# CHARACTERIZING THE ECOSYSTEM OF IDEAS IN TEXTS

*Chenhao Tan*

Department of Computer Science
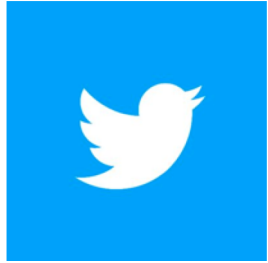University of Colorado Boulder
https://chenhaot.com
chenhao@chenhaot.com

# TEXT CORPORA ARE UBIQUITOUS

# TEXT CORPORA ARE UBIQUITOUS

# TEXT IS VALUABLE

# TEXT IS VALUABLE

- Text allows us to measure a proxy of the world
  - Topics of discussion [Grimmer et al. 2014; Nguyen et al. 2012]
  - Attitude & opinion [Dodds and Danforth 2009; Stuart and Young, 2012]
  - History of ideas [Hall et al. 2008; Uzzi et al. 2013]
  - Bias & polarization [Garg et al. 2018; Niculae et al. 2015; Tan et al. 2018]
- Text allows us to understand social interaction and human behavior
  - Persuasion [Flynn et al. 2017; Tan et al. 2016]
  - Linguistic coordination [DNM et al. 2013; Doyle et al. 2017]
  - Diffusion of norms [Amato et al. 2018; Eisenstein et al. 2014]
  - Mental health [Althoff et al. 2106; Coppersmith et al. 2014]
- Text allows us to understand language

# AUTOMATED TEXT ANALYSIS

# "QUANTITATIVE METHODS FOR TEXT AMPLIFY RESOURCES AND AUGMENT HUMANS."

Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts, Grimmer and Stewart 2013.

# AUTOMATED TEXT ANALYSIS TO AUGMENT HUMANS

what humans can do

what humans cannot do

# AUTOMATED TEXT ANALYSIS TO AUGMENT HUMANS

what humans can do

what humans cannot do

Sentiment analysis

Determine the topic of a sentence among 15 categories

# SENTIMENT ANALYSIS: IS THIS REVIEW POSITIVE OR NEGATIVE?



**Beds should look like beds**

I ordered this when I was high because I thought it was a giant ice cream sandwich. It's not. It's a bed and not the $150 ice cream sandwich I wanted

# SENTIMENT ANALYSIS: IS THIS REVIEW POSITIVE OR NEGATIVE?



⭐☆☆☆☆ **Beds should look like beds**

I ordered this when I was high because I thought it was a giant ice cream sandwich. It's not. It's a bed and not the $150 ice cream sandwich I wanted

# AUTOMATED TEXT ANALYSIS TO AUGMENT HUMANS

what humans can do

what humans cannot do

Sentiment analysis

Determine whether a tweet is going to be viral

Find the perfect argument to change Trump's view

# WHICH TWEET WAS RETWEETED MORE?

**cactus_music**
@cactus_music

Food trucks are the epitome of small independently owned LOCAL businesses! Help keep them going! Sign the petition bit.ly/P6GYCq

**cactus_music**
@cactus_music

I know at some point you've have been saved from hunger by our rolling food trucks friends. Let's help support them! bit.ly/P6GYCq

# WHICH TWEET WAS RETWEETED MORE?

**cactus_music**
@cactus_music

Food trucks are the epitome of small independently owned LOCAL businesses! Help keep them going! Sign the petition bit.ly/P6GYCq

**cactus_music**
@cactus_music

I know at some point you've have been saved from hunger by our rolling food trucks friends. Let's help support them! bit.ly/P6GYCq

# AUTOMATED TEXT ANALYSIS TO AUGMENT HUMANS

what humans can do

what humans cannot do

Sentiment analysis

Determine whether a tweet is going to be viral

Find the key ideas in a million documents

Find out idea relations in a million documents

# AUTOMATED TEXT ANALYSIS TO AUGMENT HUMANS

what humans can do          what humans cannot do

- What approach to take
- How to interpret the results
- How to evaluate the results

# *"THERE IS NO GLOBALLY BEST METHOD FOR AUTOMATED TEXT ANALYSIS."*

Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts, Grimmer and Stewart 2013.

# AUTOMATED TEXT ANALYSIS TO AUGMENT HUMANS

what humans can do

what humans cannot do

Sentiment analysis

Determine whether a tweet is going to be viral

Find out idea relations in a million documents

# AUTOMATED TEXT ANALYSIS TO AUGMENT HUMANS

what humans can do

what humans cannot do

Find out idea relations in a million documents

Ecosystems of ideas

# THE ECOSYSTEM OF IDEAS

- Competition and collaboration
  - Natural selection [Dawkins 1976]
  - Marketplace of ideas [Milton 1644; Mill 1859]
    *Tan, Card, Smith, ACL'17*

- Evolution of ideas
  - Telephone game
    *Tan, Friggeri, Adamic, ICWSM'16*

# THE ECOSYSTEM OF IDEAS

- Competition and collaboration
    - Natural selection [Dawkins 1976]
    - Marketplace of ideas [Milton 1644; Mill 1859]
    *Tan, Card, Smith, ACL'17*
- Evolution of ideas
    - Telephone game
    *Tan, Friggeri, Adamic, ICWSM'16*

# RELATIONS BETWEEN IDEAS

pro-choice    ←— **rivals** —→    pro-life

undocumented immigrants    ←— **rivals** —→    illegal alien

small government    ←— **friends** —→    free market

word alignment    ←— **friends** —→    machine translation

Chong and Druckman, 2007; Dawkins 1976; Entman, 1993; Gitlin, 1980; Lakoff, 2014; Milton 1964

# MAIN CONTRIBUTIONS

First **quantitative** framework to systematically describe relations between ideas

Demonstrate **effective explorations** with this framework on a wide range of datasets

undocumented immigrants ← rivals → illegal alien

small government ← friends → free market

# USING TEXT TO TRACE IDEAS



Hall et al. 2008

Our focus is on
relations between ideas.

We will use standard approaches
- Topics from latent Dirichlet
  allocation (Blei et al. 2003)
- Keywords (Monroe et al.
  2008)



Culturomics, Michel et al. 2011

# QUANTITATIVELY DESCRIBE RELATIONS BETWEEN IDEAS

- Given a corpus of documents over time, each document consists of a set of ideas

undocumented immigrants ←— rivals —→ illegal alien

Cooccurrence
   Pointwise mutual information
   [Church and Hanks 1990]

Rarely cooccur

# QUANTITATIVELY DESCRIBE RELATIONS BETWEEN IDEAS

- Given a corpus of documents over time, each document consists of a set of ideas
  - Cooccurrence does not capture which is winning or losing

undocumented immigrants

Pearson correlation

illegal alien

frequency

time

29

# QUANTITATIVELY DESCRIBE RELATIONS BETWEEN IDEAS

- Given a corpus of documents over time, each document consists of a set of ideas

Cooccurrence    **&**    Prevalence correlation

Within-document

Across-document

# HEAD-TO-HEAD
## (ANTI-CORRELATED, RARELY COOCCUR)



immigrant, undocumented

illegal, alien

1980    1990    2000    2010

31

# RELATIONS BETWEEN IDEAS

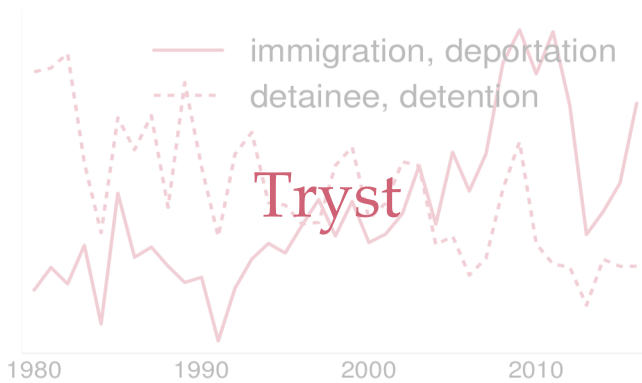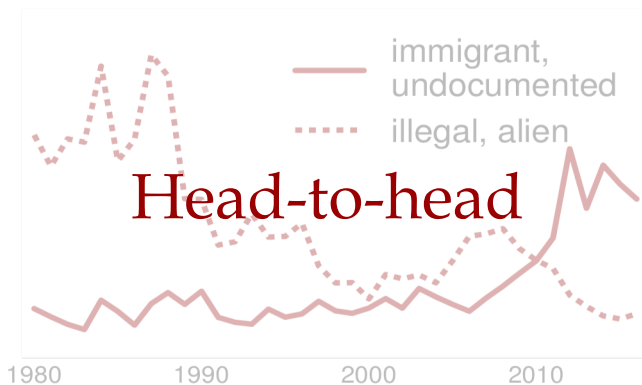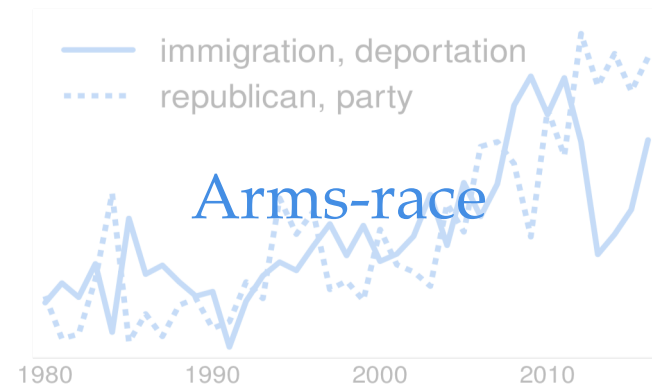Always cooccur

Anti-correlated — Correlated

Rarely cooccur

Tryst
immigration, deportation
detainee, detention

Friendship
immigrant, undocumented
obama, president

Head-to-head
immigrant, undocumented
illegal, alien

Arms-race
immigration, deportation
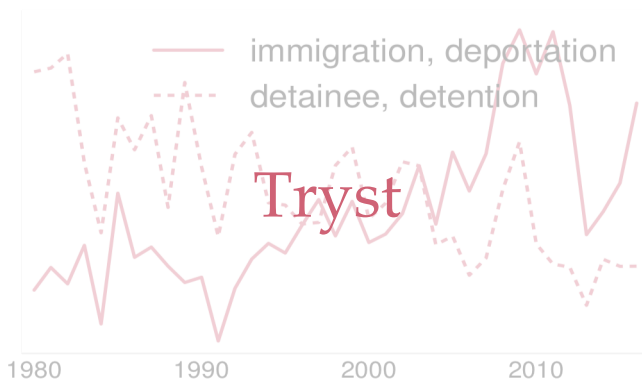republican, party

# RELATIONS BETWEEN IDEAS
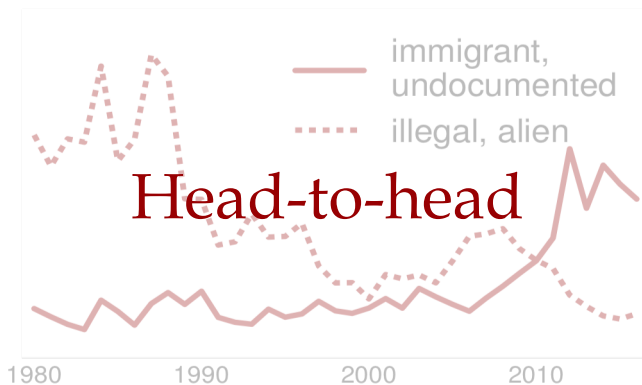
Always cooccur

Anti-correlated

Correlated

immigration, deportation
detainee, detention

Tryst

immigrant,
undocumented
obama, president

Friendship

immigrant,
undocumented
illegal, alien

Head-to-head

immigration, deportation
republican, party

Arms-race

1980    1990    2000    2010

Rarely cooccur

# FRIENDSHIP
## (CORRELATED, LIKELY TO COOCCUR)



immigrant, undocumented

····· obama, president

1980    1990    2000    2010

# RELATIONS BETWEEN IDEAS

Always cooccur

immigration, deportation
detainee, detention

Tryst

1980   1990   2000   2010

immigrant,
undocumented
obama, president

Friendship

1980   1990   2000   2010

Anti-correlated ———————————————→ Correlated

immigrant,
undocumented
illegal, alien

Head-to-head

1980   1990   2000   2010

immigration, deportation
republican, party

Arms-race

1980   1990   2000   2010

Rarely cooccur

35

# RELATIONS BETWEEN IDEAS

Always cooccur

Anti-correlated — Correlated

Rarely cooccur

Tryst

immigration, deportation
detainee, detention

Friendship

immigrant, undocumented
obama, president

Head-to-head

immigrant, undocumented
illegal, alien

Arms-race

immigration, deportation
republican, party
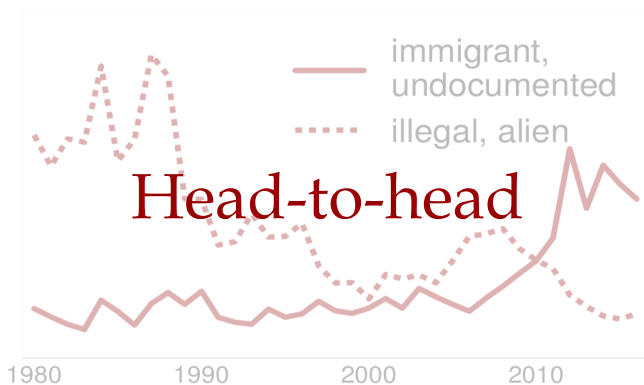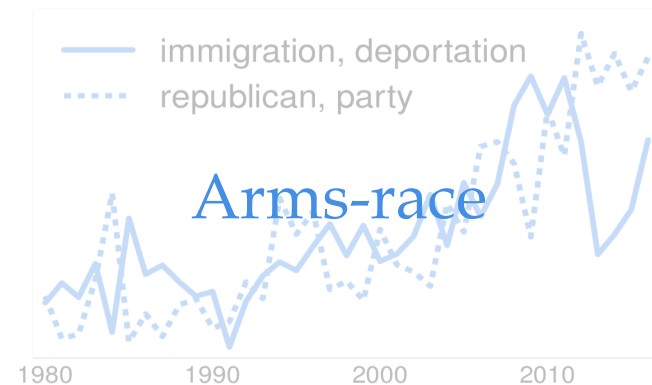
# ARMS-RACE
## (CORRELATED, RARELY COOCCUR)



immigration, deportation
republican, party

1980    1990    2000    2010

# ARMS-RACE
## (CORRELATED, RARELY COOCCUR)



— immigration, deportation

····· republican, party

2000          2010

# RELATIONS BETWEEN IDEAS
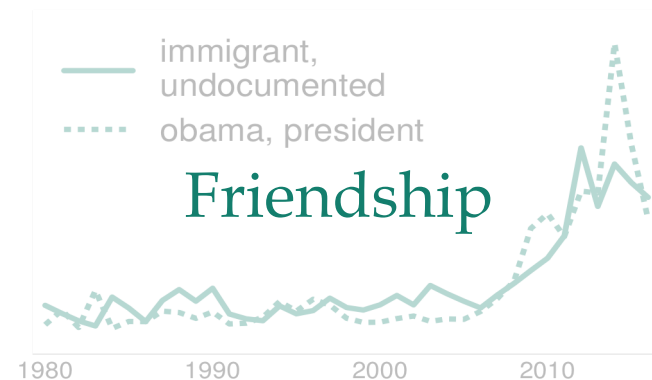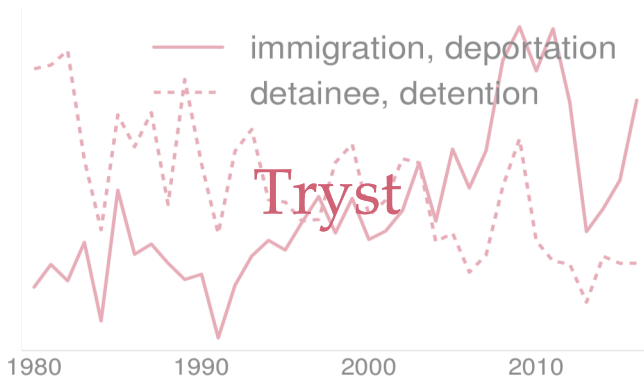
Always cooccur



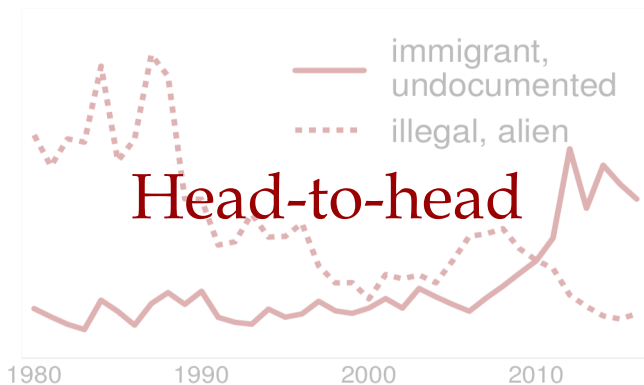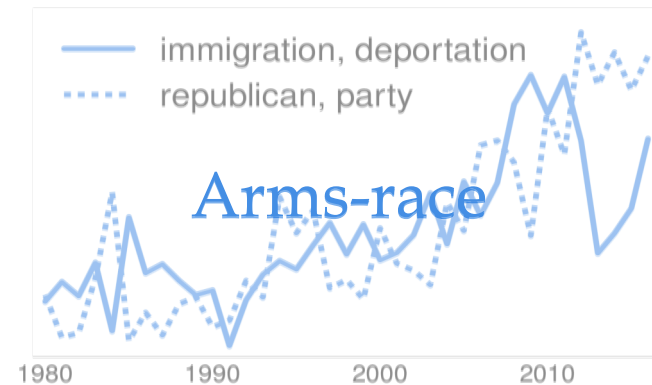Anti-correlated — Correlated

Rarely cooccur

# RELATIONS BETWEEN IDEAS

Always cooccur

Anti-correlated

Correlated

Rarely cooccur

**Tryst**
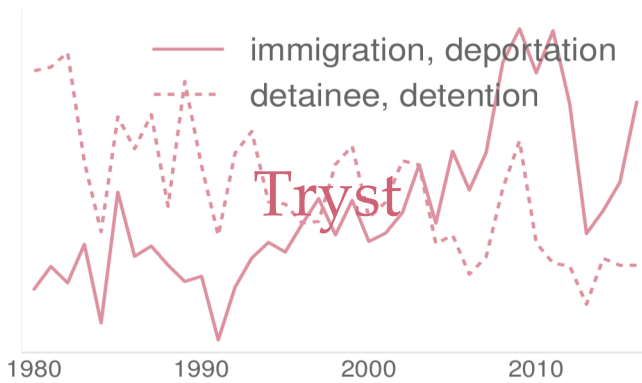immigration, deportation
detainee, detention
1980 1990 2000 2010

**Friendship**
immigrant, undocumented
obama, president
1980 1990 2000 2010

**Head-to-head**
immigrant, undocumented
illegal, alien
1980 1990 2000 2010

**Arms-race**
immigration, deportation
republican, party
1980 1990 2000 2010

# TRYST
## (ANTI-CORRELATED, LIKELY TO COOCCUR)



immigration, deportation
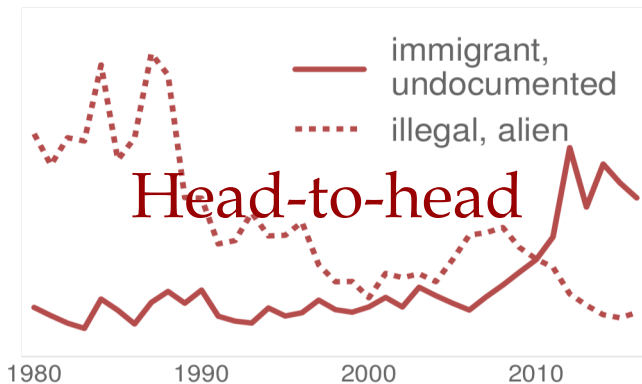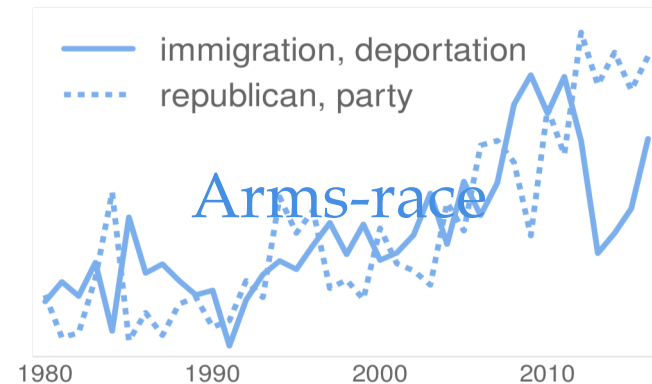detainee, detention

1980    1990    2000    2010

# RELATIONS BETWEEN IDEAS

Always cooccur

Anti-correlated

Correlated

Rarely cooccur

**Tryst** — immigration, deportation; detainee, detention

**Friendship** — immigrant, undocumented; obama, president

**Head-to-head** — immigrant, undocumented; illegal, alien

**Arms-race** — immigration, deportation; republican, party
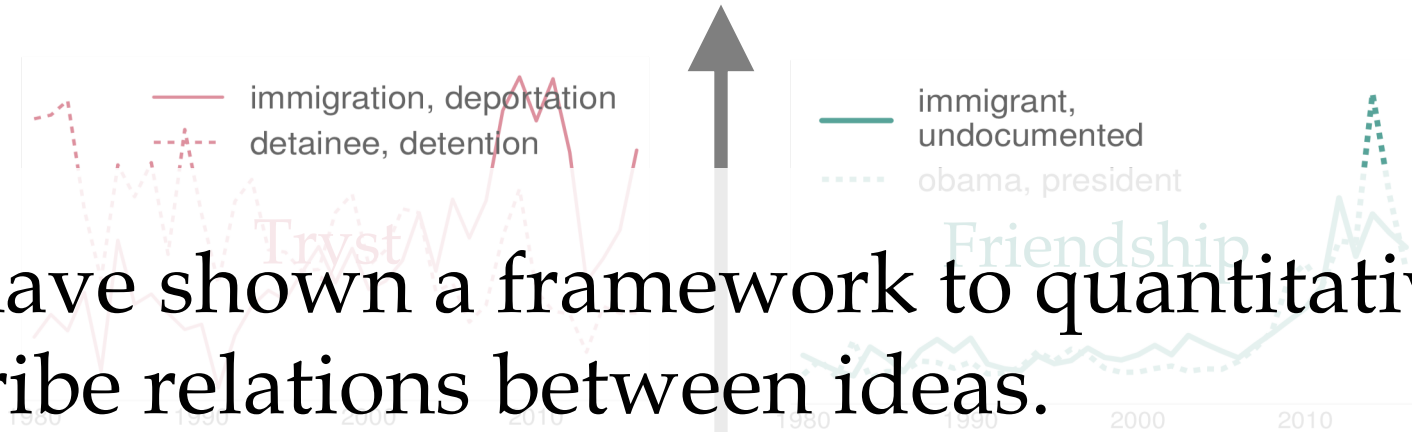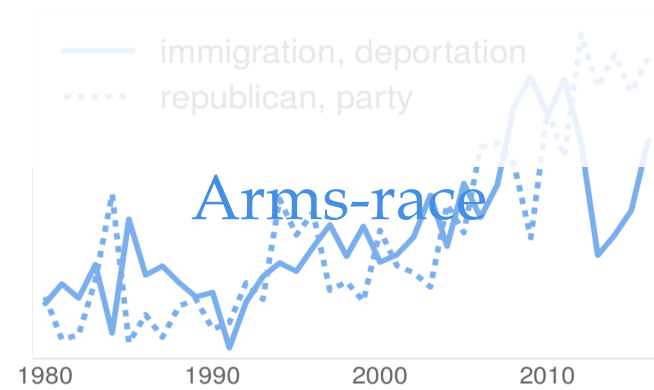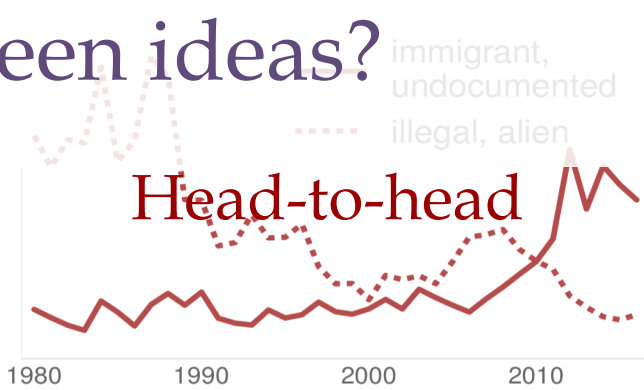
# RELATIONS BETWEEN IDEAS

Always cooccur



We have shown a framework to quantitatively describe relations between ideas.

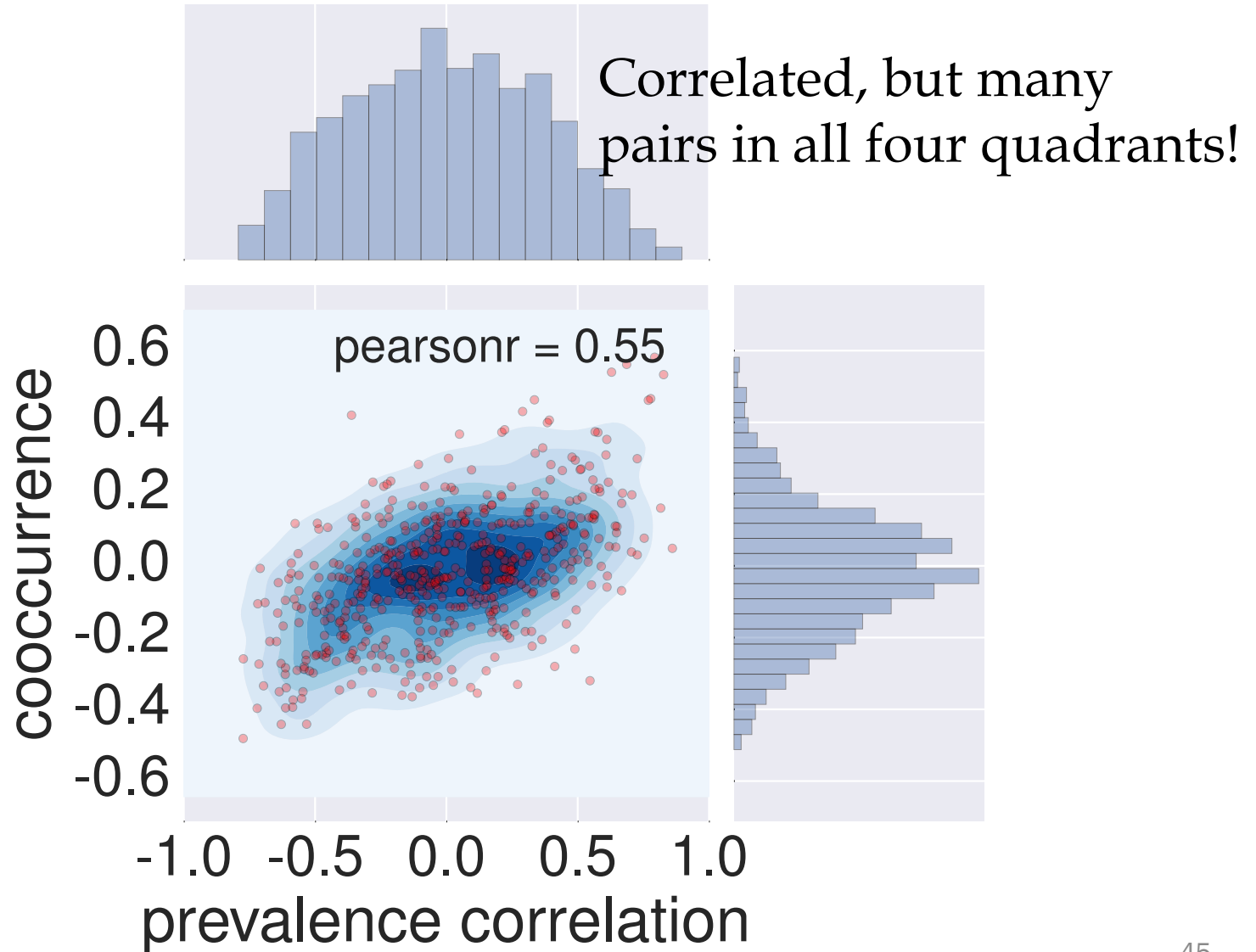Can we use them to effectively explore relations between ideas?

Head-to-head

Arms-race

Rarely cooccur

# A WIDE RANGE OF DATASETS

- Newspapers and research articles as datasets
  - Immigration
  - Terrorism
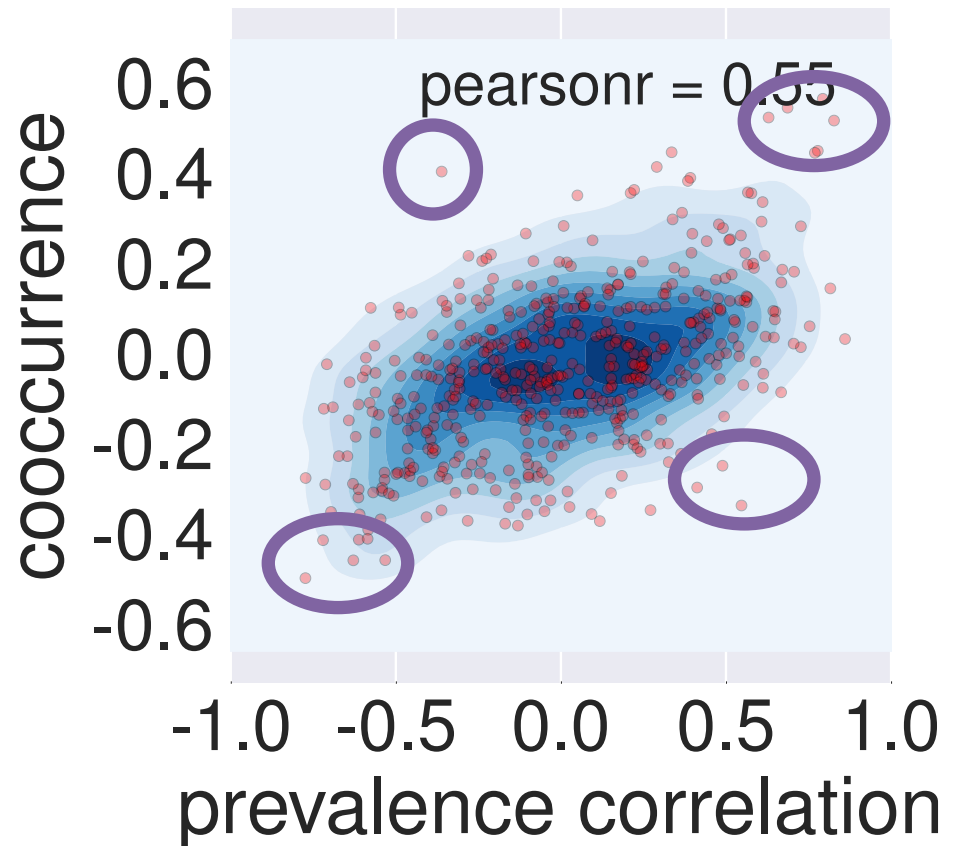  - Same-sex marriage
  - Abortion
  - Tobacco

  - ACL
  - NIPS

# JOINT DISTRIBUTIONS



Correlated, but many pairs in all four quadrants!

# THE STRENGTH OF RELATIONS

Strength = |PMI| × |correlation|

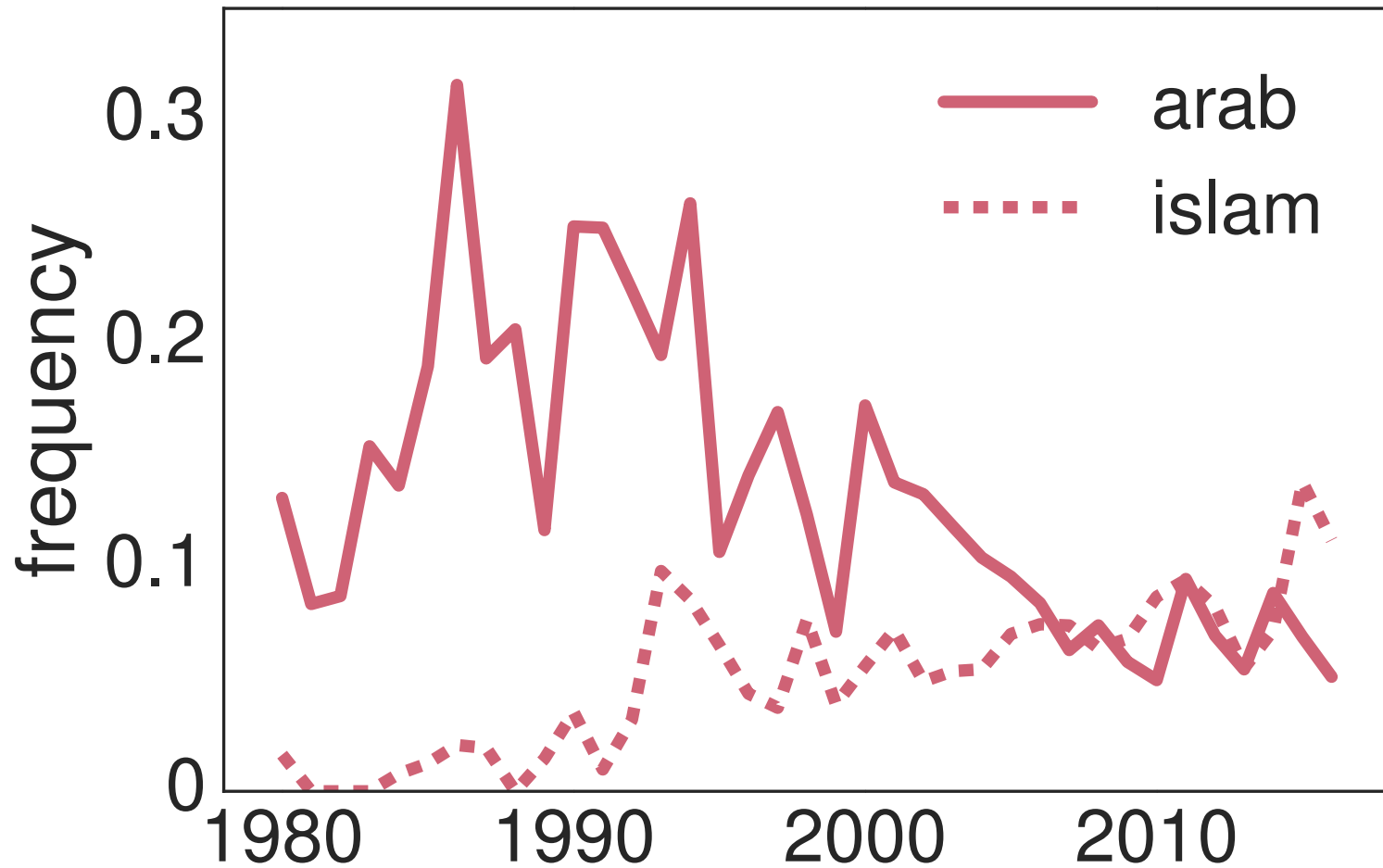Extreme pairs are
the interesting ones!

# EFFECTIVELY EXPLORE RELATIONS BETWEEN IDEAS

- Terrorism
  - Keywords
  - Topics

# EFFECTIVELY EXPLORE RELATIONS BETWEEN IDEAS
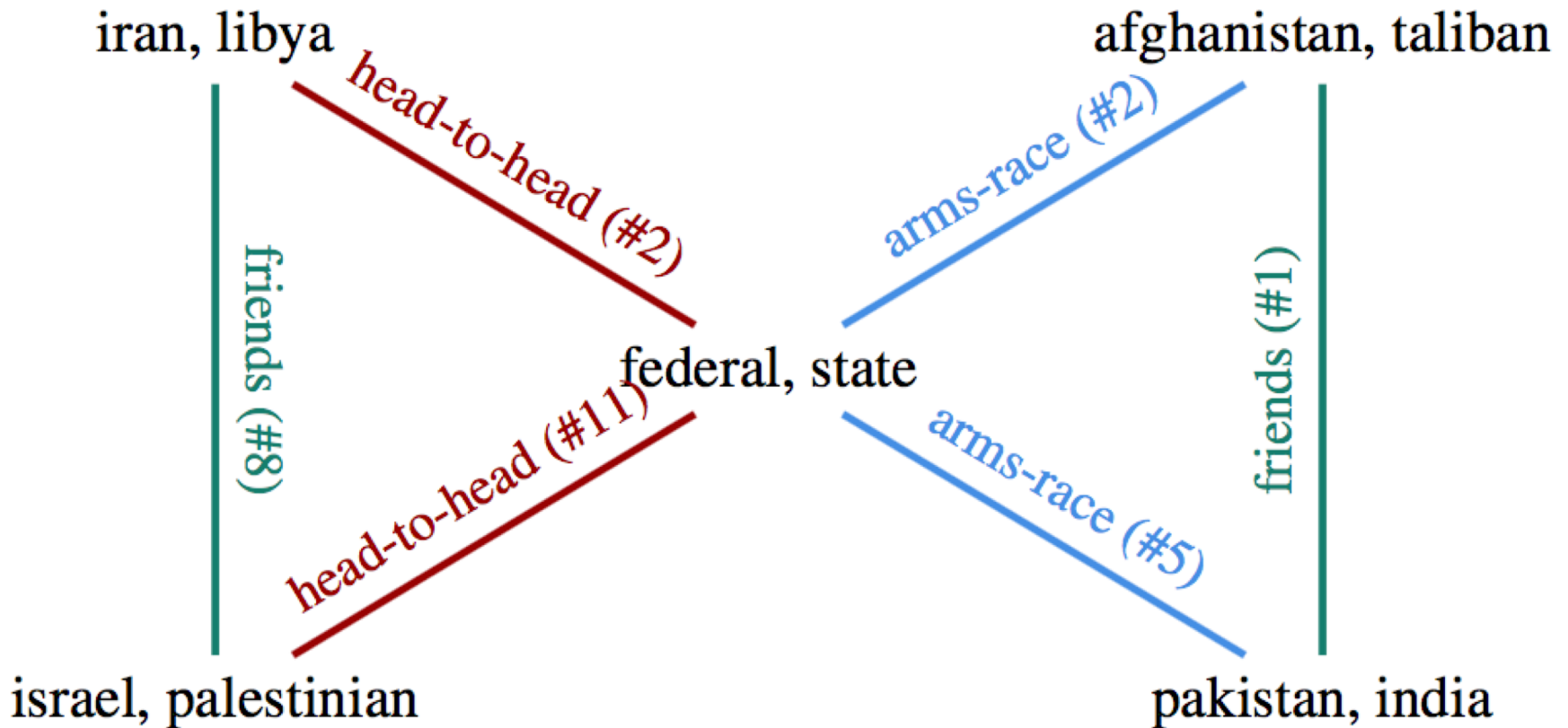
- Terrorism
  - Keywords
  - Topics

# #2 IN TRYSTS

# EFFECTIVELY EXPLORE RELATIONS BETWEEN IDEAS

- Terrorism
  - Keywords
  - Topics

# TOP RELATIONS BETWEEN IDEAS

iran, libya

afghanistan, taliban

head-to-head (#2)

arms-race (#2)

friends (#8)

federal, state

friends (#1)

head-to-head (#11)

arms-race (#5)

israel, palestinian

pakistan, india

The relations between these topics are consistent with structural balance theory: the enemy of an enemy is a friend [Cartwright and Harary, 1956; Heider, 1946]
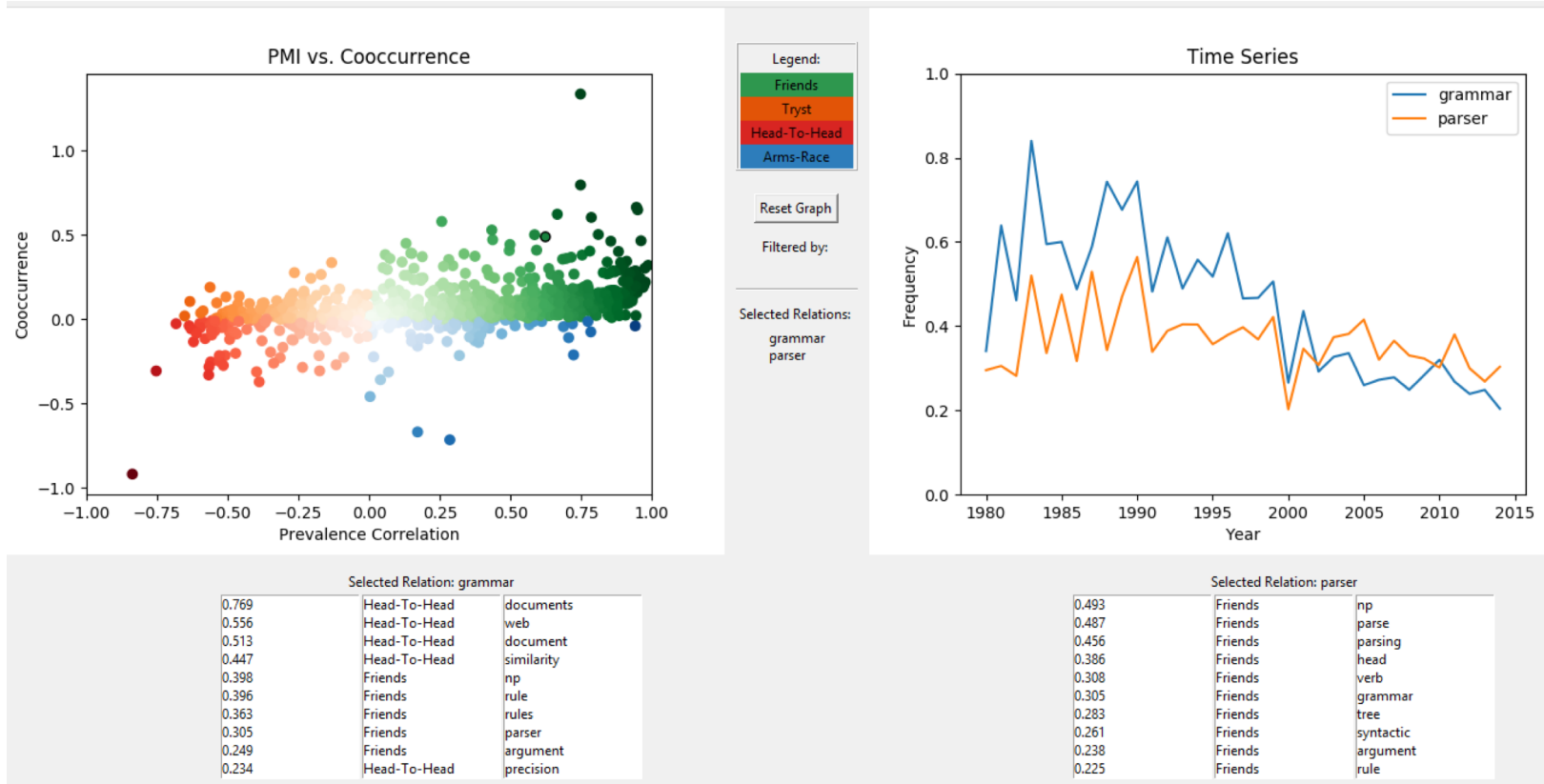
# EFFECTIVE EXPLORATIONS

Rank among all relations

| | | | PMI | Correlation | Joint |
|---|---|---|---|---|---|
| Keywords | arab | islam | 106 | 1,494 | 2 |
| Topics | federal, state | afghanistan, taliban | 43 | 99 | 2 |
| | federal, state | iran, lybia | 36 | 56 | 2 |

The "interesting" pair is ranked much higher according to our framework.

# VISUALIZATION TOOL



https://github.com/Noahs-ARK/idea_relations

https://github.com/nwrush/Visualization

# NEURAL MODELS FOR DOCUMENTS WITH METADATA

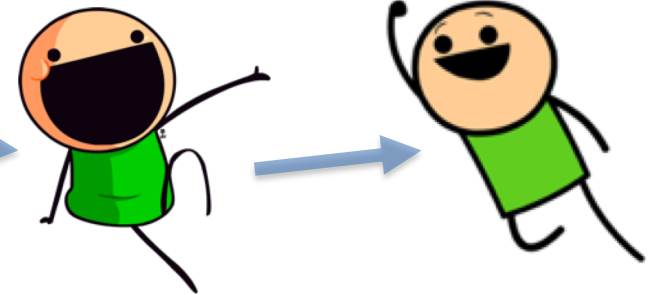| Base topics | Anti-immigration | Pro-immigration |
|---|---|---|
| ice customs agency | **criminal** customs | **detainees** detention |
| population born percent | jobs million **illegals** | english **newcomers** |
| judge case court guilty | **guilty** charges man | **asylum** court judge |
| patrol border miles | patrol border | died authorities desert |
| licenses drivers card | foreign sept visas | green citizenship card |
| island story chinese | smuggling federal | island school ellis |
| guest worker workers | bill border house | workers tech skilled |
| benefits bill welfare | republican california | law welfare students |

https://github.com/dallascard/scholar

# THE ECOSYSTEM OF IDEAS

- Competition and collaboration
  - Natural selection [Dawkins 1976]
  - Marketplace of ideas [Milton 1644; Mill 1859]
    *Tan, Card, Smith, ACL'17*
- Evolution of ideas
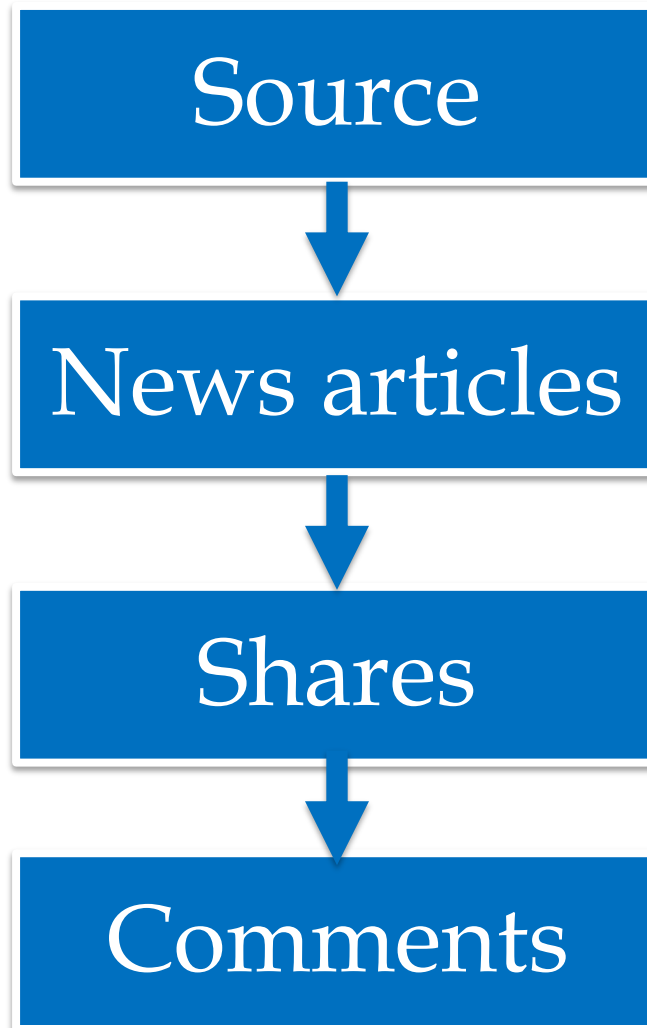  - Telephone game
    *Tan, Friggeri, Adamic, ICWSM'16*

# IDEAS REACH US
# VIA DIFFERENT WAYS
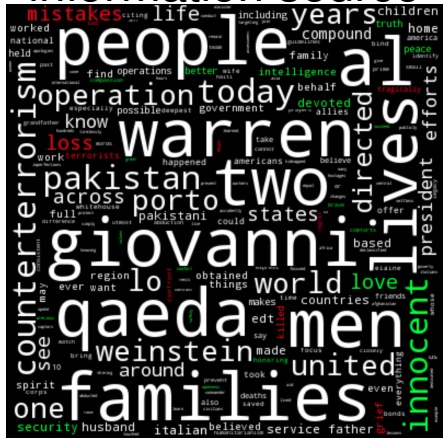
# EVOLUTION OF IDEAS

Source
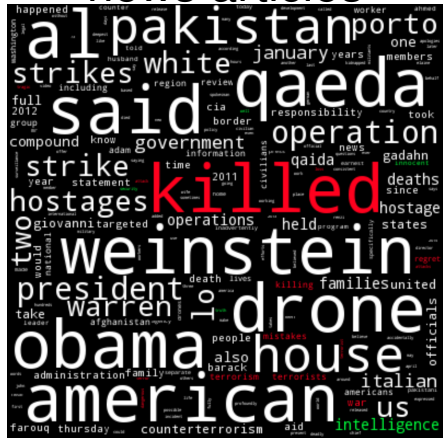
↓

News articles

↓

Shares
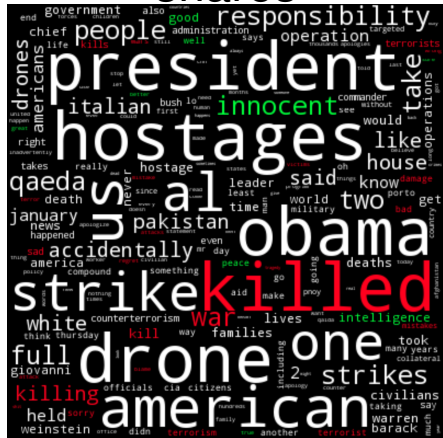
↓

Comments

# OBAMA'S SPEECH ON DEATHS OF WARREN AND GIOVANNI
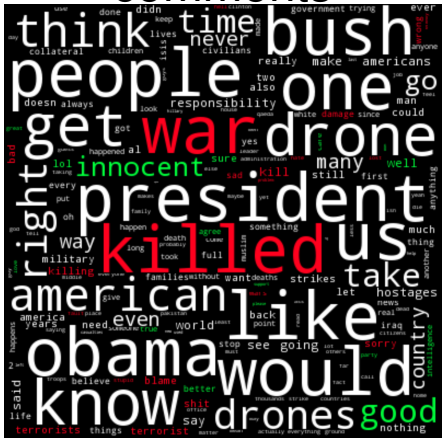
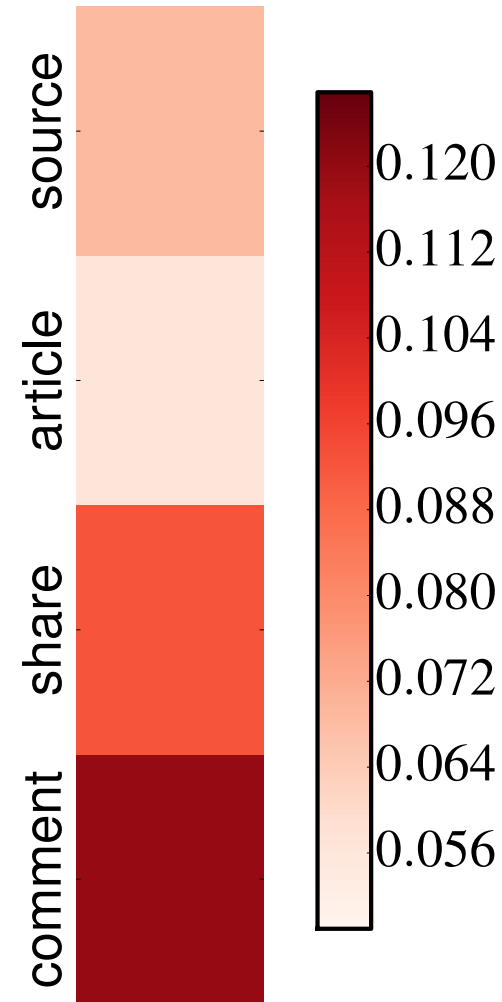

information source
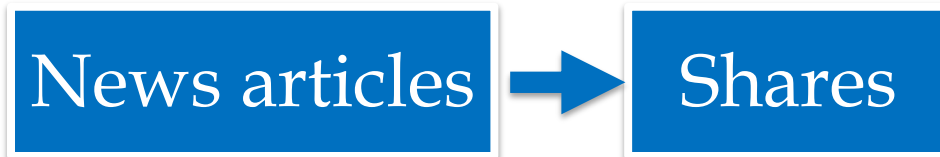
news articles

shares

comments

# SENTIMENT EVOLUTION

- Subjectivity defined as the fraction of emotional words
- News articles are less subjective than sources, but shares and comments get more and more subjective
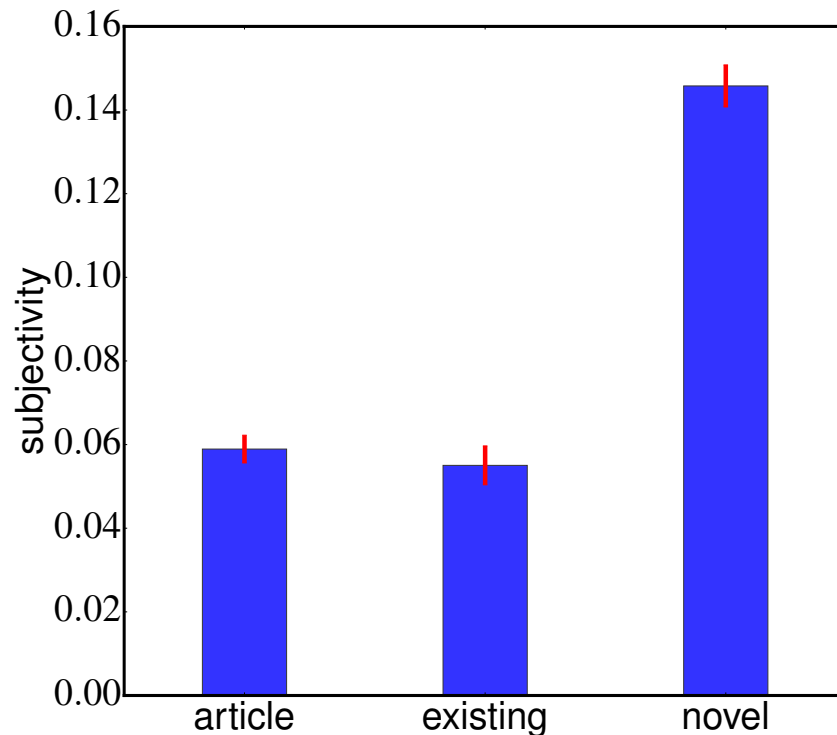


**Red: more** subjective

# WHY DOES SUBJECTIVITY INCREASE?

News articles → Shares

words already in news articles *(magnifier)*

novel words added by individuals *(creator)*

# CHARACTERIZING THE ECOSYSTEM OF IDEAS

better understand
existing ideas in the data


better generate
ideas for the future

# The beginning of an exciting journey!

Chenhao Tan
chenhao@chenhaot.com
https://chenhaot.com
@ChenhaoTan



http://www.xiangke.com/html/news/100284/61936252.html