

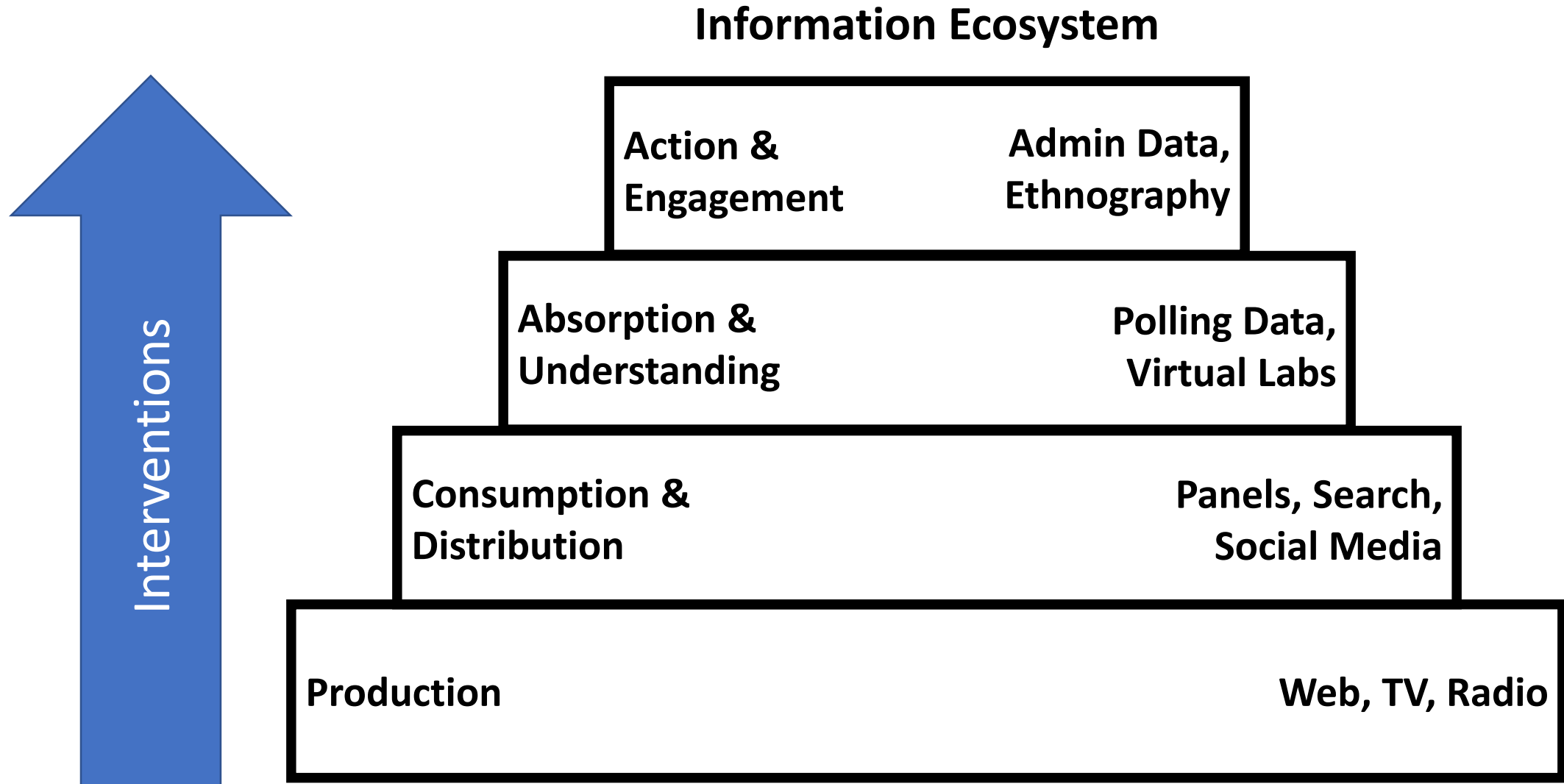
# Project Ratio: problem

- News Datasets: idiosyncratic, one-off, and often small in scale
  - Comparisons difficult across different modes of media consumption, different sample populations, and different times
- Data is scattered between private companies and academic research
  - Bringing that data together is both an infrastructural and cooperative challenge
- Research scattered across several disciplines (e.g. economics, marketing, political science, communications, psychology, sociology, computer science, and network science), each with its own set of theoretical frameworks, accepted methodologies, and publishing venues
  - Collating and reconciling results across these disciplinary boundaries is difficult and often leads to contradictory or incoherent conclusions

# Project Ratio: solution

- Build a large-scale, shared data infrastructure for studying the production, distribution, consumption, absorption, and impact of news over time and across the entire information ecosystem
- Develop academic-industry partnerships around data and solutions, and communicate these solutions to the public
- Build “many labs” research network to advance basic research and generate actionable insights that are relevant to business and public policy

# Project Ratio



# Project Ratio

- **Production:** What news is actually being produced?
  - Bias in the selection and framing of articles/events/topics
  - Provenance of articles/events
  - Networks of article/story sharing and updating
- **Consumption:** What news is actually being consumed?
  - How much news do people get
  - Where do they get their news
  - How much of it is mainstream, fake, commentary, hyper-partisan
  - Do people live in filter bubbles
- **Absorption:** What news are people actually absorbing?
- **Action:** How does absorbing news change actions or engagement?

# production data

News Scraping (since 8/2014)  
(3,000 US and 4,500 European pub.)



Metadata Generation  
(events, provenance)



Research

Web Tools

- 7,500 publishers/ 500,000 articles (top 300 publishers' home pages every 30 mins)
- Landing pages/placement of each article on the page/font-size
- Recording transcripts of all national TV channel and relevant local markets ✓

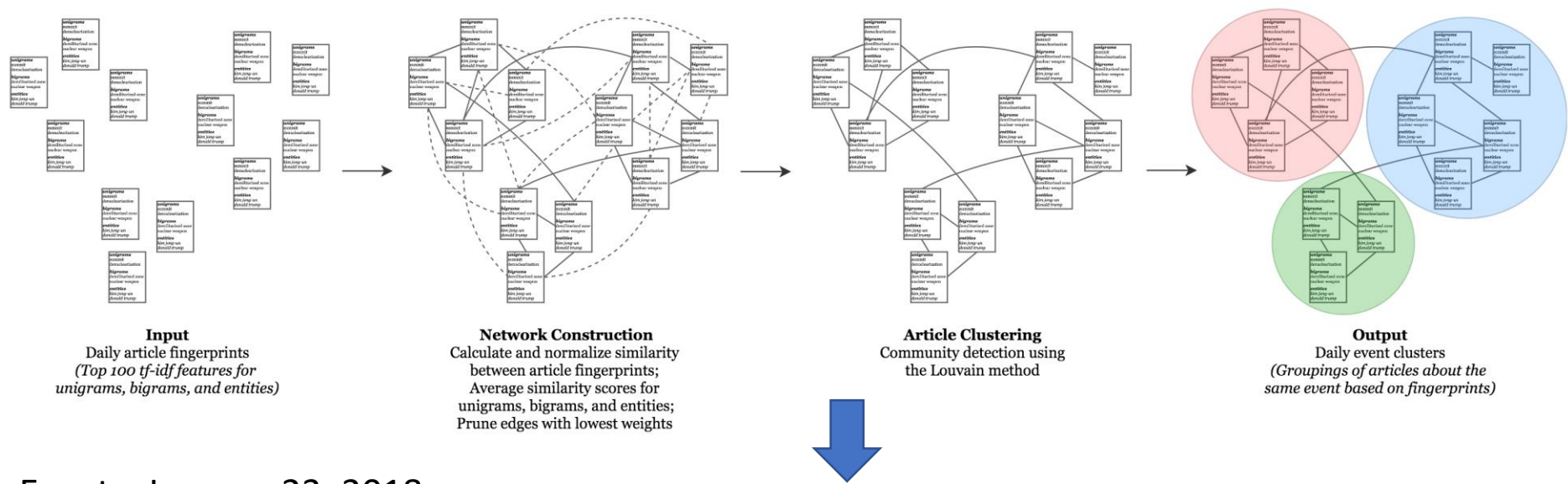
- We partition daily newly published articles into “events”
  - Within events, we analyze the provenance of articles:
    - Text from agencies such as AP and Syndicated articles
    - Original content
- ✓

- Annotated raw data and meta data becomes available for research.
- Aggregated meta data is available through a public website. ✓

# production pipeline



# production pipeline



We partition articles into daily events by fingerprinting them and then clustering them.

## Events: January 23, 2018

|    | A       | B                                  | C          | D           | E           | F             | G  | H | I | J | K | L |
|----|---------|------------------------------------|------------|-------------|-------------|---------------|--|---|---|---|---|---|
| 1  | eventID | label1                             | label2     | label3      | article sha | words         | bigrams  |   |   |   |   |   |
| 2  | 0       | 'Logan' Oscar Nominations: Film    | Here Are t | Oscars no   | 3.32%       | nomin, os     | phantom thread, ladi bird, oscar nomin, eb missouri, film acade      |   |   |   |   |   |
| 3  | 1       | Kentucky shooting: Bevin to discu  | Kentucky s | Kentucky s  | 2.64%       | school, sh    | marshal counti, school shoot, benton ky, fatal school, paducah s     |   |   |   |   |   |
| 4  | 2       | Senate Votes Overwhelmingly to     | Senate vo  | Senate Co   | 2.05%       | senat, deng   | govern shutdown, white hous, border wall, feb 8, senat democr        |   |   |   |   |   |
| 5  | 3       | National Tsunami Center cancels    | Alaska ear | Alaska ear  | 1.79%       | tsunami, a    | tsunami warn, kodiak island, higher ground, british columbia, cc     |   |   |   |   |   |
| 6  | 4       | Sentencing nears end for former    | USA Gymr   | 'Enjoy hell | 1.55%       | nassar, gy    | usa gymnast, larri nassar, matthew dae, michigan state, victim ir    |   |   |   |   |   |
| 7  | 5       | Mueller's office spoke with Sessic | Mueller's  | Mueller's   | 1.12%       | mueller, si   | russia investig, justic depart, attorney gener, session interview, n |   |   |   |   |   |
| 8  | 6       | Tecmo Super Bowl predicts the S    | Philadelph | Patriots-E  | 1.01%       | bowl, sup     | super bowl, eagl fan, bowl lii, tom bradi, nfl footbal, england pai  |   |   |   |   |   |
| 9  | 7       | Trump slaps new trade tariffs on   | Trump sig  | Trump app   | 0.89%       | solar, tarif  | wash machin, solar panel, solar cell, solar power, solar industri,   |   |   |   |   |   |
| 10 | 8       | UN readies aid as Turkey attacks   | Turkey ba  | Erdogan S   | 0.85%       | syria, syria  | chemic weapon, northern syria, syrian kurdish, chemic attack, le     |   |   |   |   |   |
| 11 | 9       | FBI Wray Threatens To Resign If S  | Under Crit | Trump Rej   | 0.83%       | fbi, mccab    | fbi director, jame comei, andrew mccabe, white hous, act fbi, jin    |   |   |   |   |   |
| 12 | 10      | Trump on FBI Agent Strzok's Miss   | Sessions L | #DeepStat   | 0.81%       | fbi, text, st | text messag, secret societi, peter strzok, fbi agent, miss text, gut |   |   |   |   |   |

Search through these events to find school shootings.

# production of real news





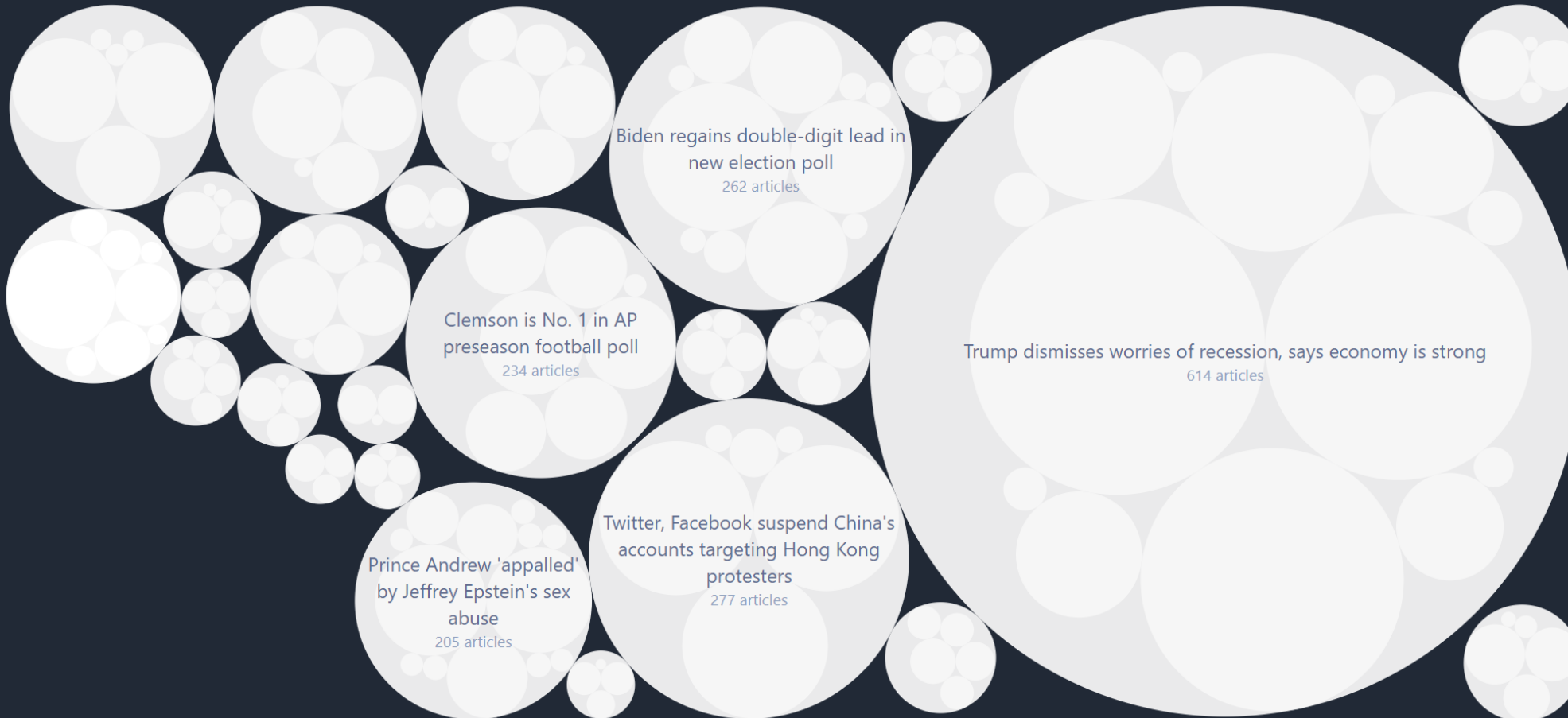
# production dashboard

Project Ratio [About this project](#)



► August 19, 2019

Top 25 News Events  
3572 articles



← Expand Sidebar

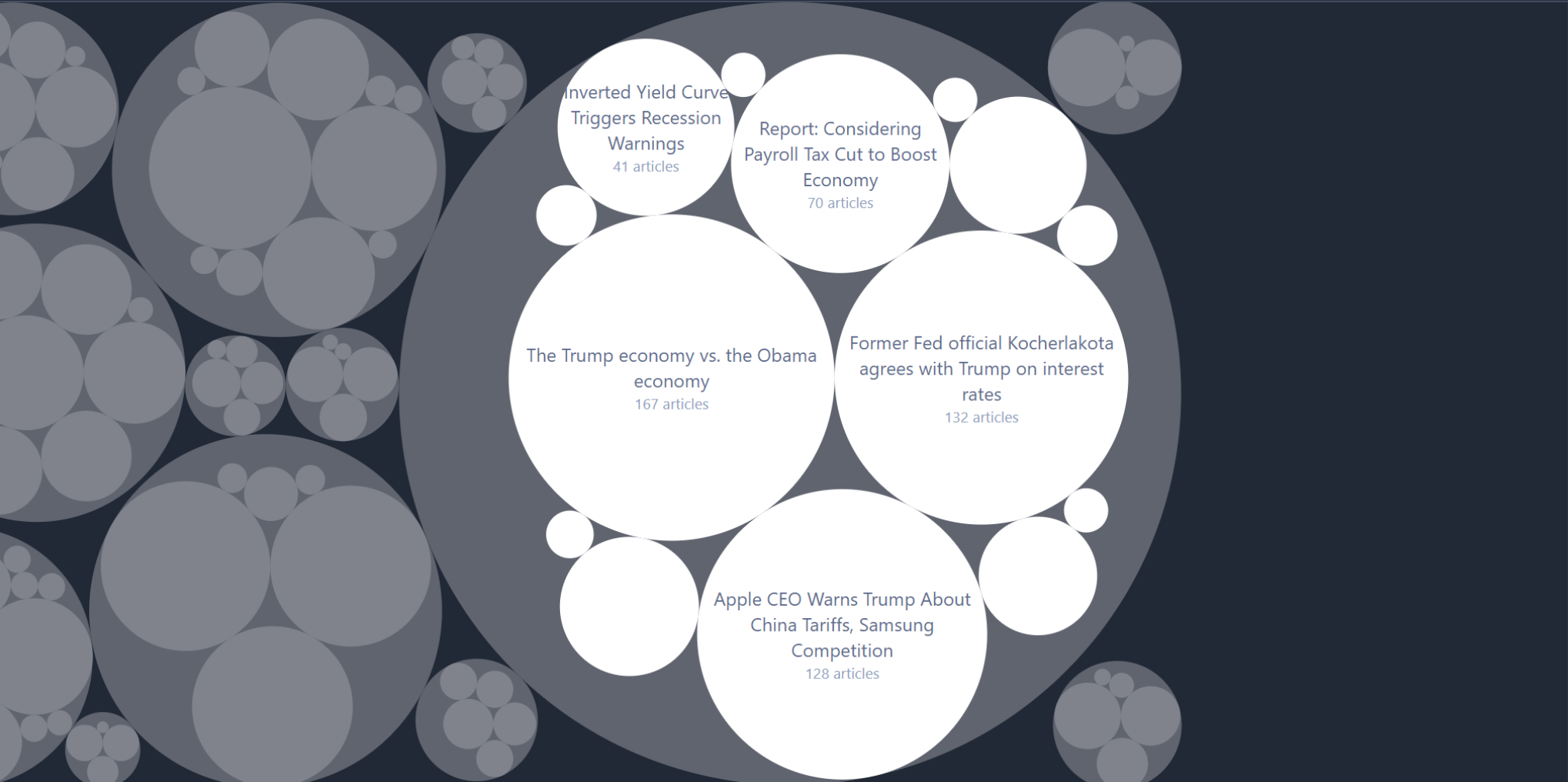
- 614 Trump dismisses worries of recession, s...
- 277 Twitter, Facebook suspend China's acco...
- 262 Biden regains double-digit lead in new el...
- 234 Clemson is No. 1 in AP preseason footb...
- 205 Prince Andrew 'appalled' by Jeffrey Epst...
- 180 Omar, Tlaib to discuss Israel, Palestine t...
- 177 NYPD fires Officer Daniel Pantaleo over ...
- 167 UK Government to End EU Open Border...
- 151 Planned Parenthood leaves family plann...
- 139 Elizabeth Warren apologizes for heritage...
- 105 49ers vs. Broncos : NFL preseason game
- 101 Netflix loses \$300 handle on Apple TV+ ...
- 96 William Barr appoints new prisons direct...
- 89 Antonio Brown's helmet saga continues t...
- 86 El Paso Mass Shooting Suspect On Sui...
- 84 Iran tanker row: Detained ship sets sail f...

# production dashboard

Project Ratio [About this project](#)



▶ August 19, 2019    Top 25 News Events    > Trump dismisses worries of recession, says economy is strong  
3572 articles    EVENT – 614 articles



← Expand Sidebar

|     |  |
|-----|--|
| 614 | Trump dismisses worries of recession, s...   |
| 167 | The Trump economy vs. the Obama e...         |
| 132 | Former Fed official Kocherlakota agre...     |
| 128 | Apple CEO Warns Trump About Chin...          |
| 70  | Report: Considering Payroll Tax Cut to...    |
| 41  | Inverted Yield Curve Triggers Recessi...     |
| 25  | Compensation Growth Supporting Str...        |
| 22  | Georgia voters challenge legality of ne...   |
| 12  | China could crush U.S. military in Pacific   |
| 4   | Jon Voight declares Trump is 'greatest...    |
| 4   | The Next Recession Will Hurt Some I...       |
| 3   |  |
| 2   | Consumer Reports: How to protect yo...       |
| 2   | Life expectancy drops in Wisconsin du...     |
| 2   | Newsletter: The RV Industry Signals a...     |
| 277 | Twitter, Facebook suspend China's acco...    |
| 262 | Biden regains double-digit lead in new el... |
| 234 | Clemson is No. 1 in AP preseason footb...    |

# consumption data



**Nationally broadcast TV content.** Nielsen's nationally representative TV panel ( $N = 100,000$ ), from which we compute the number of minutes per person per day devoted to watching specific programs.



**Local TV news.** Subset of the national panel ( $N = 50,000$ ) sampled from the 25 largest local markets.



**Online content (desktop).** Nielsen's nationally representative web panel ( $N = 65,000$ ), which records the complete browser history of its panelists including the specific URLs visited as well as the length of the visit.



**Online content (mobile).** Comscore aggregated metrics, which track total time spent and total online viewers per month for each website for mobile (including browser and app usage) and desktop respectively by age group.

# consumption pipeline

## Television news consumption

- Time consuming any of the roughly 400 programs that are classified by Nielsen as “news”
- Upper Bound: Magazine news (e.g. Inside Edition, Dateline), morning shows (e.g. Good Morning America, Today Show), entertainment news (e.g. TMZ, Access Hollywood), and late-night shows (e.g. The Daily Show with Trevor Noah, the Late Show with Stephen Colbert)

## Online news consumption

- Time consuming any article published in more than 800 websites, adapted from (Athey et al 2017): primarily cover “hard” news topics like politics, business, and US and international affairs.
- Fraction of time consuming social media sites (Facebook, Twitter, Reddit), search engines (Google, Bing, and Yahoo) using referral links
- Fraction of time consuming YouTube (classified 10,000 videos internally as “news and politics”)
- 100% of time consuming portals (MSN, Yahoo, AOL)

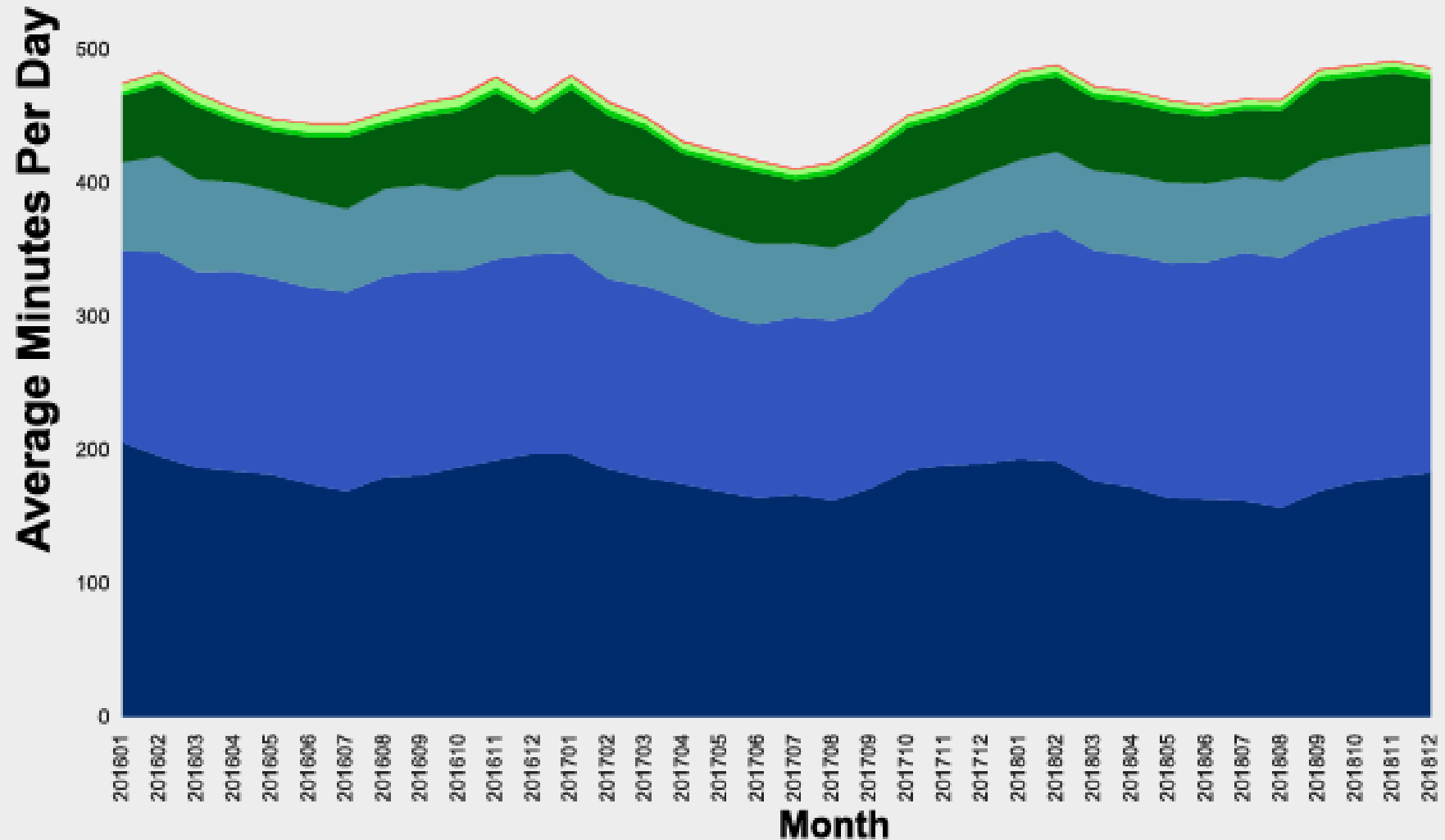
## Fake news consumption

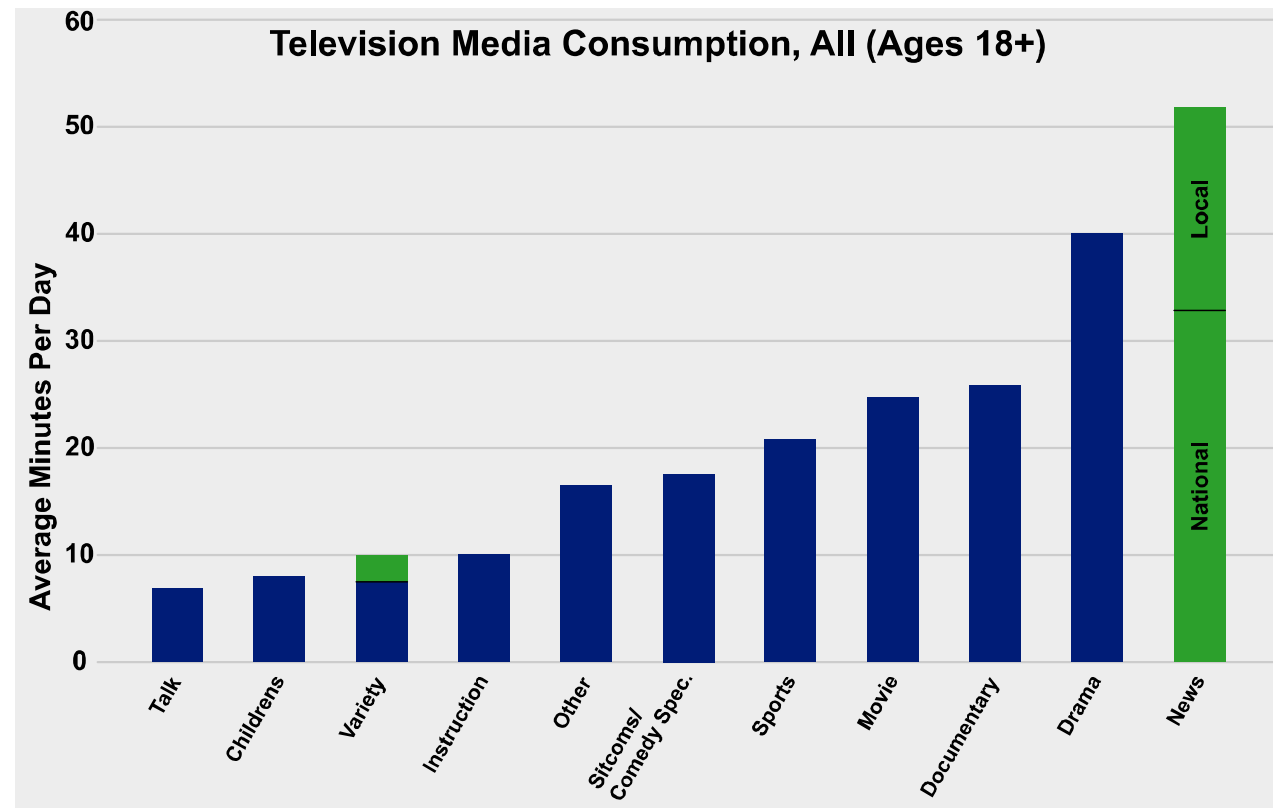
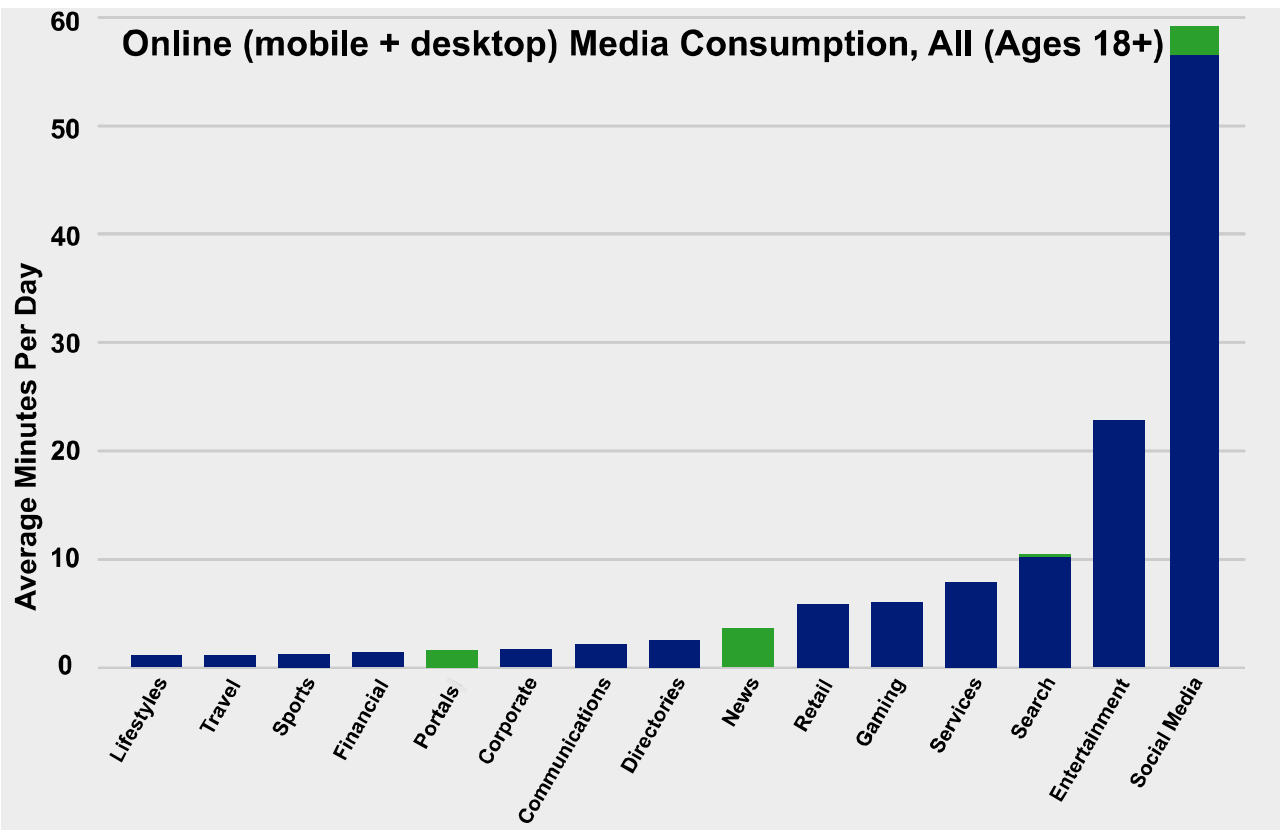
- Time consuming 1 of 98 websites previously identified by researchers (Grinberg et al 2017), professional fact checkers, journalists as: fake, deceptive, low quality, or hyperpartisan news.
- Also attribute some fraction to social media and search as above
- YouTube hand classified

Both news and fake news definitions intended to overestimate real consumption

TV: Not News    TV: News    N/A TV: Fake News  
Mobile: Not News    Mobile: News    Mobile: Fake News  
Desktop: Not News    Desktop: News    Desktop: Fake News

## Media Consumption, All (Ages 18+)





# consumption: platforms, social media, search

|   | MSN<br>Portal | Yahoo!<br>Portal | AOL<br>Portal | Bing<br>Search | Google<br>Search | Yahoo!<br>Search | You-<br>Tube | Face-<br>book | Reddit | Twitter |
|---|---------------|------------------|---------------|----------------|------------------|------------------|--------------|---------------|--------|---------|
| % of overall online<br>Media Consumption<br>that is on this<br>Platform | 0.2%          | 0.4%             | 0.1%          | 0.2%           | 2.1%             | 0.1%             | 9.4%         | 10.6%         | 0.2%   | 0.5%    |
| % of Consumption<br>on Platform that is<br>News                         | 100%          | 100%             | 100%          | 4%             | 5%               | 3%               | 3%           | 8%            | 5%     | 11%     |
| Minutes of<br>News/Day on this<br>Platform                              | 0.52          | 0.84             | 0.16          | 0.02           | 0.22             | 0.01             | 0.60         | 1.93          | 0.03   | 0.12    |
| % of News on<br>Platform that is<br>Swampy                              | 0%            | 0%               | 0%            | 4%             | 4%               | 7%               | 38%          | 15%           | 6%     | 6%      |
| Minutes of Real<br>News/Day on this<br>Platform                         | 0.52          | 0.84             | 0.16          | 0.02           | 0.21             | 0.01             | 0.37         | 1.63          | 0.02   | 0.11    |

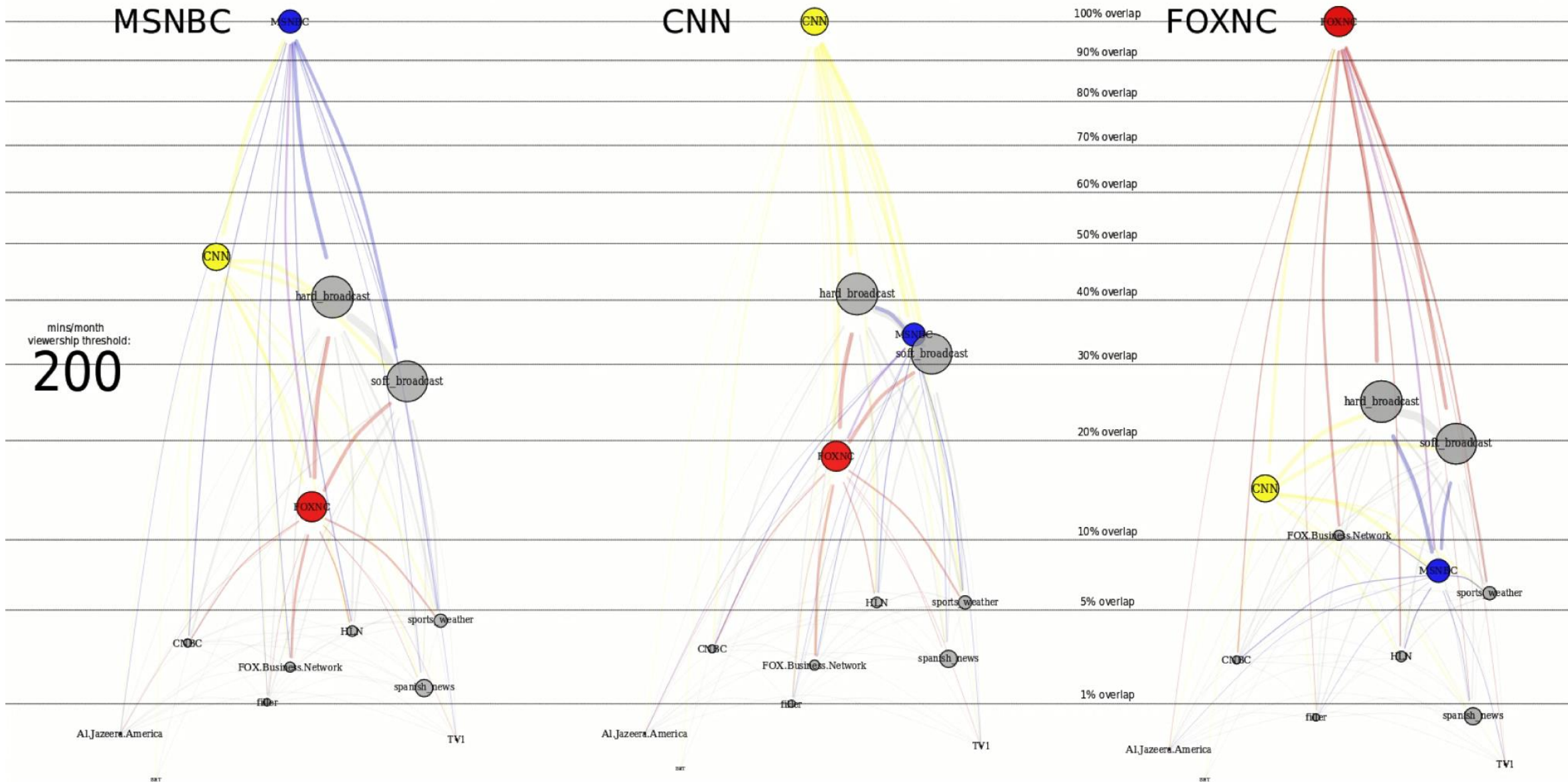
# consumption: networks

Audience overlap with MSNBC, CNN, FOXNC, with broadcast, threshold changing from 0 to 300 mins/month

4% of all viewers

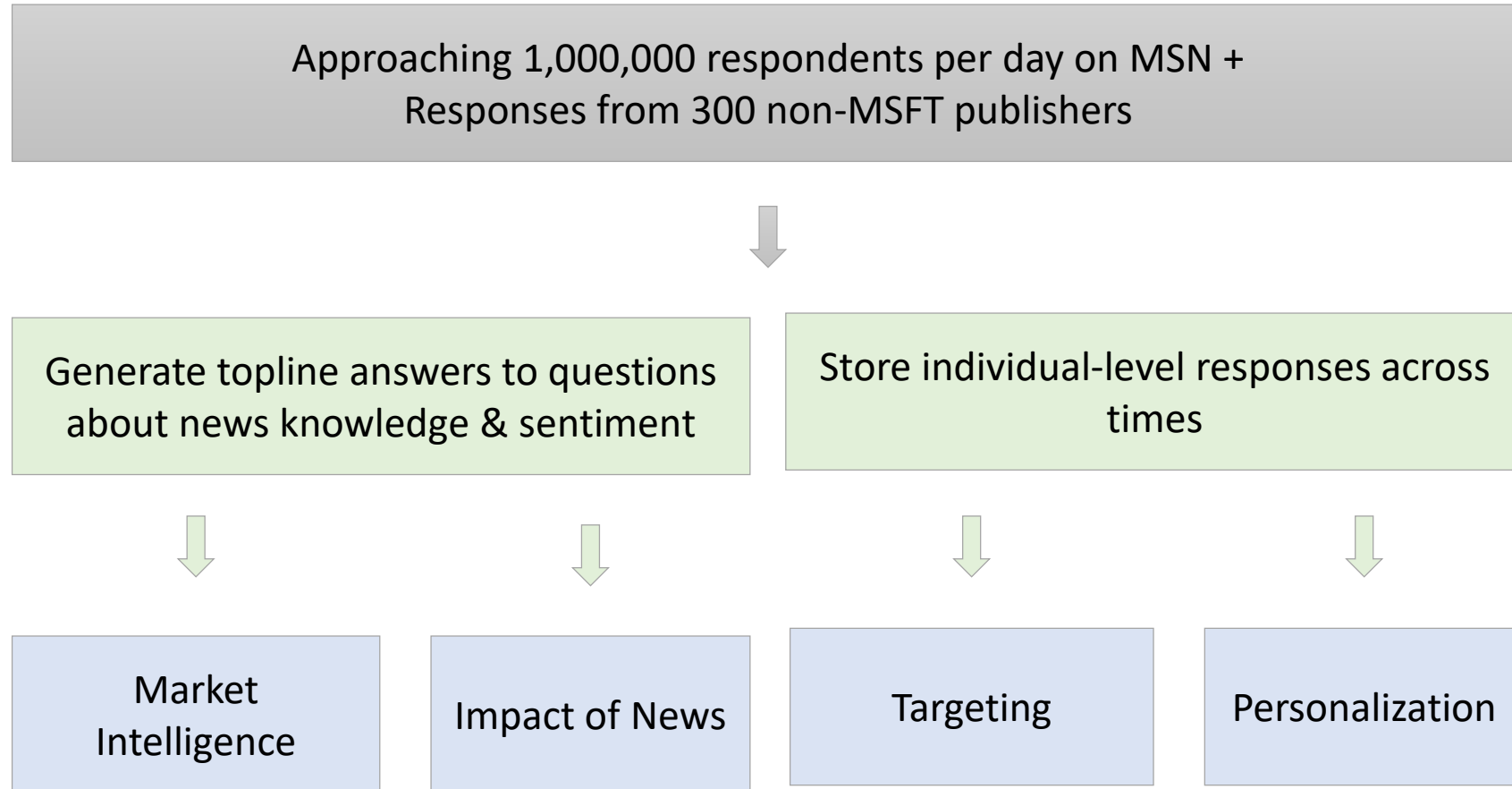
6% of all viewers

7% of all viewers





# absorption data



# absorption data

