



Fighting Disinformation Using NLP

Jeff Kao | Minneanalytics Big Data Conference | 6/5/2018



Jeff Kao

Researcher,  **DATA FOR
DEMOCRACY**

Machine Learning Engineer,  atrium

https://github.com/j2kao/fcc_nn_research/



METIS



COLUMBIA
LAW SCHOOL

UNIVERSITY OF
WATERLOO





Astroturf (n.)

An artificially-manufactured political movement designed to give the appearance of grassroots activism.*

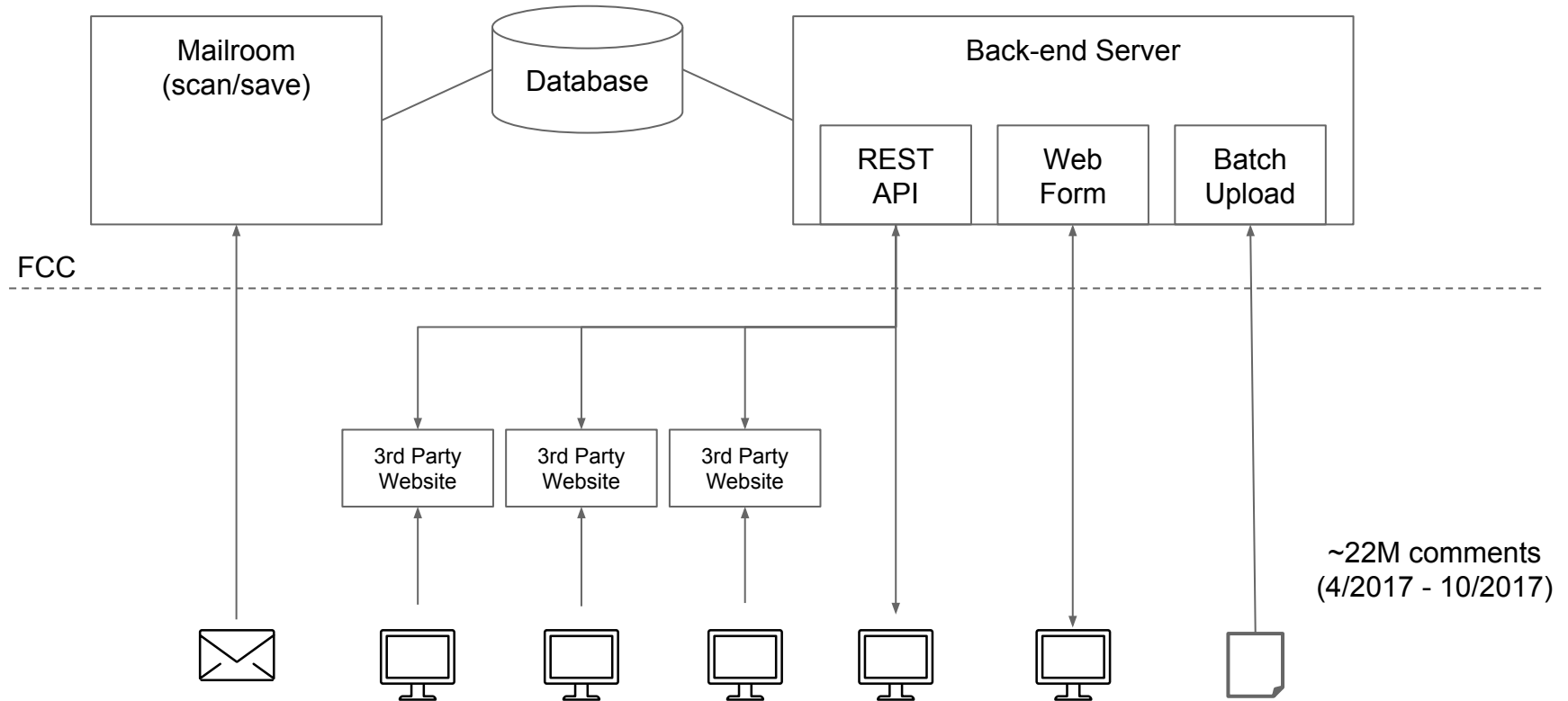


1. Background
2. Data Analysis (+ short jupyter notebook session)
3. Survey/Validation
4. Discussion/Applications



1. Background

FCC Public Comment System



ECFS Express

1
COMMENT

2
REVIEW

3
CONFIRMATION

Proceeding(s)

* Required - Type to search proceedings. Press ENTER key after each selection

Name(s) of Filer(s)

*Required - Add Filer(s). Press ENTER key after each entry

**Primary Contact
Email**

example@email.com

Address

☐ International

* Required

Address 2

City

* Required

State

* Required ▼

ZIP

* Required

Brief Comments

Note: You are filing a document into an official FCC proceeding. All information submitted, including names and addresses, will be publicly available via the web.

☐ **Email Confirmation**

Continue to review screen ➞

Reset Form

Filing Detail

ID	1051157755251	Proceeding	17-108
Name of Filer	<u>Barack Obama</u>		
Type of Filing	COMMENT	Filing Status	DISSEMINATED
Viewing Status	Unrestricted		
Date Received	May 11, 2017	Date Posted	May 12, 2017
Address	1600 Pennsylvania Ave NW	City	Washington
		State	DC
ZIP	20500		
Brief Comment	The unprecedented regulatory power the Obama Administration imposed on the internet is smothering innovation, damaging the American economy and obstructing job creation. I urge the Federal Communications Commission to end the bureaucratic regulatory overreach of the internet known as Title II and restore the bipartisan light-touch regulatory consensus that enabled the internet to flourish for more than 20 years. <u>The plan currently under consideration at the FCC to repeal Obama's Title II power grab is a positive step forward and will help to promote a truly free and open internet for everyone.</u>		



22M comments
(4/2017- 10/27/2017)

3M unique comments

3M 300-D vectors

scrape to mongoDB

migrate important
fields to postgresql

hash text; remove
duplicates;
preserve count

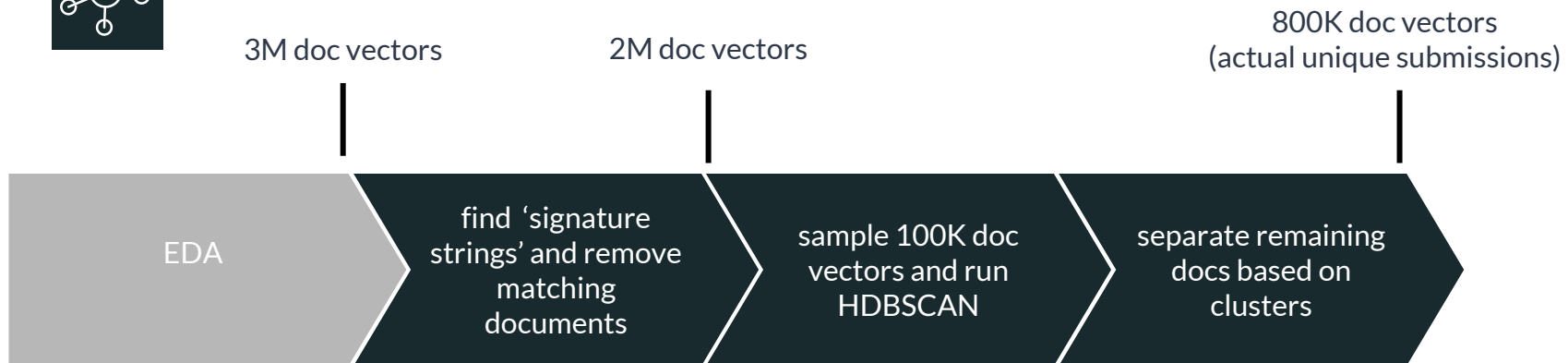
encode 300-D
vector for each
unique comment

ECFS API
(JSON)

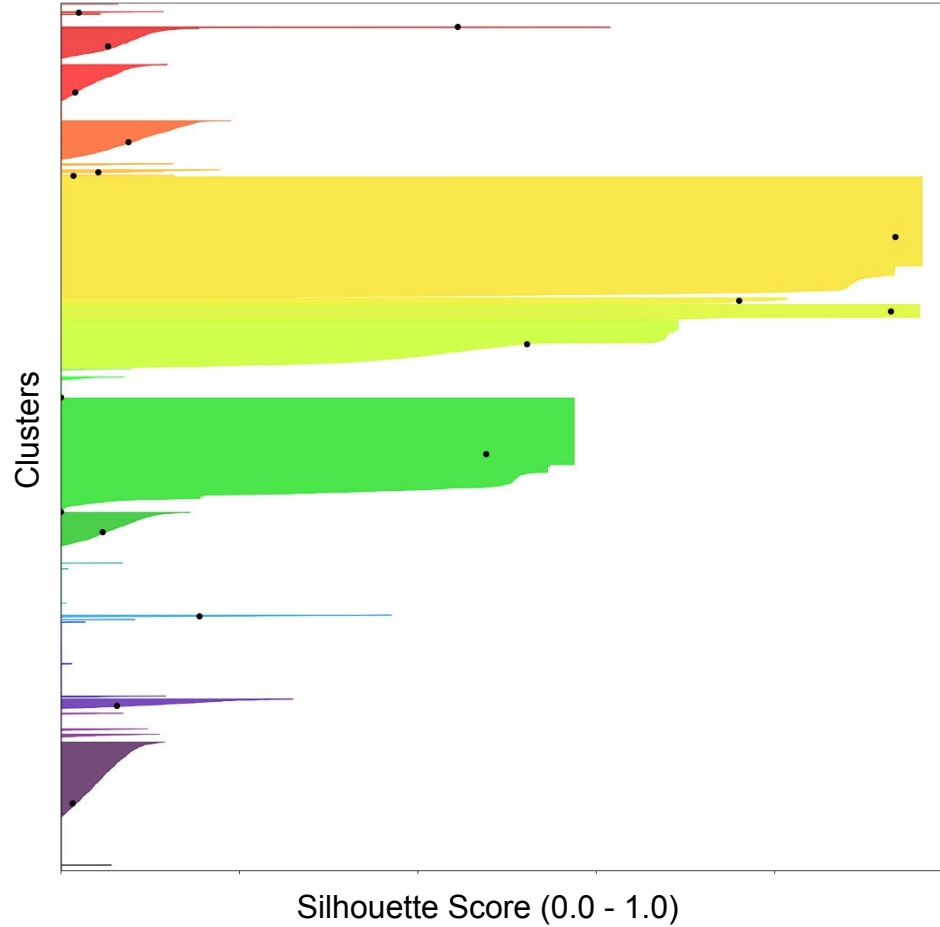


2. Data Analysis

https://github.com/j2kao/fcc_nn_research/



Silhouette Plot (EDA)





2. Data Analysis:

“Mad-lib” Bot Comments

"In the matter of restoring Internet freedom. I'd like to recommend the commission to undo The Obama/Wheeler power grab to control Internet access. Americans, as opposed to Washington bureaucrats, deserve to enjoy the services they desire. The Obama/Wheeler power grab to control Internet access is a distortion of the open Internet. It ended a hands-off policy that worked exceptionally successfully for many years with bipartisan support.",

"Chairman Pai: With respect to Title 2 and net neutrality. I want to encourage the FCC to rescind Barack Obama's scheme to take over Internet access. Individual citizens, as opposed to Washington bureaucrats, should be able to select whichever services they desire. Barack Obama's scheme to take over Internet access is a corruption of net neutrality. It ended a free-market approach that performed remarkably smoothly for many years with bipartisan consensus.",

"FCC: My comments re: net neutrality regulations. I want to suggest the commission to overturn Obama's plan to take over the Internet. People like me, as opposed to so-called experts, should be free to buy whatever products they choose. Obama's plan to take over the Internet is a corruption of net neutrality. It broke a pro-consumer system that performed fabulously successfully for two decades with Republican and Democrat support.",

"Mr Pai: I'm very worried about restoring Internet freedom. I'd like to ask the FCC to overturn The Obama/Wheeler policy to regulate the Internet. Citizens, rather than the FCC, deserve to use whichever services we prefer. The Obama/Wheeler policy to regulate the Internet is a perversion of the open Internet. It disrupted a market-based approach that functioned very, very smoothly for decades with Republican and Democrat consensus.",

"FCC: In reference to net neutrality. I would like to suggest Chairman Pai to reverse Obama's scheme to control the web. Citizens, as opposed to Washington bureaucrats, should be empowered to buy whatever products they prefer. Obama's scheme to control the web is a betrayal of the open Internet. It undid a hands-off approach that functioned very, very successfully for decades with broad



2. Data Analysis:

Clustering Comment Campaigns



3M doc vectors

2M doc vectors

800K doc vectors
(actual unique submissions)



HDBSCAN:

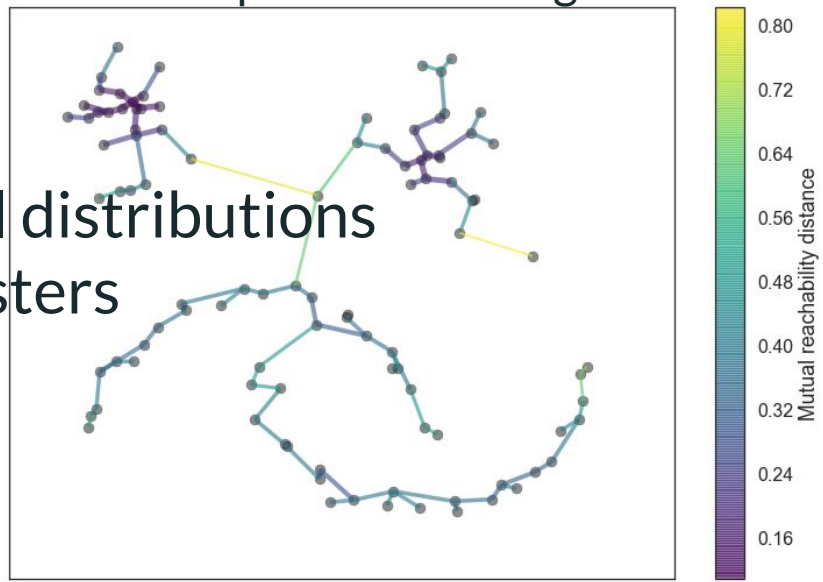
- hierarchical clustering based on DBSCAN
 - single linkage; build minimum spanning tree based on mutual reachability distance
 - hierarchically cluster by condensing the tree
 - progressively increase distance threshold to capture clusters greater than `minimum_cluster_size`

Advantages:

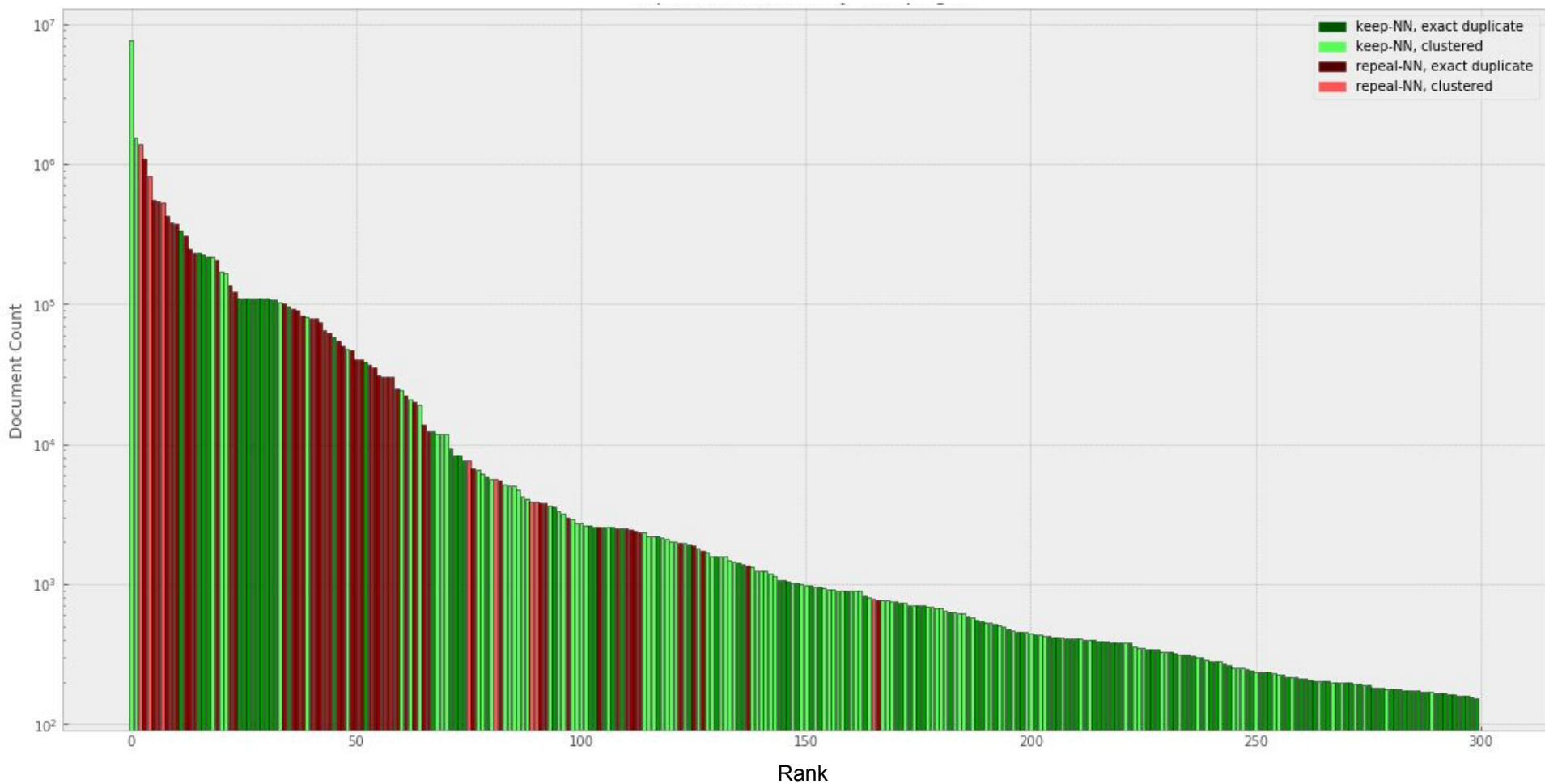
- deals well with outliers, long tail distributions
- arbitrary shape / number of clusters

Disadvantages:

- cluster drift



Top 300 Net Neutrality 'Campaigns'





2. Data Analysis:

Analyzing Unique Comments



3M doc vectors

2M doc vectors

800K doc vectors
(actual unique submissions)



What proportion of “unique” outlier comments were pro-net neutrality?

Analysis	% Keep-NN	Method
Fossett (5/2017) **	97%	Sampling + Hand Labeling
Emprata (8/2017) *	98.7%	Supervised learning (LSTM Classifiers)
Kao (11/2017)	99.6%	Sampling + Hand Labeling

* <http://www.emprata.com/emp2017/wp-content/uploads/2017/08/FCC-Restoring-Internet-Freedom-Comments-Analysis.pdf>

** <http://jeffreyfossett.com/2017/05/13/fcc-filings.html>

What proportion of “unique” outlier comments were pro-net neutrality?

“I am 82, handicapped, and home bound, but not lonely, because I have the free internet. I can roam the world. use Facebook to visit family friends. I can sell my work on Etsy without fear of Amazon getting preference should the 2015 law be repealed. If you (The FCC) no longer had oversight, my ISP could raise its prices so that I couldn’t afford to have the Internet at all! I am relying on the FCC to protect me and others like me.”



3. Survey/Validation

- More Direct Measurement
- 40+ largest campaigns identified
- 450,000 emails sent
- 40% bounce rate overall
- 3-4% response rate overall

STARTUP
POLICY LAB

Address:
[1355 Market Street](#)
[San Francisco, California](#)

Hello,

You are receiving this email because, according to the public records of the Federal Communications Commission (FCC), you submitted comments in FCC Proceeding 17-108 regarding net neutrality on 08/11/2017.

We are a team of public interest researchers from Startup Policy Lab conducting a one-time survey to verify FCC submissions. Your participation in this one-question survey will help confirm your submission.

To support our research project, would you please indicate whether you submitted the following comment:

This is the comment that was attributed to Jeffery Kao with email address at jeffery.kao@gmail.com:

"In 2015 Chairman Tom Wheeler's Federal Communications Commission (FCC) imposed restrictive Title II utility-style regulations u..."

*Note: Only the first 255 characters of the comment are included here.

**Can you please confirm that you
submitted this comment?**

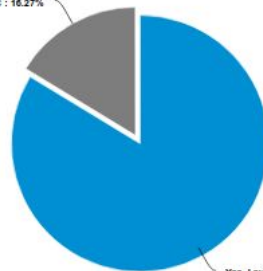
Federal Communications Commission (FCC) - Das...



Can you please confirm that you submitted this comment?



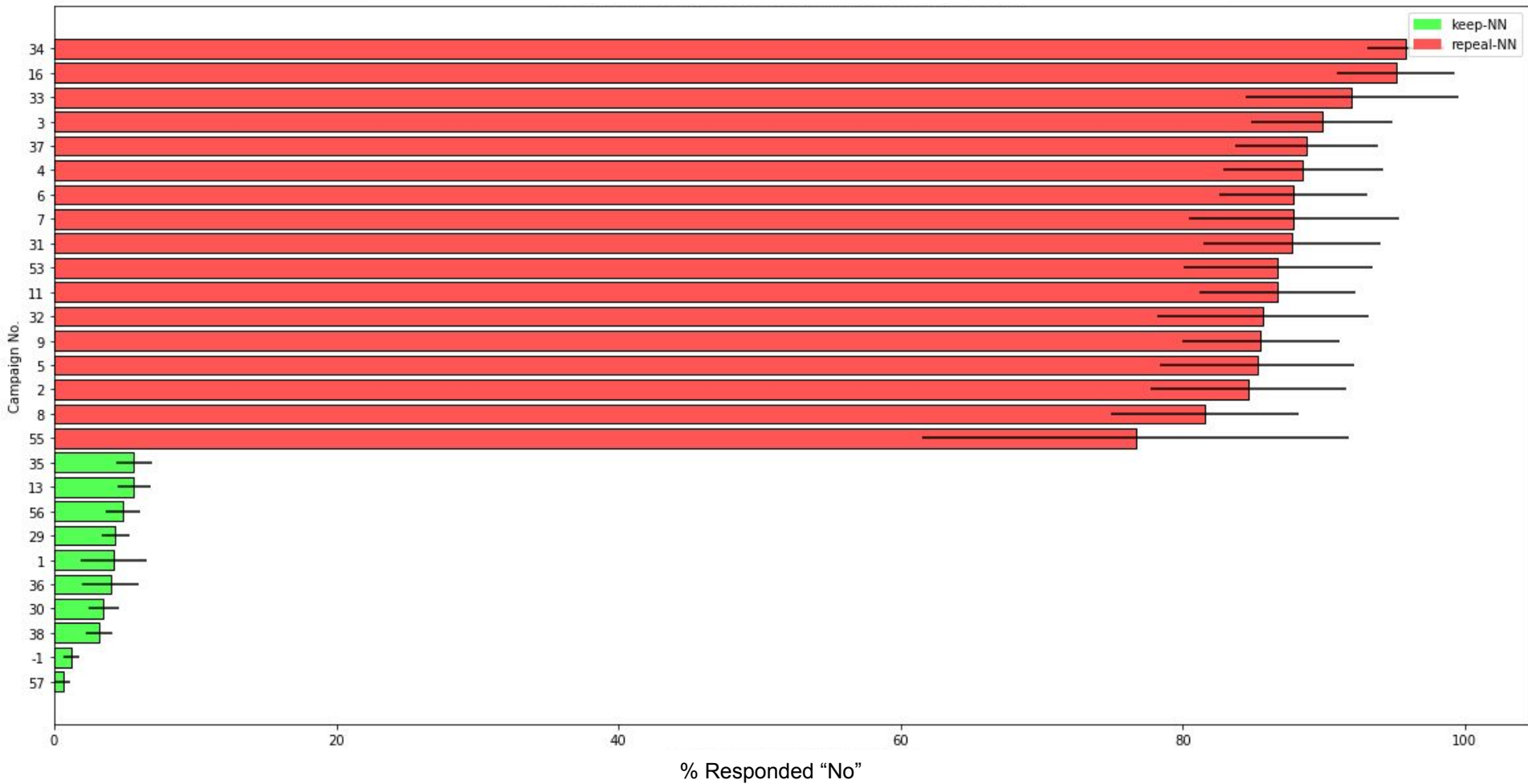
No, I did not submit the comment above to the FCC : 16.27%



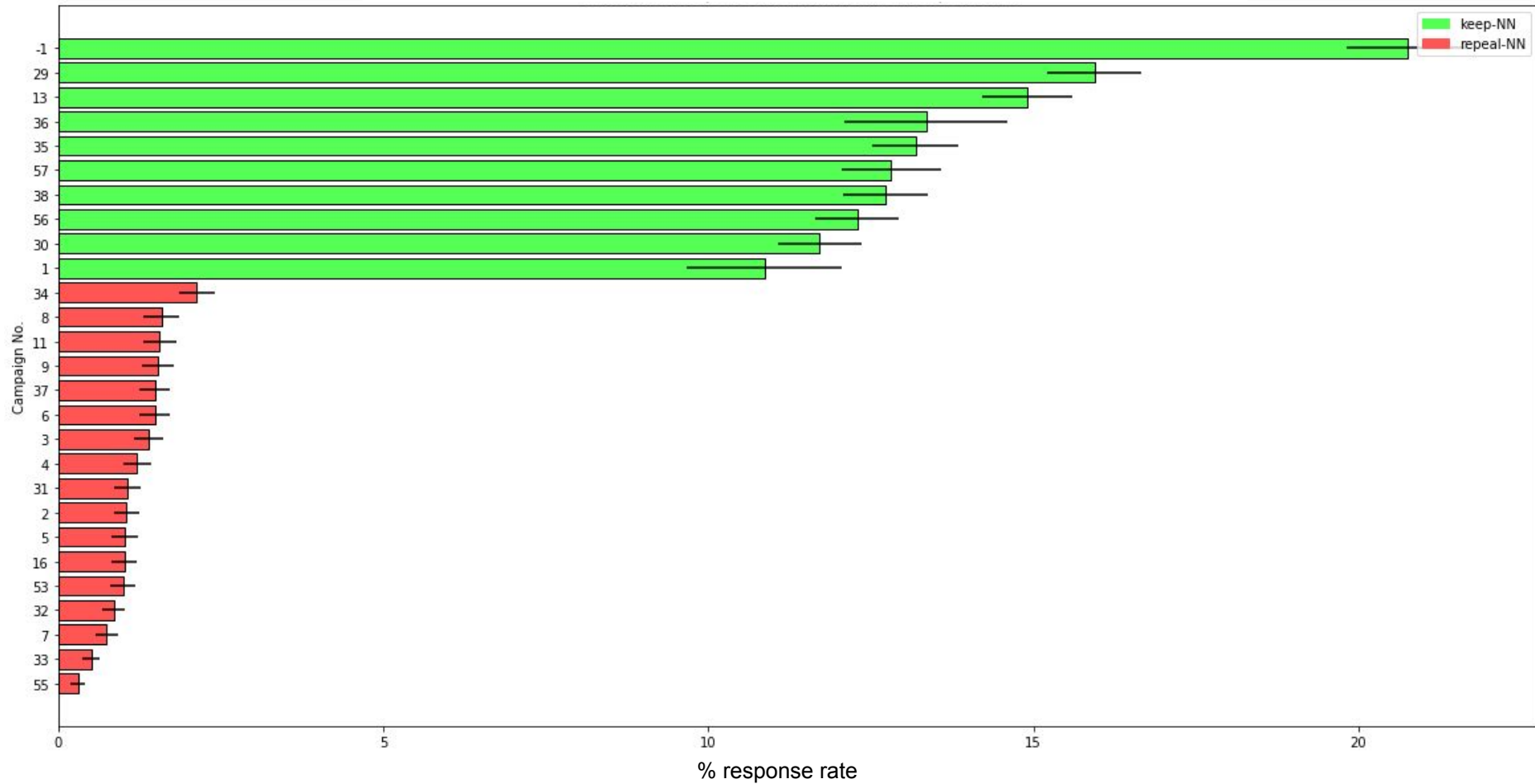
Yes, I submitted the comment above to the FCC : 83.73%

Answer	Count	Percent	20%	40%	60%	80%	100%
Yes, I submitted the comment above to the FCC	11305	83.73%	<div></div>				
No, I did not submit the comment above to the FCC	2197	16.27%	<div></div>				
Total	13502	100 %					

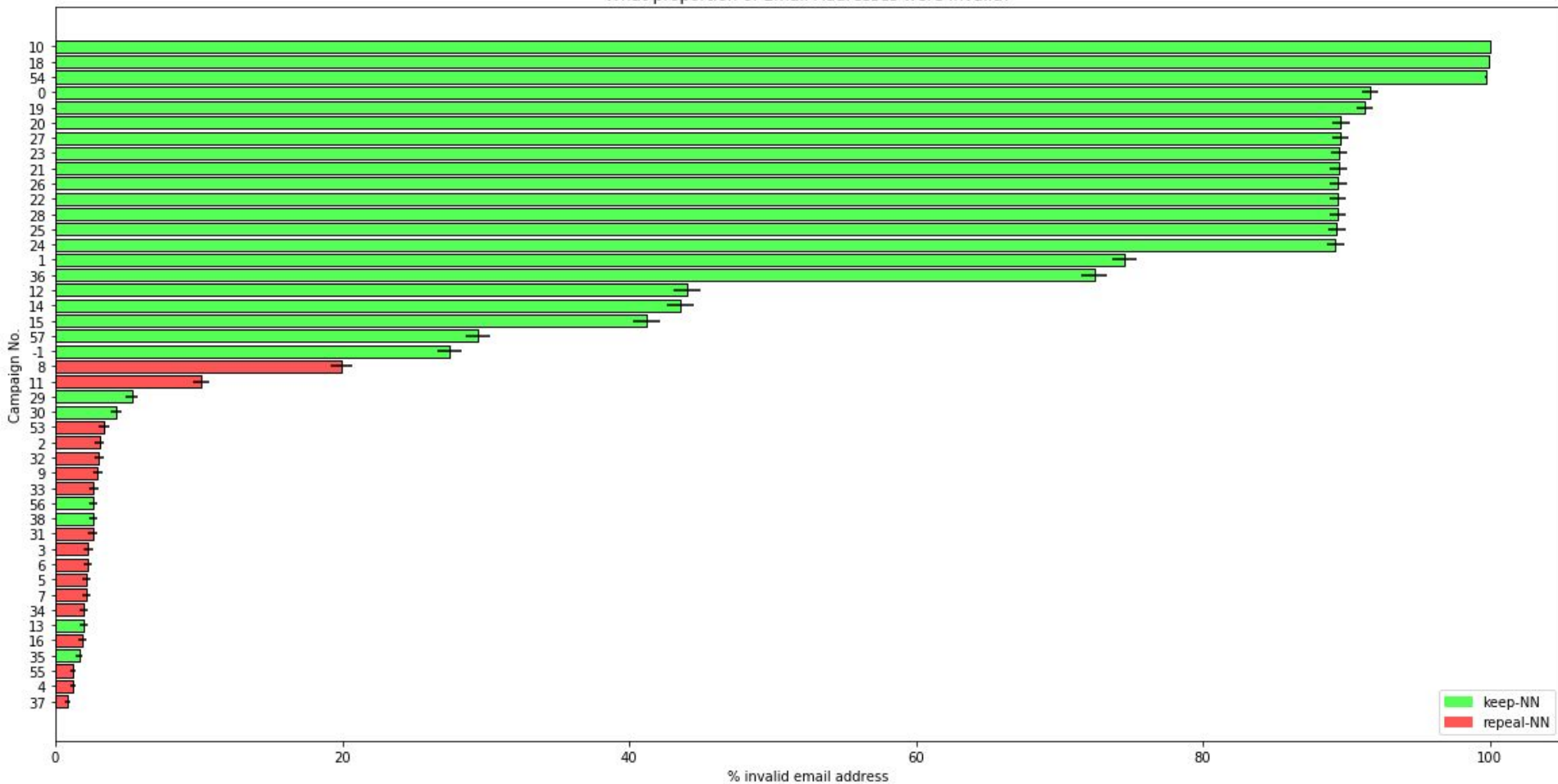
Did you submit this comment? (min. 20 responses)



Commenter Response Rate (min. 20 responses)



What proportion of Email Addresses were Invalid?

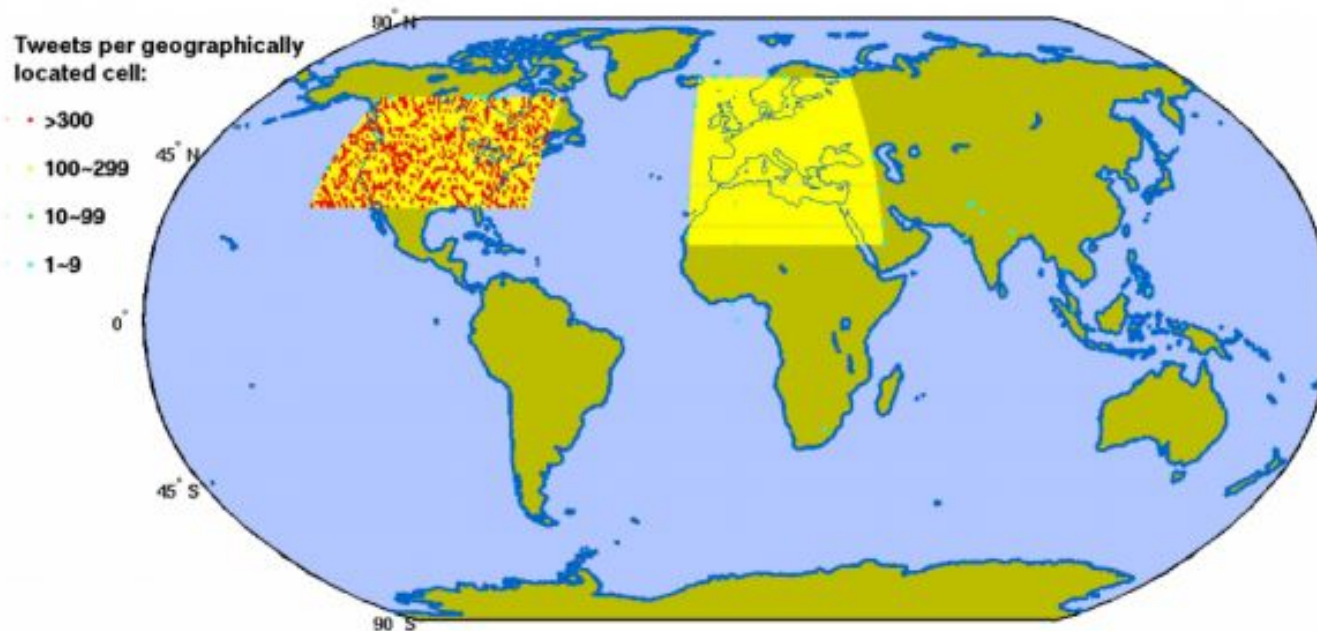




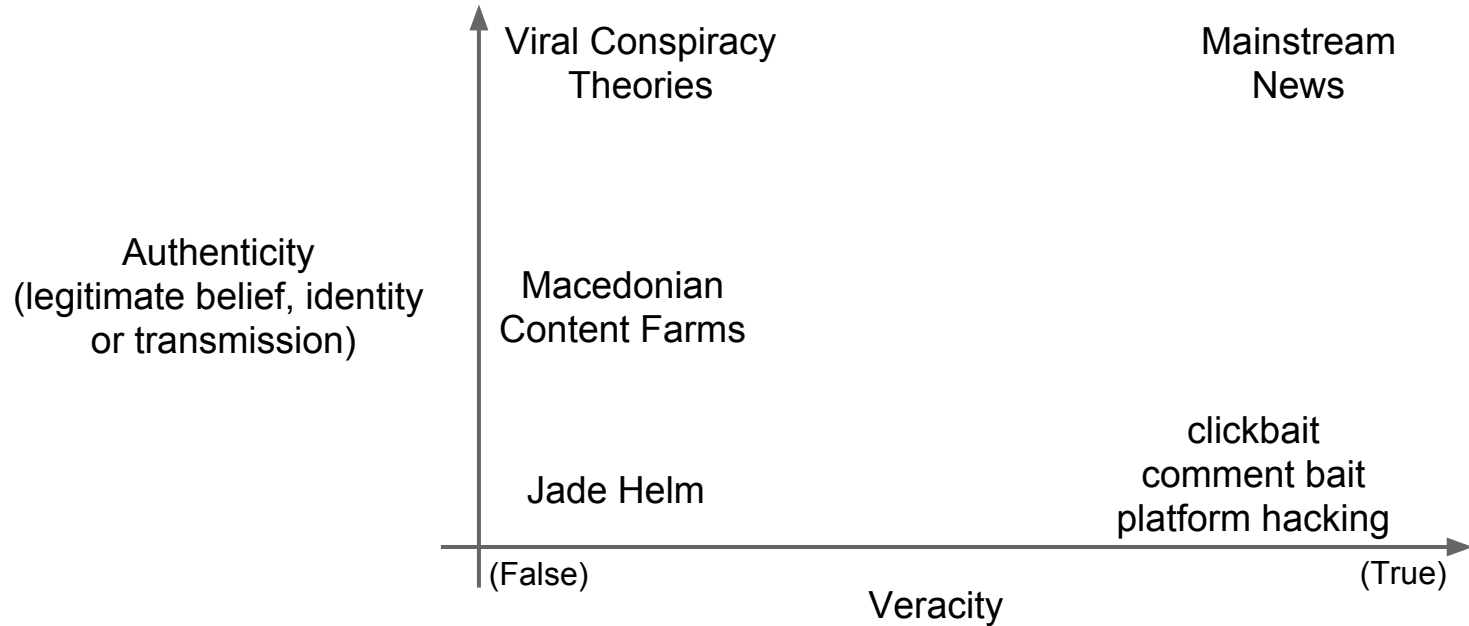
4. Discussion/Applications

- Vector embeddings can help detect astroturf not caught by bag-of-words and string matching techniques
- Distributions over semantic meaning can be used as a feature for determining authenticity
- While unsupervised learning can help, analysis is highly context dependent (public availability of data for researchers)
- Machine learning is not a silver bullet: what are you actually fitting to?
(<https://www.theverge.com/2018/5/22/17380630/twitter-moderation-cyrillic-russian-bots>)

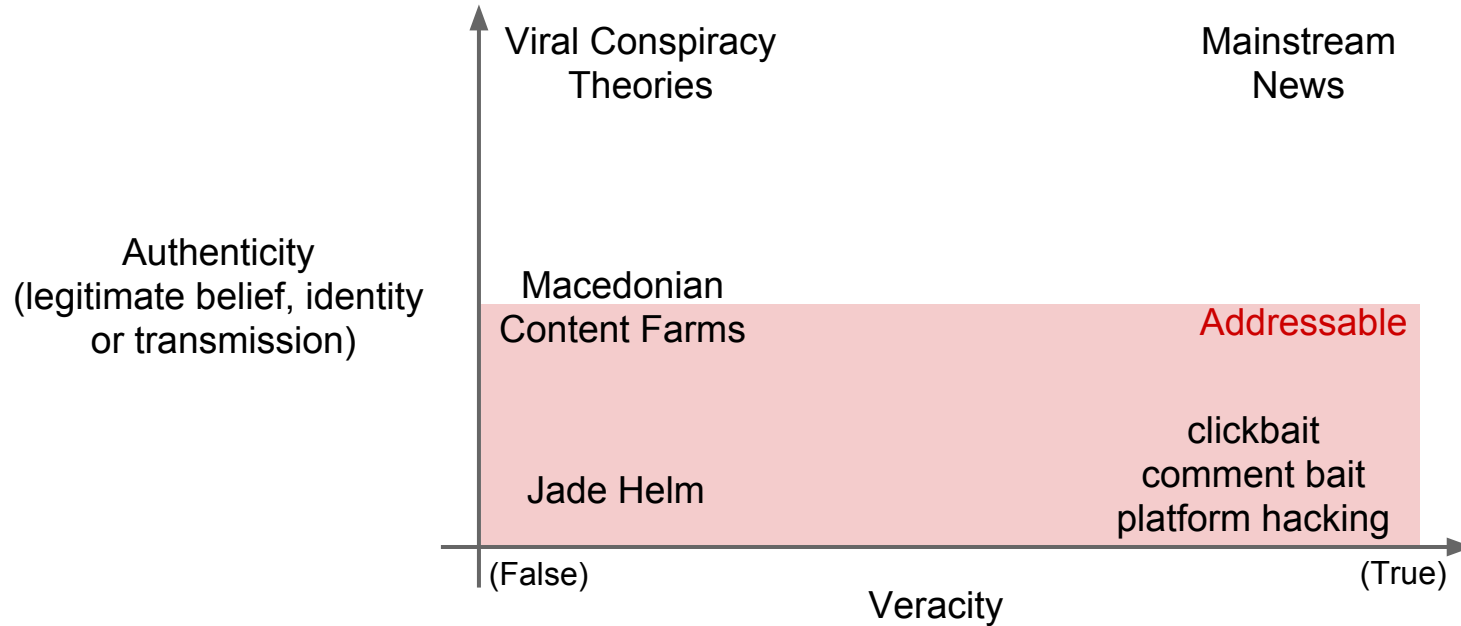
Example: Star Wars Twitter botnets



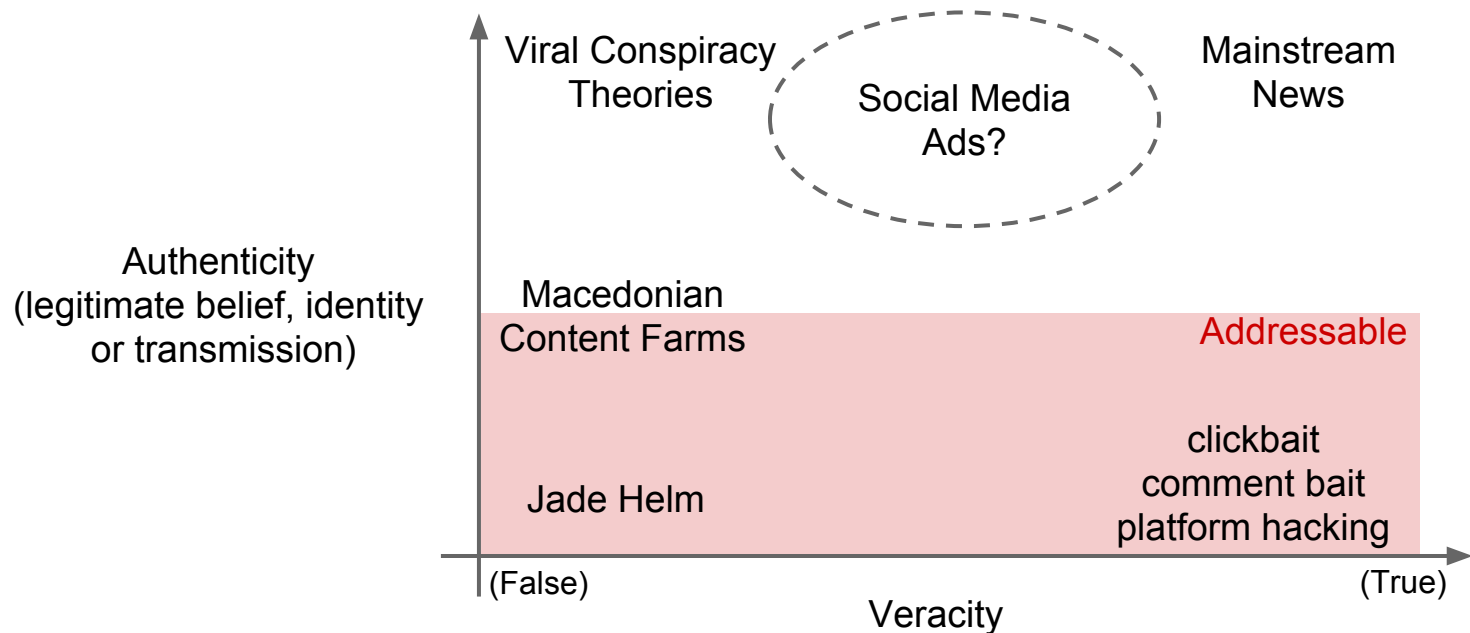
What do we mean by “Fake News”?



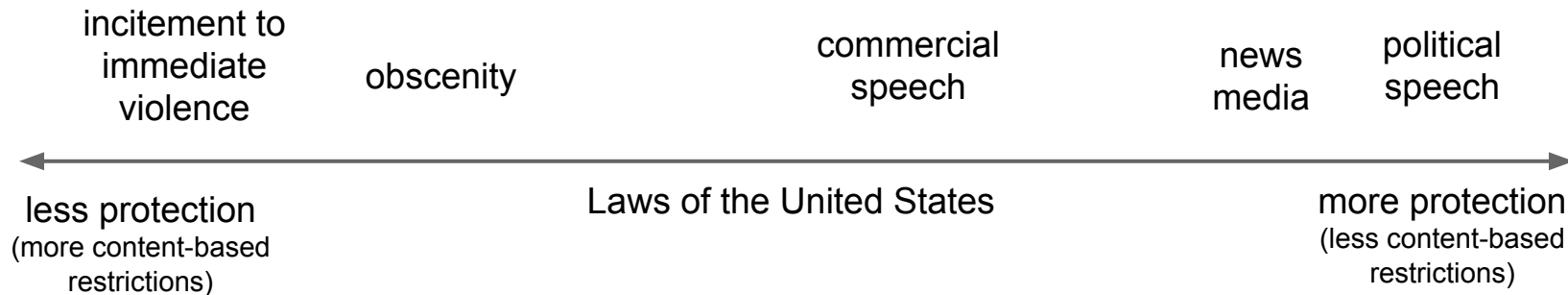
What do we mean by “Fake News”?



What do we mean by “Fake News”?



Can we expect to regulate “Fake News”?



Thanks!



jeffykao.com



jeffery.kao@gmail.com



linkedin.com/in/jeffykao



@jeffykao



j2kao



