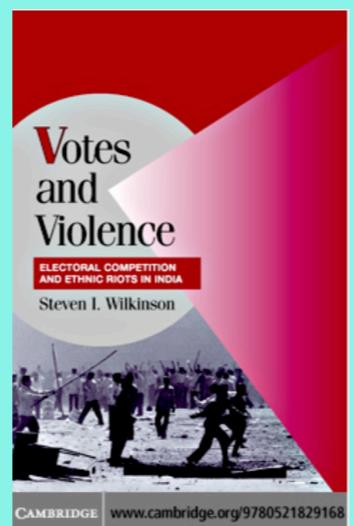
Andrew Halterman
Political Science,
MIT



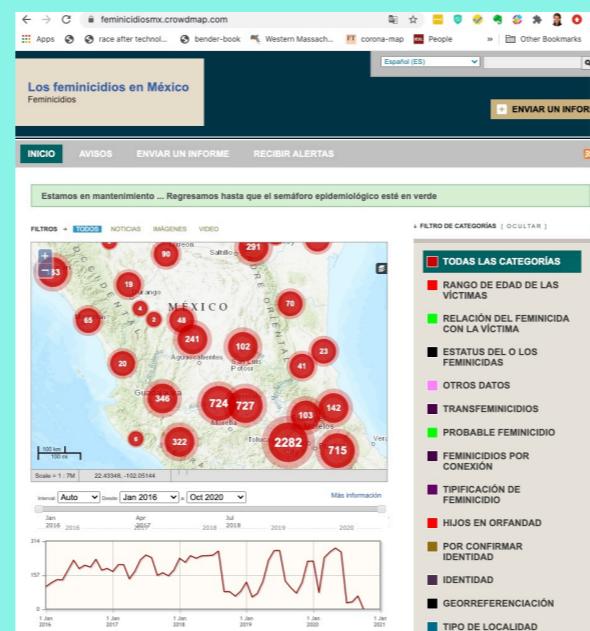
Sheikh Sarwar
Computer Science,
UMass Amherst

My ongoing and future work on supporting *counterdata* collection

Using news articles to automatically detect police actions during communal violence in India



Building a broad set of NLP tools to augment human *counterdata* collection



kaggle Open Data Research Grant



Andrew Halterman
Political Science,
MIT

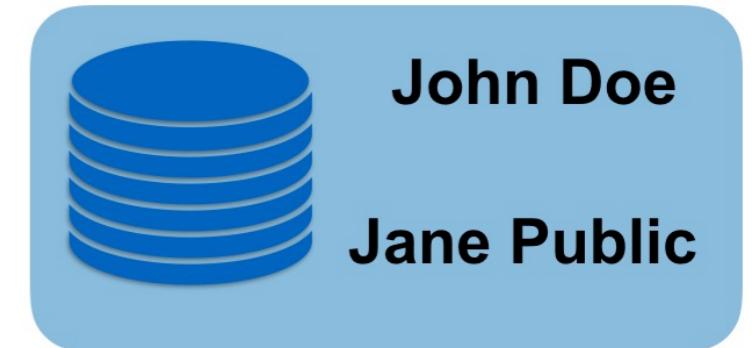


Sheikh Sarwar
Computer Science,
UMass Amherst

e.g. Maria Salguero
manually maps **femicides**
in Mexico

<https://feminicidiosmx.crowdmap.com>

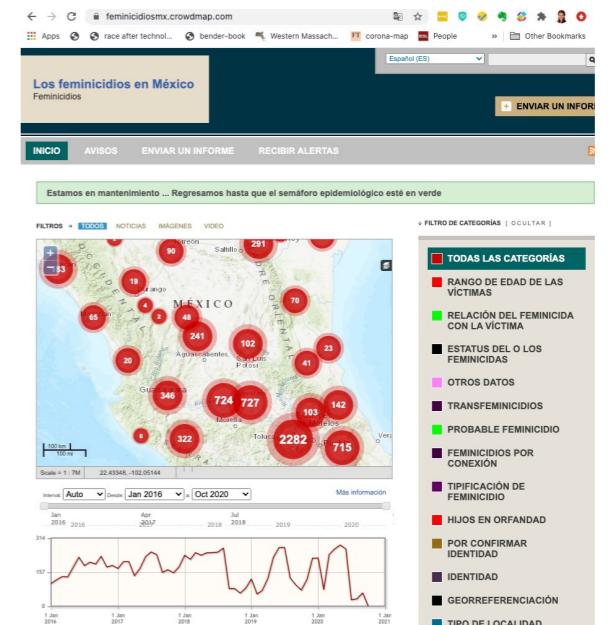
We can reduce the **annotation burden** via methods such as distant supervision and its variants



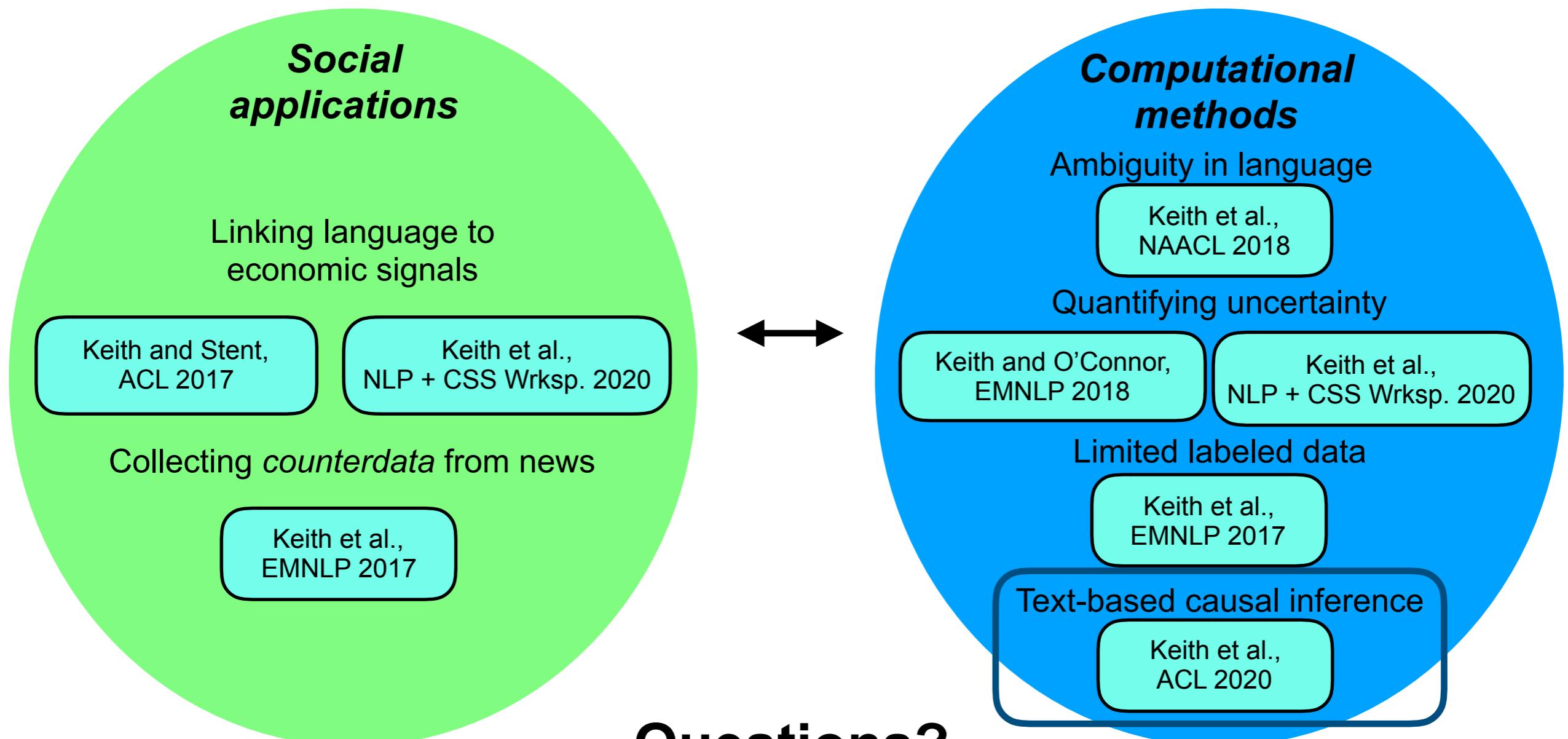
could help augment



Other counterdata collection efforts that are **currently done manually**



My Research Agenda



How can **data science** contribute to **social impact?**

How can we improve outcomes in the world?

How can **data science** contribute to **social impact?**

How can we *improve outcomes* in the world?

values +
measurement

How can **data science** contribute to **social impact?**

How can we improve outcomes in the world?

causal

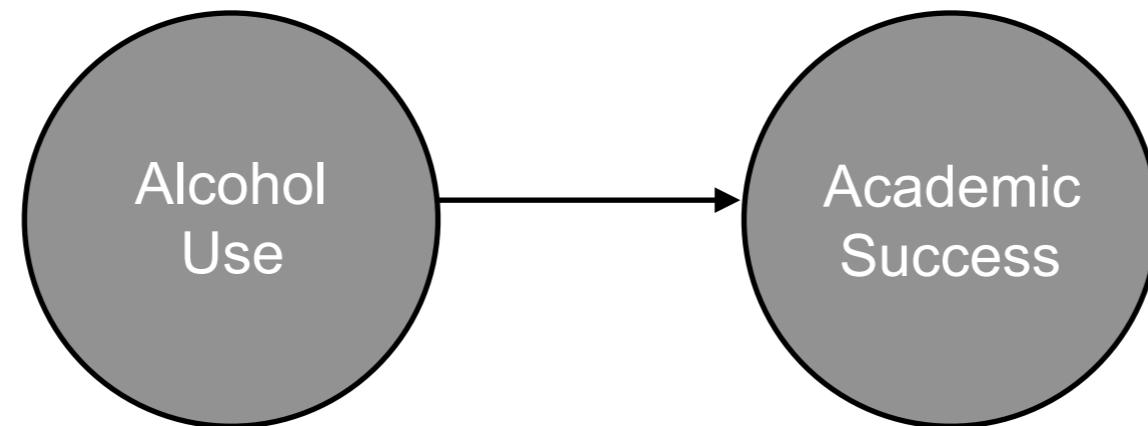
Causality



How can we improve outcomes in the world?

How can we improve outcomes in the world?

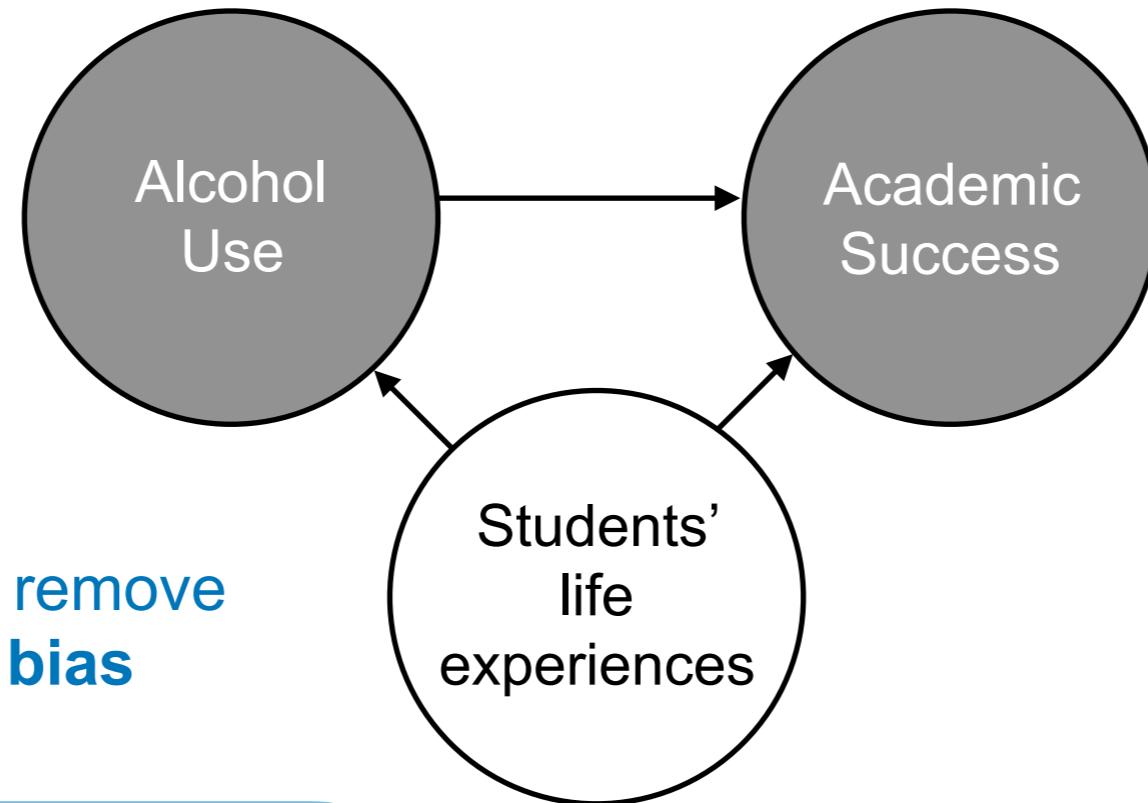
For college students, what is the effect of alcohol use on academic success?



(Kiciman et al. Using longitudinal social media analysis to understand the effects of early college alcohol use. ICWSM, 2020)

How can we improve outcomes in the world?

For college students, what is the effect of alcohol use on academic success?



Problem: need to remove
confounding bias

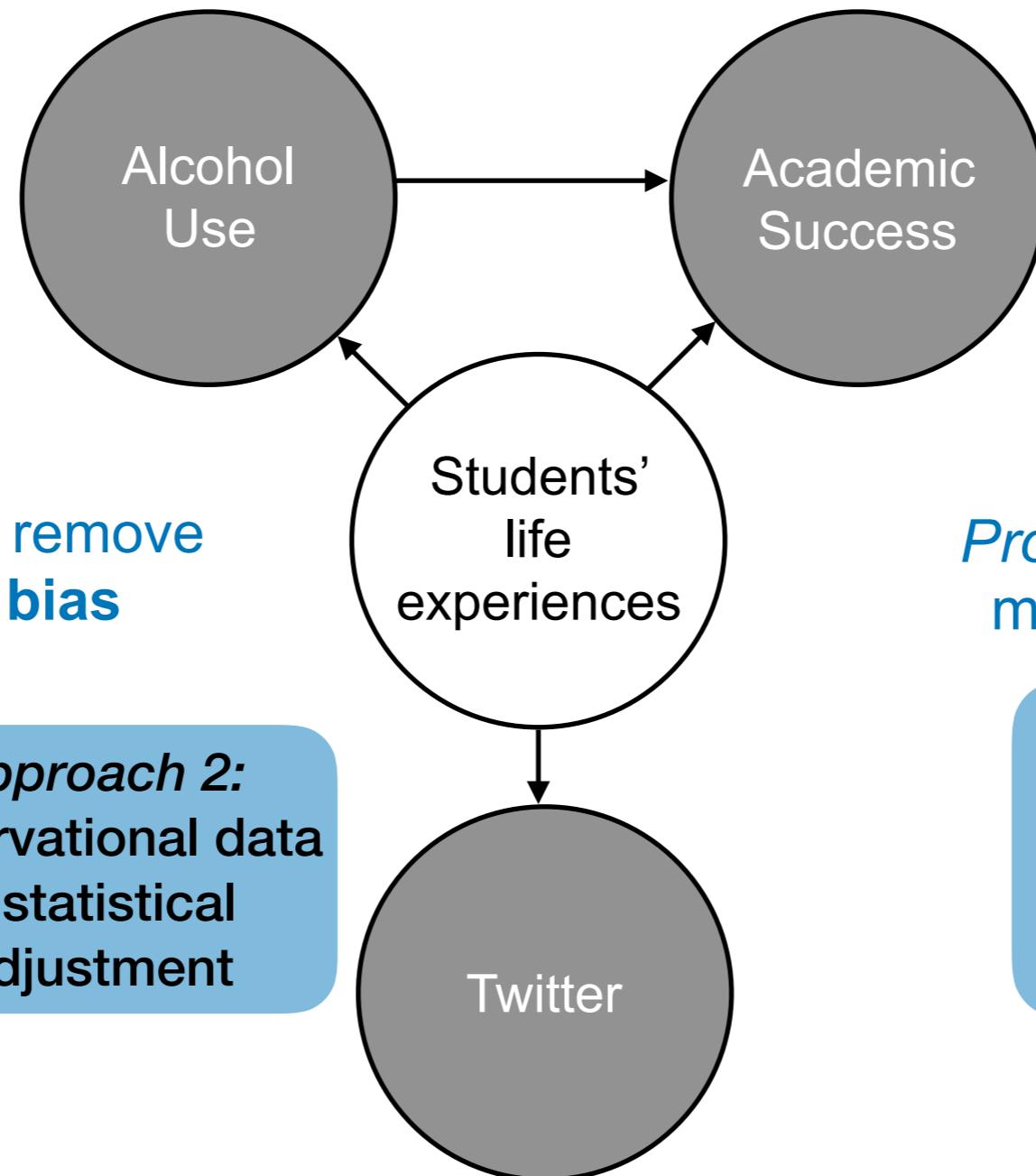
Approach 1:
Intervention

Approach 2:
observational data
+ statistical
adjustment

(Kiciman et al. Using longitudinal social media analysis to understand the effects of early college alcohol use. ICWSM, 2020)

How can we improve outcomes in the world?

For college students, what is the effect of alcohol use on academic success?



Problem: need to remove
confounding bias

Approach 1:
Intervention

Approach 2:
observational data
+ statistical
adjustment

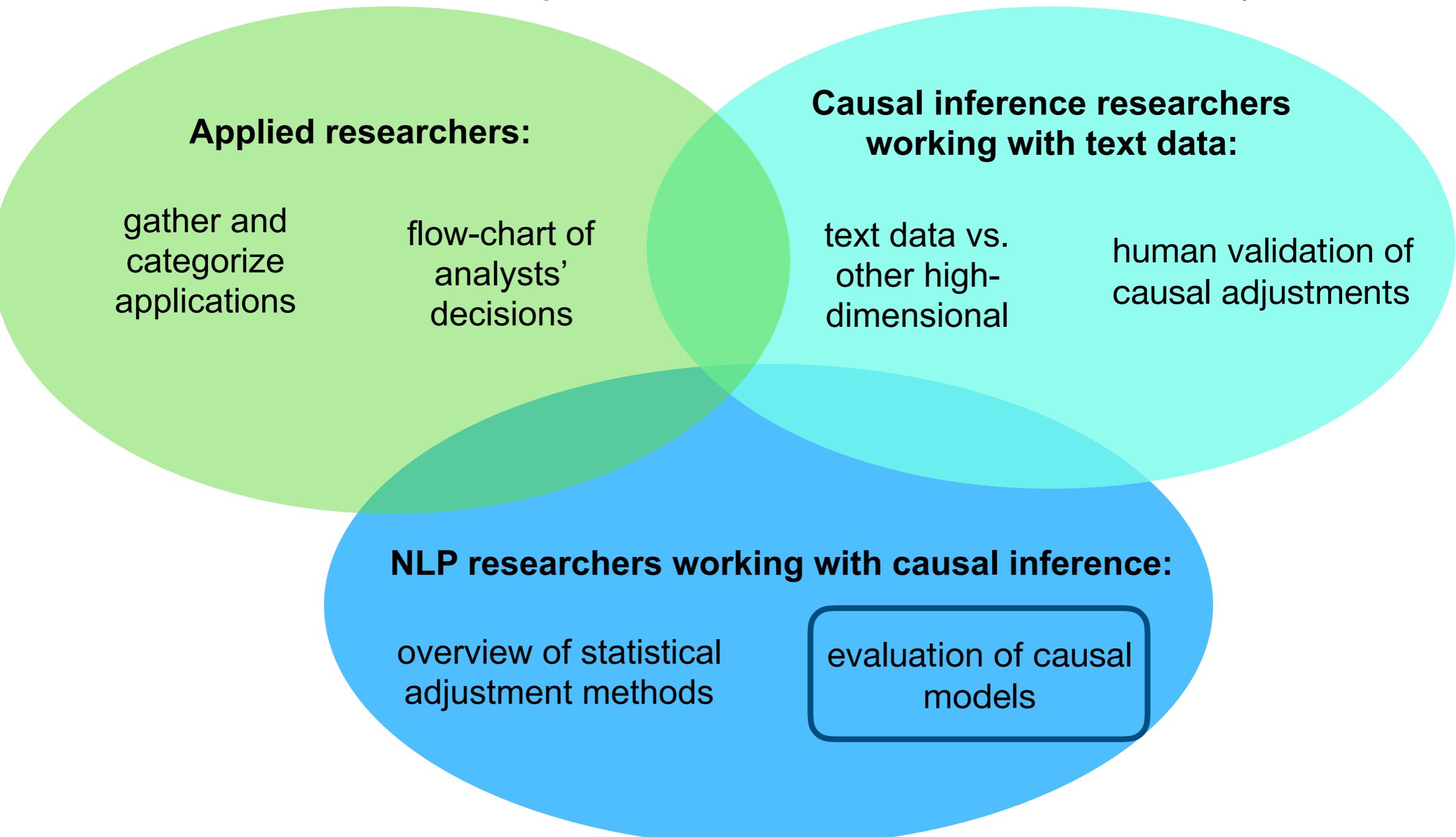
Problem: cannot directly
measure confounders

One approach:
Use text as a
surrogate for
confounders

(Kiciman et al. Using longitudinal social media analysis to understand the effects of early college alcohol use. ICWSM, 2020)

How does one use text to adjust for confounding?

(Keith et al. "Text and Causal Inference: A Review of Using Text to Remove Confounding from Causal Estimates." ACL, 2020)



Future work: evaluating text-based causal methods

Problem Type

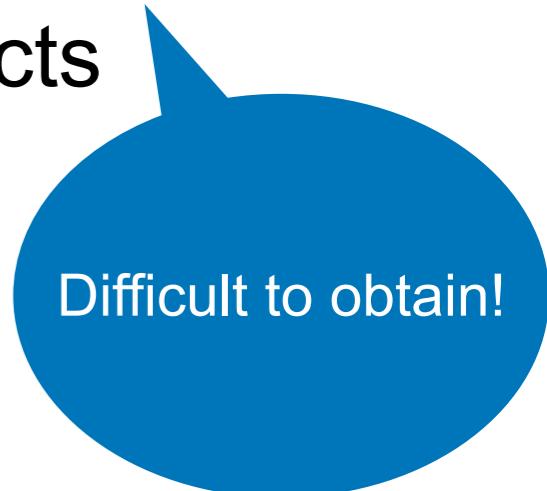
Predictive

Causal

Evaluation

Predictive performance
(e.g. accuracy) on a
held-out test set

Estimated vs. true
causal effects



Difficult to obtain!

Future work: evaluating text-based causal methods

(A) Constructed observational studies

Randomized

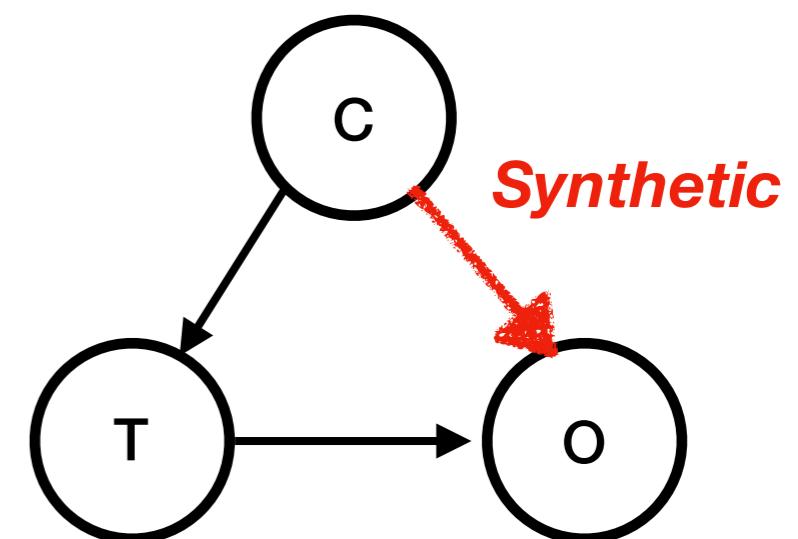


Non-randomized



In other social sciences:
(LaLonde (1986); Shadish et al. (2008); Glynn and Kashin (2013))

(B) Semi-synthetic datasets



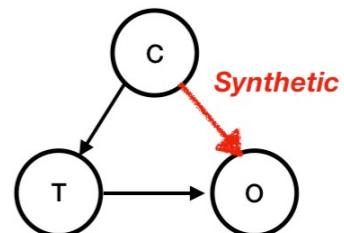
With text to remove confounding:
(Johansson et al. 2016; Veitch et al. 2019; Roberts et al. 2020)

Many **open problems** in text-based causal inference

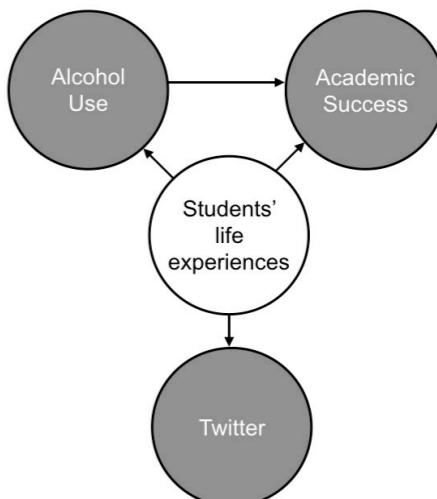


needed to answer

Causal evaluation



Socially impactful text-based causal questions



How can data science contribute to **social impact**?

How can we *improve outcomes* in the world?

causal

values +
measurement

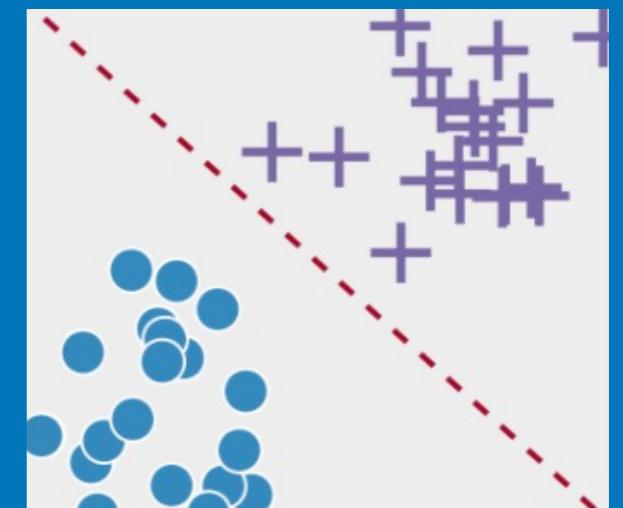
(1) Interdisciplinary
collaborations



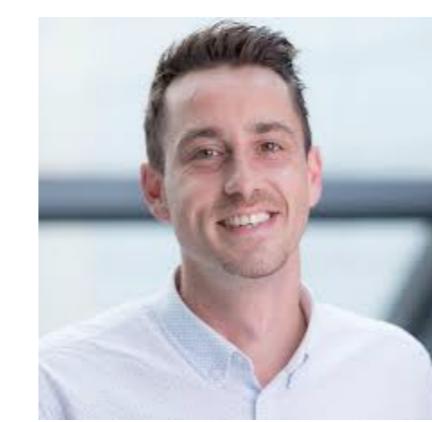
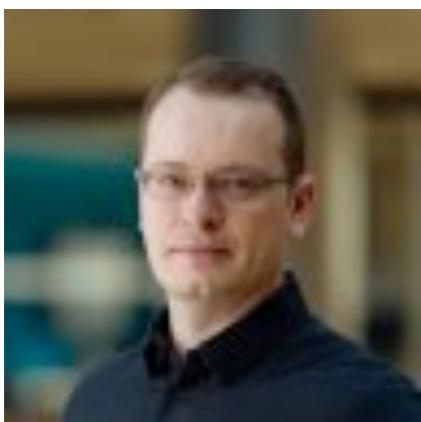
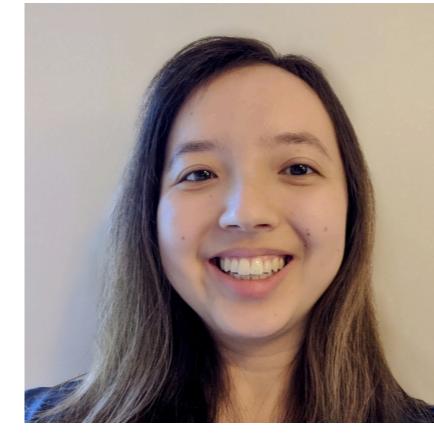
(2) Gathering
additional text
data sources



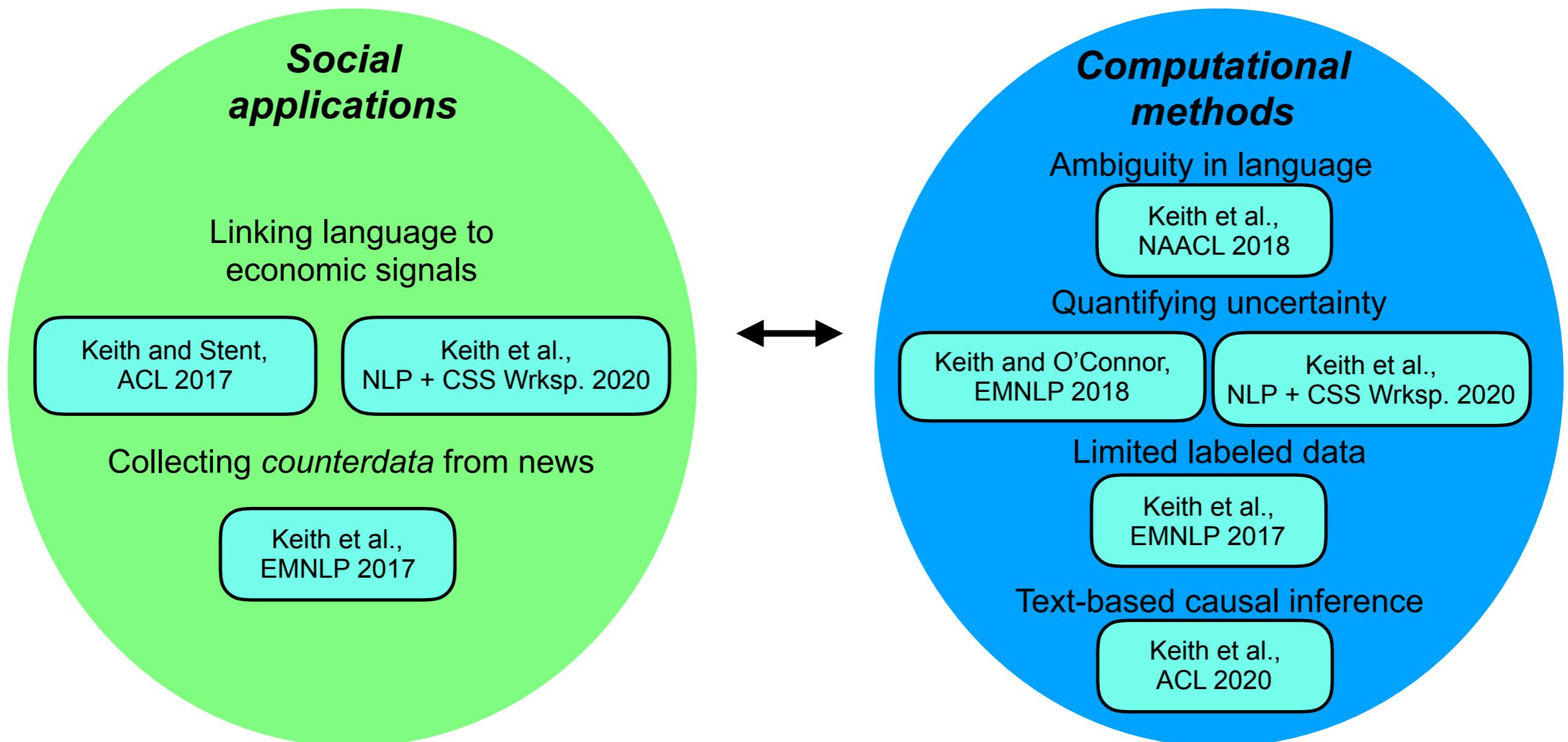
(3) Improving
computational methods



Thank you to my co-authors!



My Research Agenda



Thank you! Questions?

