# Hunting Malicious Bots on Twitter: An Unsupervised Approach

## Zhouhan Chen

Advisor: Devika Subramanian,
Committee Members: Dan Wallach and Lin Zhong

April 6, 2018

Computer Science Department

RICE®

MIND    HEALTH    TECH    SUSTAINABILITY    EDUCATION    VIDEO    PODCASTS    B

COMPUTING

# How Twitter Bots Help Fuel Political Feuds

# Agenda

- Problem: Identifying Twitter bots and spammers who create those bots
    a. Significance of problem
    b. Existing approaches
    c. Shortcomings of existing approaches

- Contribution: Designed and implemented a **group based unsupervised** algorithm that effectively detects bots and **spam campaigns**

- Results and findings

- Use Cases
    1. Hong Kong #UmbrellaRevolution
    2. #ReleaseTheMemo

- Future work

# Why focus on Twitter bots

## Impersonation: 2016 US Election tweet collection



**Link to malware**

# Severity of the problem

## Bots are a major presence on Twitter

- 9-15% of Twitter accounts are bots[1]
- 50% of tweet traffic generated by bots[2]

## Bots violate Twitter's terms of service

- Send spam (click-bait, affiliate marketing)[3]
- Send malware
- Interfere with elections[4]

[1] O. Varol, E. Ferrara, C. A. Davis, F. Menczer, and A. Flammini, "Online human-bot interactions: Detection, estimation, and characterization," 2017.

[2] Z. Gilani, J. Crowcroft, R. Farahbakhsh, and G. Tyson, "The implications of twitterbot generated data traffic on networked systems," in Proceedings of the SIGCOMM Posters and Demos, ser. SIGCOMM Posters and Demos '17. New York, NY, USA. 2017

[3] Bessi, A., & Ferrara, E. (2016). Social bots distort the 2016 U.S. Presidential election online discussion. First Monday. 2016.

[4] K. Thomas, C. Grier, D. Song, and V. Paxson, "Suspended Accounts in Retrospect: An Analysis of Twitter Spam," in Proceedings of the 2011 ACM SIGCOMM Conference on Internet Measurement Conference, IMC '11, (New York, NY, USA), pp. 243–258, ACM, 2011.

# Existing bot detection methods

## Two major approaches

- **Supervised approaches** that learn to classify bots based on a number of structural and behavioral features of bots.
- **Unsupervised approaches** that use a programmed protocol based on pre-defined behavioral features.

## Features for detection

- Behavioral: temporal tweeting patterns[1][2]
- Structural: number of tweets[1], shortened URL usage[3]

......

[1] Davis, C. A., Varol, O., Ferrara, E., Flammini, A., & Menczer, F. (2016, April). BotOrNot: A system to evaluate social bots. In Proceedings of the 25th International Conference Companion on World Wide Web (pp. 273-274). International World Wide Web Conferences Steering Committee Available: https://arxiv.org/abs/1602.00975

[2] Z. Chu, S. Gianvecchio, H. Wang, and S. Jajodia, "Detecting automation of twitter accounts: Are you a human, bot, or cyborg?" IEEE Transactions on Dependable and Secure Computing, vol. 9, no. 6, pp. 811–824, Nov 2012.

[3] D. Wang, S. B. Navathe, L. Liu, D. Irani, A. Tamersoy, and C. Pu, "Click traffic analysis of short url spam on twitter," in 9th IEEE International Conference on Collaborative Computing: Networking, Applications and Worksharing, Oct 2013, pp. 250–259.
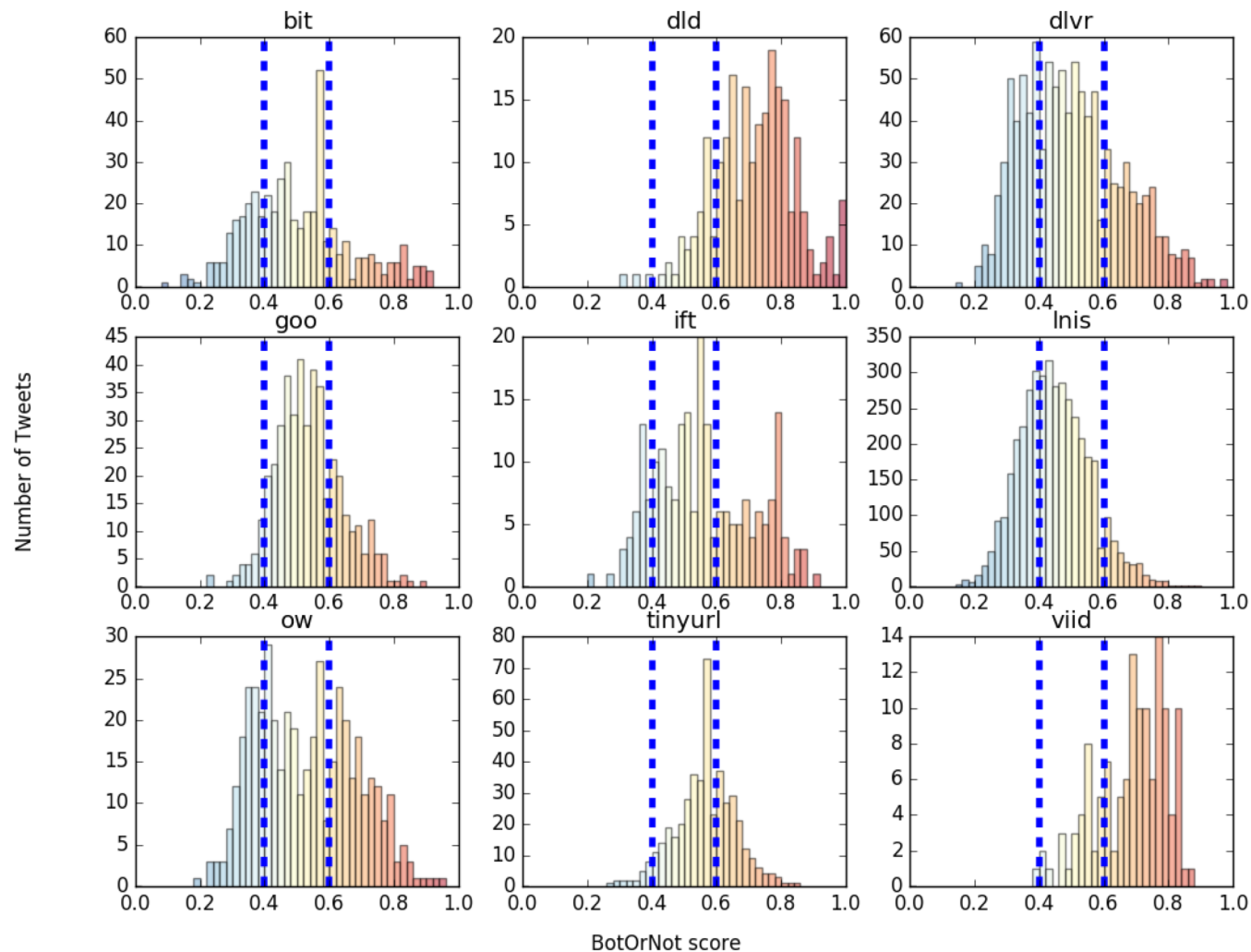
# Existing bot detection methods

| | Supervised |
|---|---|
| Human intervention | Yes |
| Unit of detection | Individual |
| State-of-art Application | BotOrNot[1] |

[1] Davis, C. A., Varol, O., Ferrara, E., Flammini, A., & Menczer, F. (2016, April). BotOrNot: A system to evaluate social bots. In Proceedings of the 25th International Conference Companion on World Wide Web (pp. 273-274).

[2] N. Chavoshi, H. Hamooni, and A. Mueen, "Debot: Twitter bot detection via warped correlation," in 2016 IEEE 16th International Conference on Data Mining (ICDM), Dec 2016.

# BotOrNot: Ambiguous probability model



Most scores fall in the range of 0.4 to 0.6 (uncertainty)

# Existing bot detection methods

| | Supervised |
|---|---|
| Human intervention | Yes |
| Unit of detection | Individual |
| State-of-art Application | BotOrNot[1] |
| Bot accounts overlapped with our protocol | N/A |

[1] Davis, C. A. et al, 2016
[2] N. Chavoshi et al, 2016

8

# Existing bot detection methods

| | **Supervised** | **Unsupervised** ✓ |
|---|---|---|
| Human intervention | Yes | No |
| Unit of detection | Individual | Group |
| State-of-art Application | BotOrNot[1] | DeBot[2] |
| Overlap with bots detected by our protocol | N/A | Mean 11.69% Std 7.48% |

[1] Davis, C. A. et al, 2016
[2] N. Chavoshi et al, 2016

# New unsupervised approach

Detect **groups** of accounts tweeting **similar texts** over a long period of time

Why? Duplicate tweeting is widely used to send spam, to bait user into visiting sides and to inflate SEO results.

Collect tweets with **embedded (shortened) URLs**



**Enhanced Bet Offers** @enhancedoffers
888 Bet £10 & Get £30 in Free Bets
Use code 30F
New customers only
T&Cs apply,18+
JOIN HERE
bit.ly/88830fr

**Boost A Bet** @boostbets · 1h
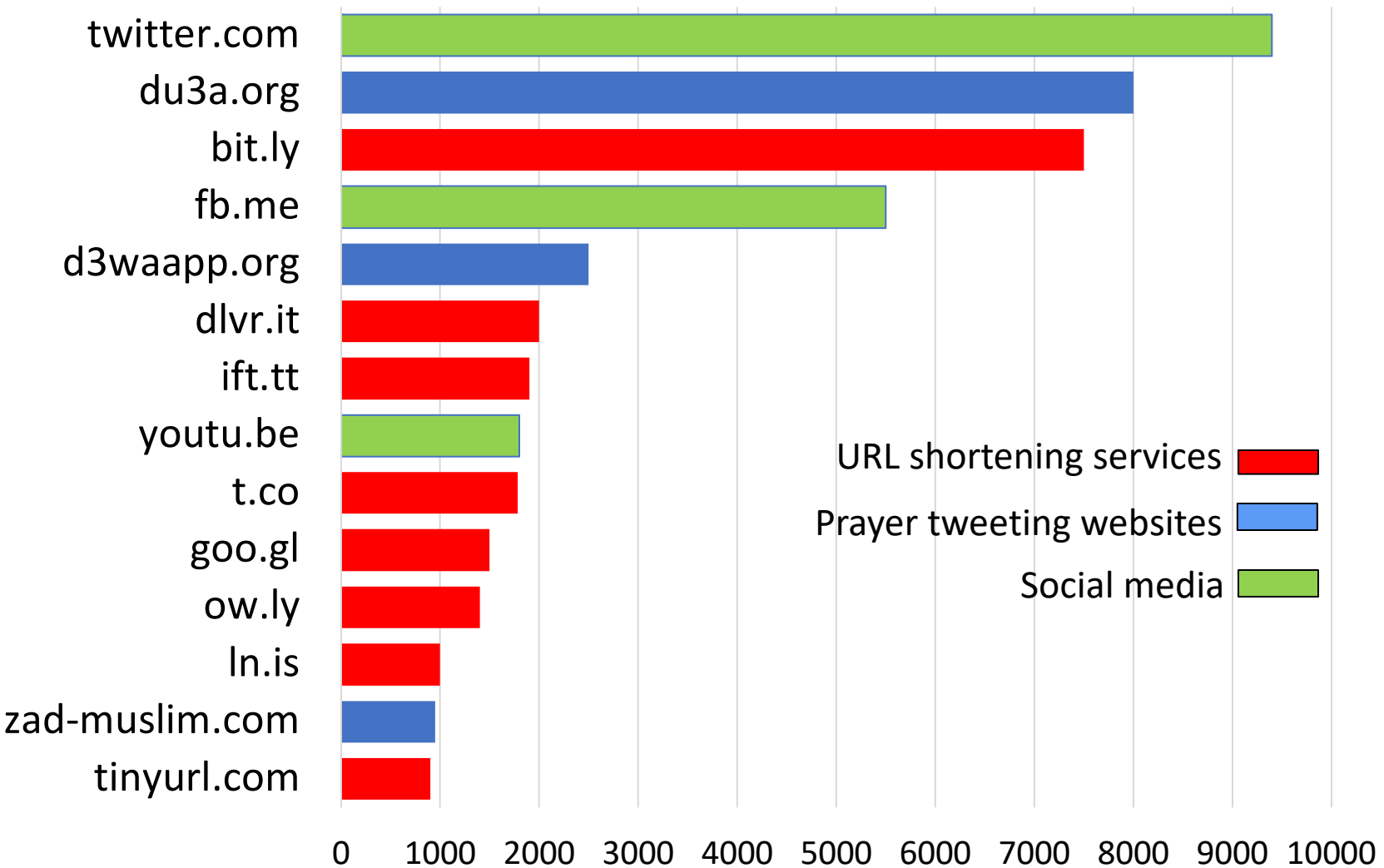888 Bet £10 & Get £30 in Free Bets
Use code 30F
New customers only
T&Cs apply,18+
JOIN HERE
bit.ly/88830fr
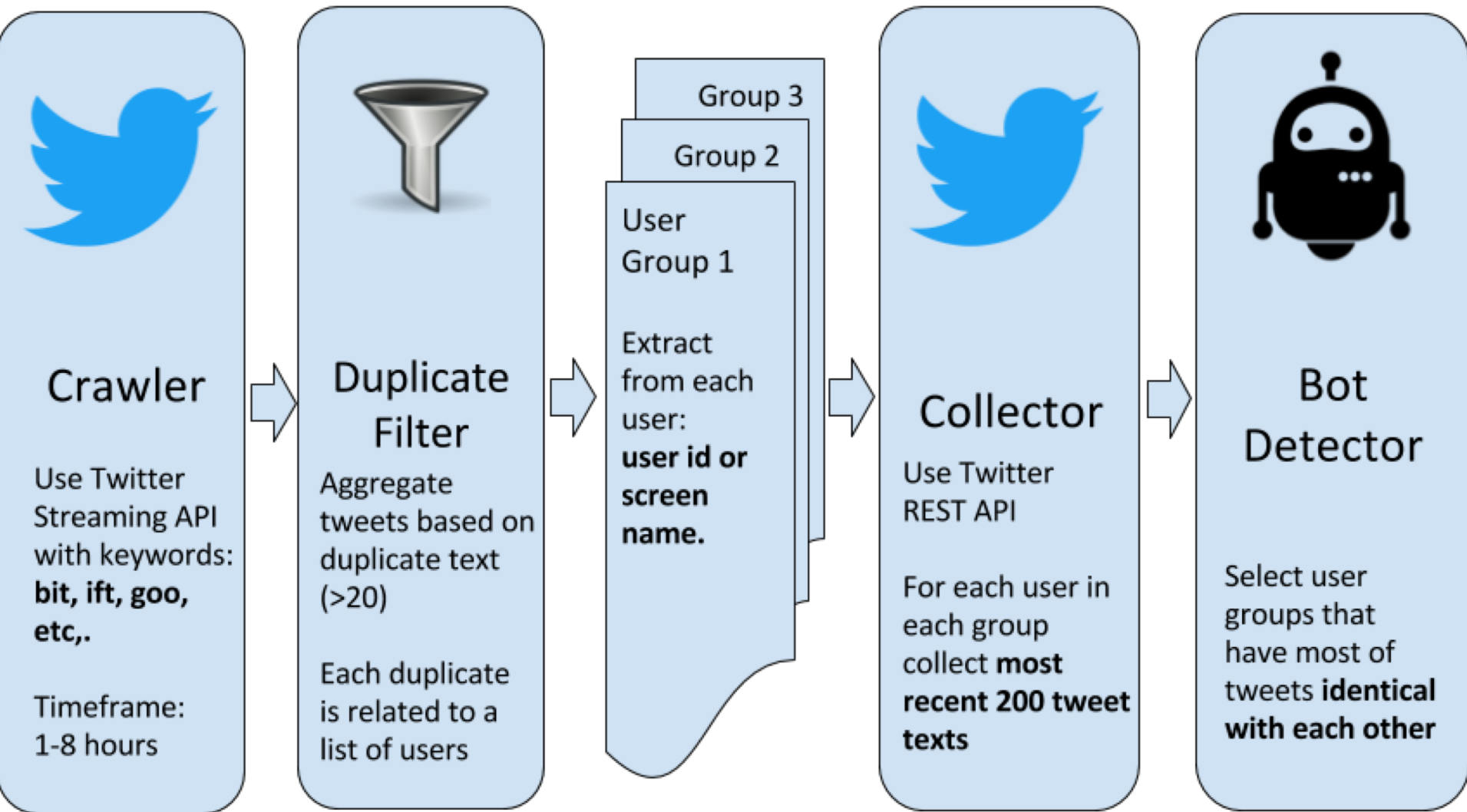
# Why focus on shortened URLs?



Real-time Trending URLs on Twitter

Legend:
- URL shortening services (red)
- Prayer tweeting websites (blue)
- Social media (green)

Bar chart data (approximate values):
- twitter.com — 9400 (Social media)
- du3a.org — 8000 (Prayer tweeting websites)
- bit.ly — 7500 (URL shortening services)
- fb.me — 5500 (Social media)
- d3waapp.org — 2500 (Prayer tweeting websites)
- dlvr.it — 2000 (URL shortening services)
- ift.tt — 1900 (URL shortening services)
- youtu.be — 1800 (Social media)
- t.co — 1800 (URL shortening services)
- goo.gl — 1500 (URL shortening services)
- ow.ly — 1400 (URL shortening services)
- ln.is — 1000 (URL shortening services)
- zad-muslim.com — 950 (Prayer tweeting websites)
- tinyurl.com — 900 (URL shortening services)

X-axis: 0, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000, 9000, 10000

11

# How system detects bots

**Crawler**

Use Twitter Streaming API with keywords: **bit, ift, goo, etc,.**

Timeframe: 1-8 hours

**Duplicate Filter**

Aggregate tweets based on duplicate text (>20)

Each duplicate is related to a list of users

Group 3

Group 2

User Group 1

Extract from each user: **user id or screen name.**

**Collector**

Use Twitter REST API

For each user in each group collect **most recent 200 tweet texts**

**Bot Detector**

Select user groups that have most of tweets **identical with each other**

# How system detects bots

**Algorithm 1** Algorithm for detecting botnets

**Input:** $\alpha$ (minimum duplicate factor), $\beta$ (overlap ratio),

a group $G$ of $n$ accounts $a_1, \ldots, a_n$,

sets $T(a_1), \ldots, T(a_n)$ of tweets where $T(a_i) = \{t_{i1}, \ldots, t_{i200}\}$ of the 200 most

recent tweets for each account $a_i, 1 \leq i \leq n$

1: $C = \emptyset$ /* most frequent tweet set */

2: $S = \emptyset$ /* bot account set */

3: **for** each user $a_i \in G$ **do**

4:     **if** $(|\{i \mid t \in T(a_i); 1 \leq i \leq n\}| \geq \alpha)$ **then**

5:         $C = C \cup \{t\}$

6:     **end if**

7: **end for**

8: **for** each user $a_i \in G$ **do**

9:     **if** $(a_i \in S \iff \frac{|T(a_i) \cap C|}{|T(a_i)|} \geq \beta)$ **then**

10:         $S = S \cup \{a_i\}$

Step 1: construct a set of common tweets (α)

Step 2: find users whose tweets overlap with the common set (β)

13

# Experimental Results[1] (500,000 tweets/URL)

| URL Shortening Services | Total # of accounts | Total # of bots | % bots suspended by Twitter until 6/10/17 | % bots suspended by Twitter until 7/17/17 | % bots suspended by Twitter until 9/25/17 |
|---|---|---|---|---|---|
| bit.ly | 28964 | 696 | 3.74% | 4.74% | 8.9% |
| ift.tt | 12543 | 321 | 2.80% | 9.97% | 10.59% |
| ow.ly | 28416 | 894 | 45.30% | 48.21% | 48.43% |
| tinyurl.com | 20005 | 705 | 5.39% | 7.66% | 12.34% |
| dld.bz | 6893 | 304 | 8.22% | 11.84% | 18.75% |
| viid.me | 2605 | 129 | 38.76% | 55.81% | 63.57% |
| goo.gl | 11250 | 710 | 0.42% | 3.24% | 7.04% |
| dlvr.it | 15122 | 1194 | 7.37% | 9.13% | 9.46% |
| ln.is | 25384 | 5857 | 1.11% | 1.25% | 1.50% |

[1] Z. Chen, R. S. Tanash, R. Stoll, and D. Subramanian, Hunting Malicious Bots on Twitter: An Unsupervised Approach. Cham: Springer International Publishing, 2017, pp. 501–510. [Online]. Available: https://doi.org/10.1007/978-3-319-67256-4 40

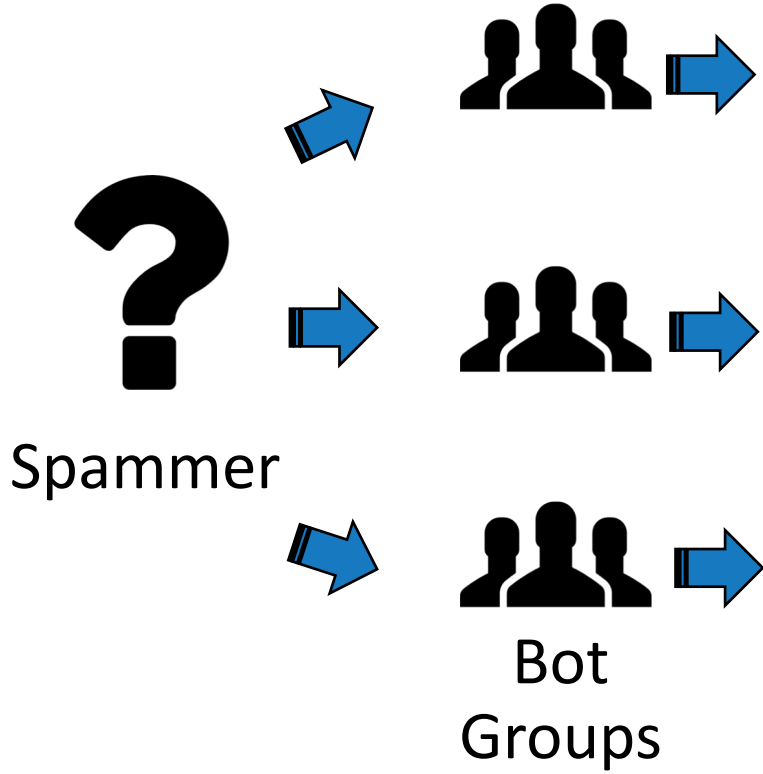# New Data Collection (2 month study)

- 70+ days (09/02/2017 to 11/14/2017)
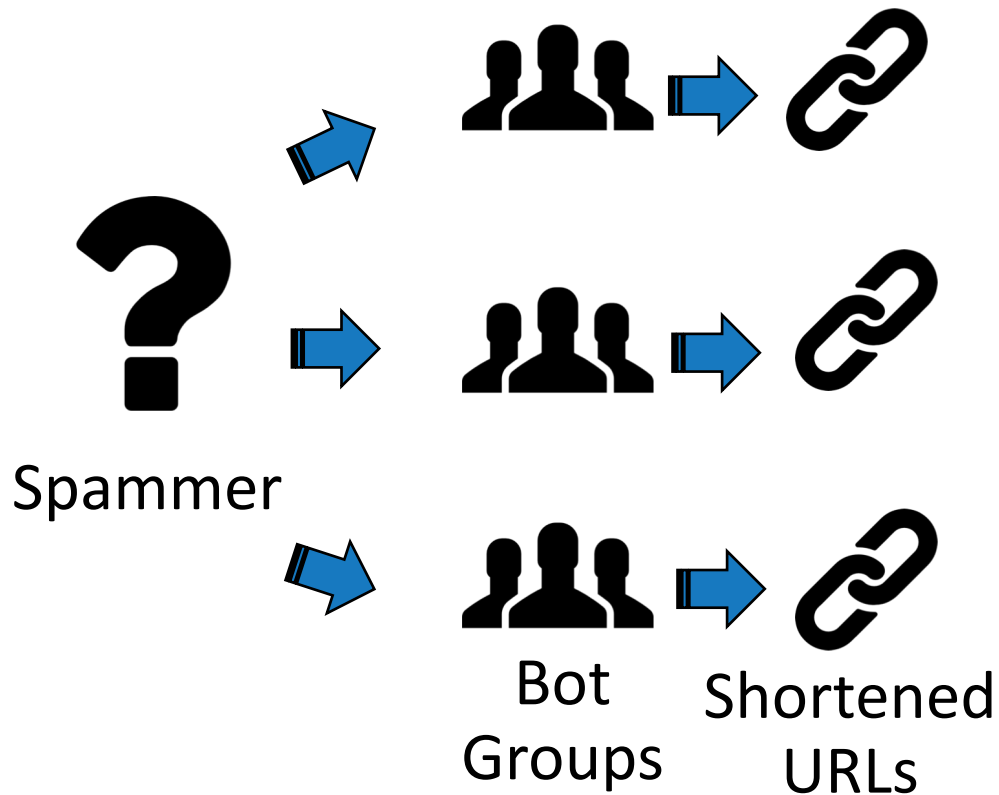- 7 URL shortening services
- 30000 tweets collected per service per day

# Bot traffic accounts for **10-50%** of tweets with shortened URLs
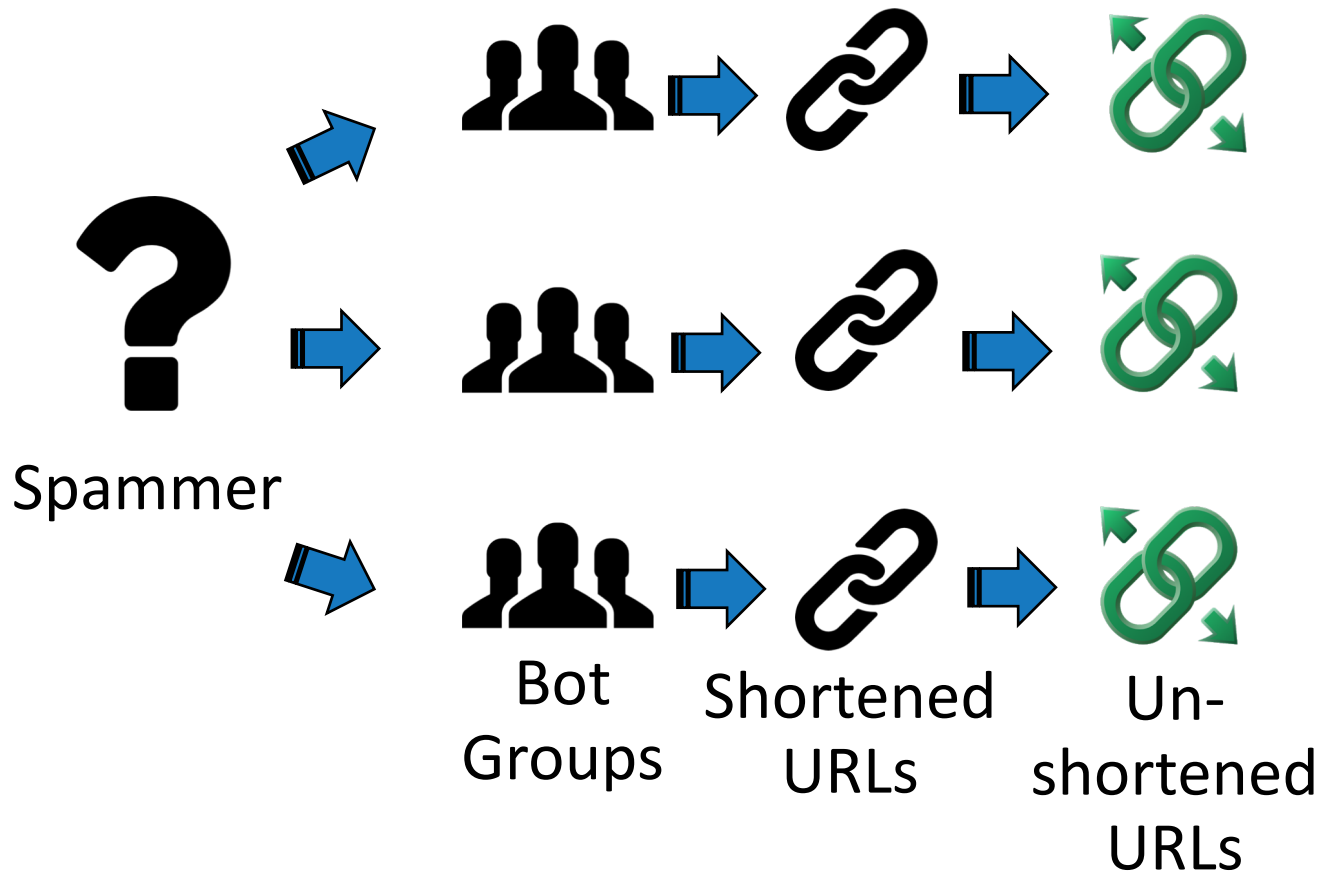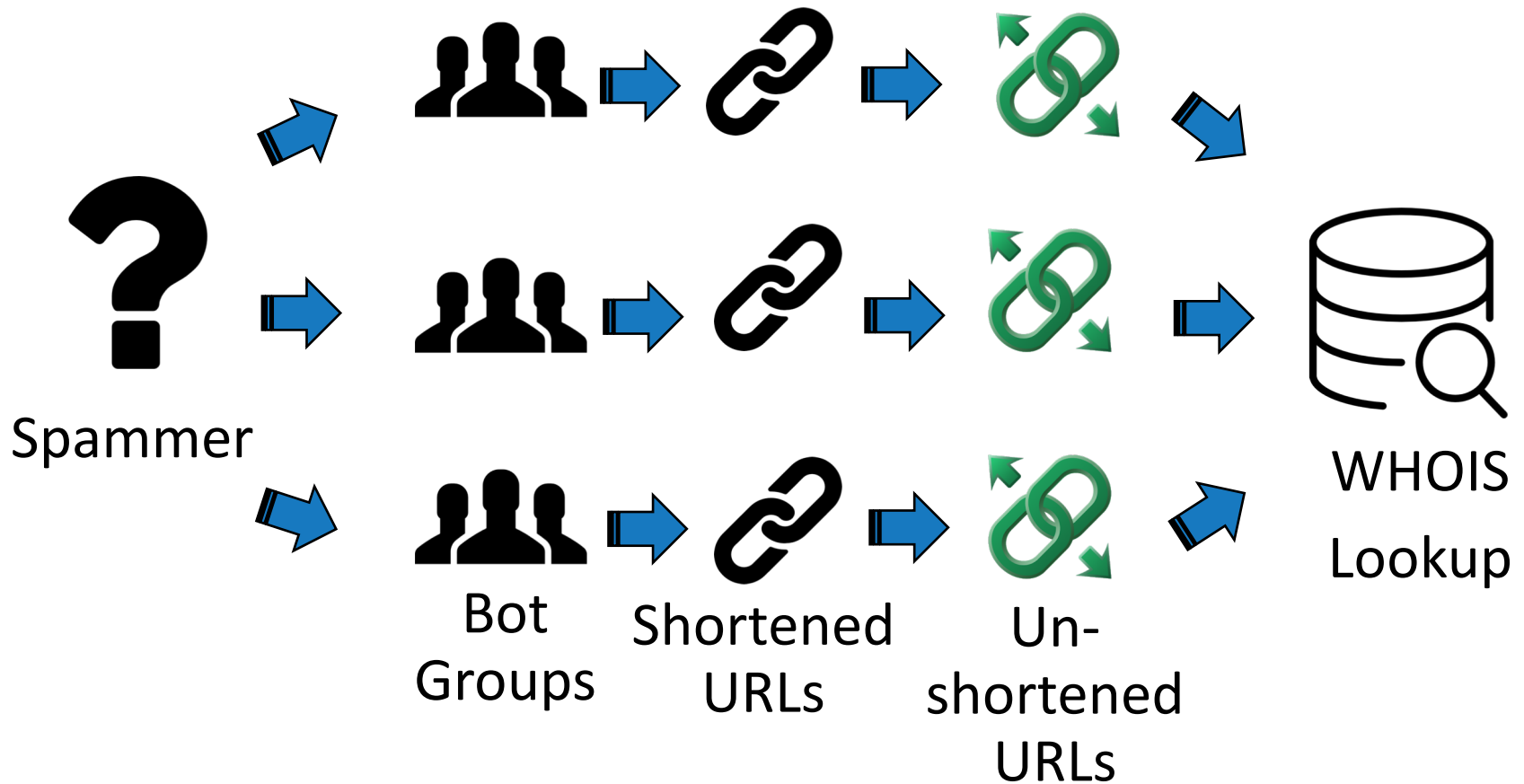
# From bot group to spam campaign



Spammer
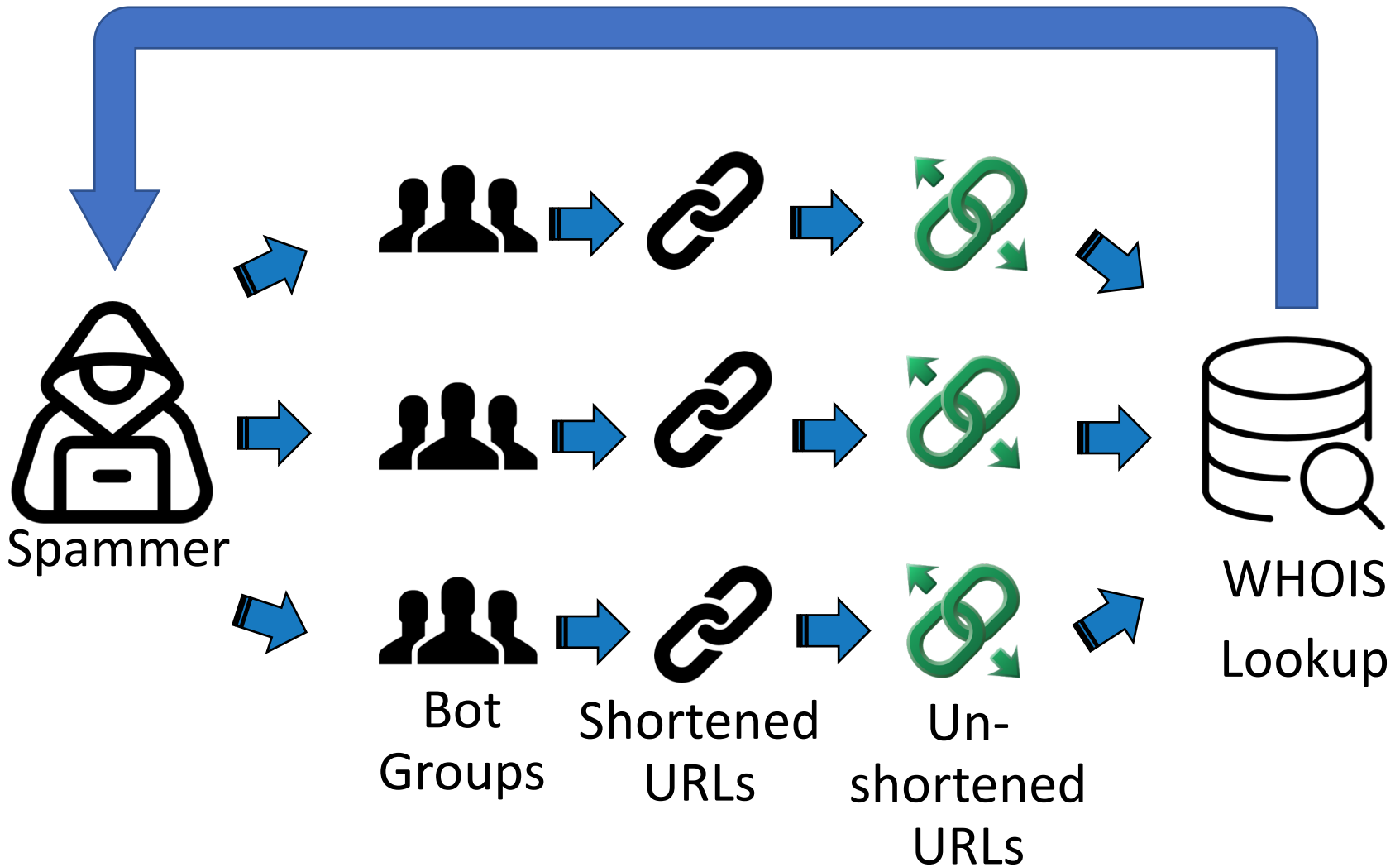
Bot
Groups

# From bot group to spam campaign



Spammer

Bot Groups

Shortened URLs

# From bot group to spam campaign

Spammer

Bot Groups

Shortened URLs

Un-shortened URLs

# From bot group to spam campaign



Spammer

Bot Groups

Shortened URLs

Un-shortened URLs

WHOIS Lookup

# From bot group to spam campaign



Spammer → Bot Groups → Shortened URLs → Un-shortened URLs → WHOIS Lookup

21

# From botnet to spam campaign

| From 09/02/2017 to 11/14/2017 | |
|---|---|
| # bot accounts identified | 200,379 |
| # bot groups | 7,350 |
| # suspicious registrants | 848 |

| Giuseppe Malfitano | Shashank Vaishnav | Proxy Server |
|---|---|---|
| i5-**news**.com | **awesome**nature.info | newbuye.**review** |
| a6-**news**.com | **awesome**pix.info | vidisp.**review** |
| a8-**news**.com | **awesome**post.info | superdoppy.**review** |
| i5-**news**.com | **awesome**stuff.info | situari.**review** |
| i7-**news**.com | **awesome**thingz.info | sacraffm.**review** |

# From botnet to spam campaign



Input URL

IFTTT

bitly dlvr.it

Shortened URL

Register URL shortening services

Spammer → Create bots

Register URL/domain

tweet

tweet

tweet

Redirect to malicious websites

Case study 1: #UmbrellaRevolution
          Remove bots for
          community detection

Case study 2: #ReleaseTheMemo
          Track how bots interfere
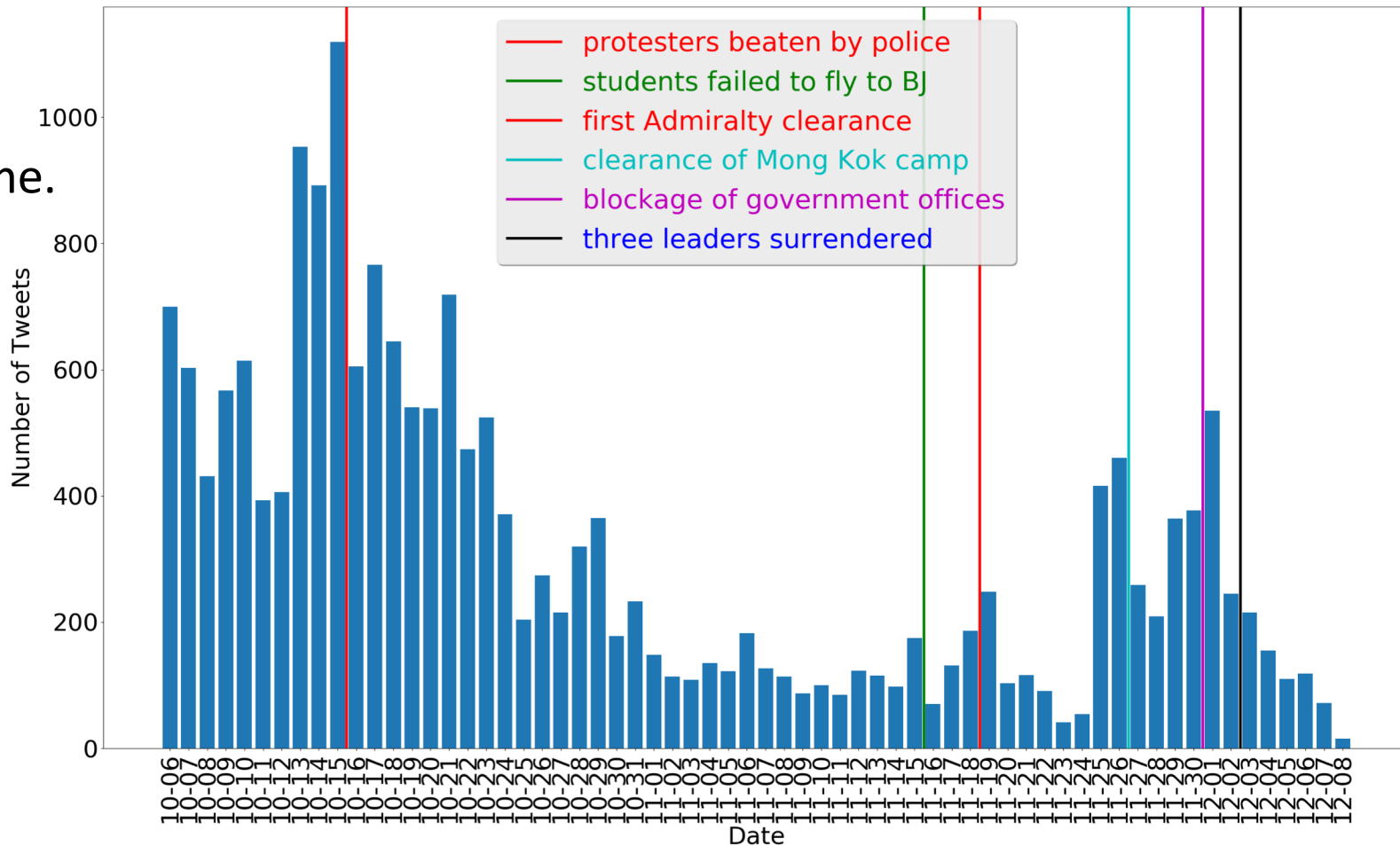          with political discussions

# Case study 1: #UmbrellaRevolution

▪ **Background**: The Umbrella Revolution was a large scale social movement in Hong Kong started in late September 2014 and ended in December 2014.

▪ **Goal**: Understand human interaction on social media.

▪ **Challenge**:
  ➢ Design a filtering mechanism to remove bots.
  ➢ Community detection using tweet-retweet graph

# Case study 1: # UmbrellaRevolution

- Collected live tweets from Streaming API
- Time collected:        10/06/2014 – 12/08/2014
- # tweets collected:   1,062,606

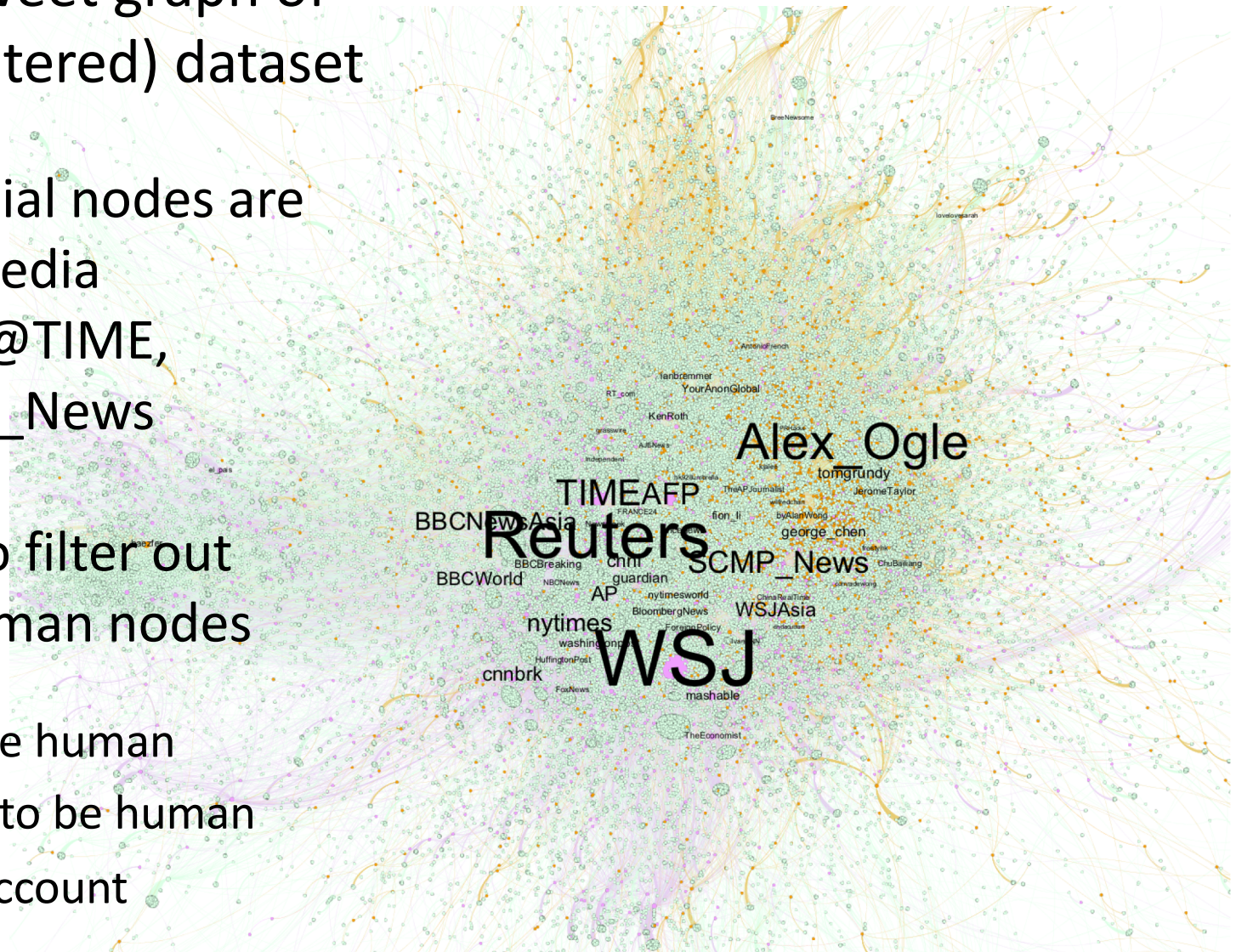Right: daily tweet volume. Peaks correspond with major events.



Legend:
- protesters beaten by police
- students failed to fly to BJ
- first Admiralty clearance
- clearance of Mong Kok camp
- blockage of government offices
- three leaders surrendered

# Case study 1: # UmbrellaRevolution

Tweet-retweet graph of
raw (not filtered) dataset

1. Influential nodes are
   news media
   @WSJ, @TIME,
   @SCMP_News

2. Need to filter out
   non-human nodes

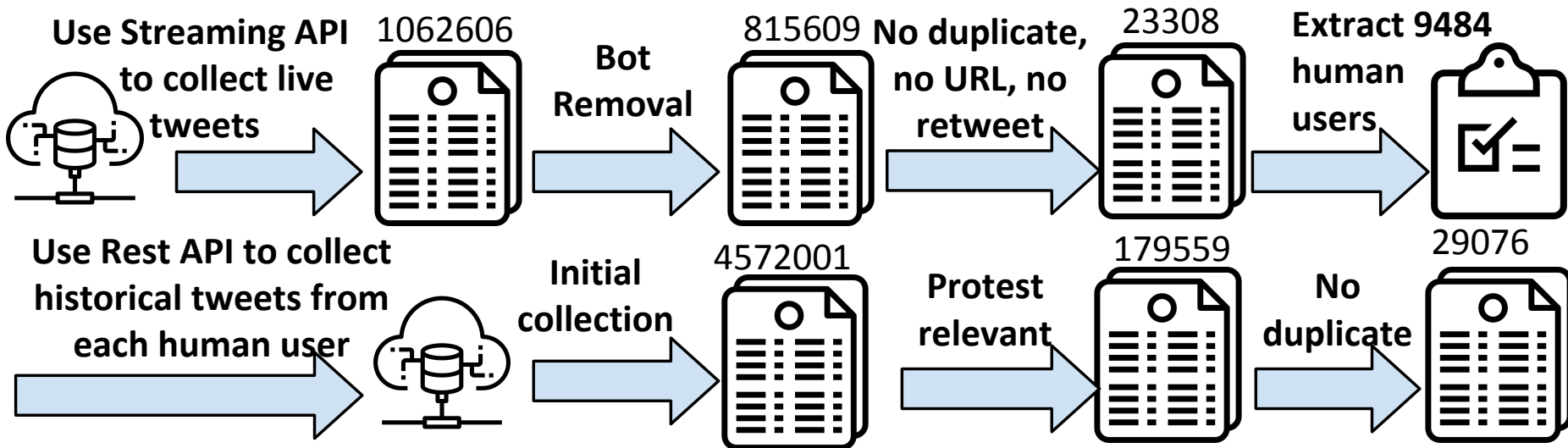🟠 likely to be human

🟢 not likely to be human

🟣 verified account

# Case study 1: # UmbrellaRevolution

Data processing pipeline:
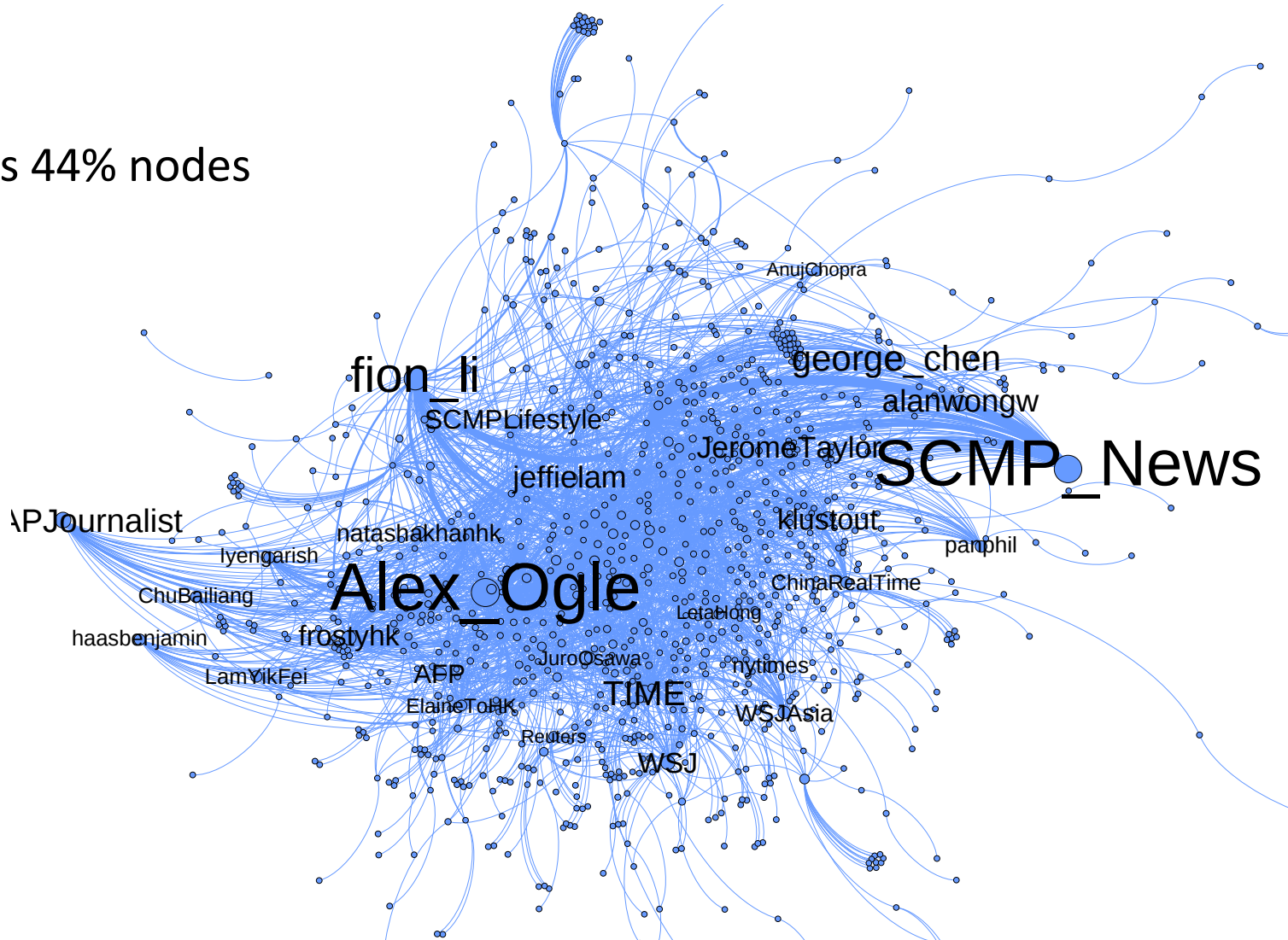
      stage 1: filter out bots

      stage 2: collect more human tweets

**Use Streaming API to collect live tweets** 1062606 **Bot Removal** 815609 **No duplicate, no URL, no retweet** 23308 **Extract 9484 human users**

**Use Rest API to collect historical tweets from each human user** **Initial collection** 4572001 **Protest relevant** 179559 **No duplicate** 29076

# Case study 1: # UmbrellaRevolution

Tweet-retweet graph of journalist community



Journalists 44% nodes

# Case study 1: # UmbrellaRevolution

Tweet-retweet graph of activist community

Activists 56% nodes

# Case study 1: # UmbrellaRevolution

Tweet-retweet graph of both communities

- 🟣 Journalists 44% nodes
- 🟢 Activists 56% nodes

# Case study 1: # UmbrellaRevolution

What we learn:

There are two major communities discussing this event on Twitter

Top three news accounts (**journalist**)

The Wall Street Journal ✔
@WSJ

SCMP News ✔
@SCMP_News

Alex ✔
@Alex_Ogle

Top three pro-protest accounts (**activist**)

Umbrella Movement
@hk928umbrella

HKDemoNow
@hkdemonow

學聯 HKFS
@HKFS1958

# Case study 2: #ReleaseTheMemo

Goal: Track activities of political bots

**#ReleaseTheMemo** exploded on Twitter

On Feb. 2, 2018, the United States House Intelligence Committee Chairman Devin Nunes, released a controversial memo

| 2018-1-18 | 2018-1-23 | 2018-2-2 | 2018-2-3 |
|-----------|-----------|----------|----------|

| Datasets | pre memoday | memoday | post memoday |
|----------|-------------|---------|--------------|
| # Tweets | 99999 | 253383 | 54424 |
| Duration | 3h:25m:16s | 4h:10m:12s | 4h:28m:04s |

# Case study 2: #ReleaseTheMemo

| Datasets | #accounts | #bot accounts | %bot accounts | %bot tweets |
|---|---|---|---|---|
| pre memoday | 36347 | 4030 | 11.1 | 18.9 |
| memoday | 67654 | 11254 | 13.1 | 26.7 |
| post memoday | 30764 | 3718 | 12.1 | 15.9 |

Number of bots and bot tweets in three dataset. Bot activities peaked on **memoday**

# Case study 2: #ReleaseTheMemo

Bots retweet from



Verified account
@dbongino

Parody account
@sean_spicier

Influential bots
@DanCovfefe1

# #ReleaseTheMemo – Pre memoday



Normal Accounts 87.07%

# #ReleaseTheMemo – Pre memoday



Normal Accounts 87.07%
Normal Bots        12.36%
Influential Bots      0.14%

# times retweeted > 50

# #ReleaseTheMemo – Pre memoday

|  | Pro-Trump | Anti-Trump |
|---|---|---|
| 🟠 Normal Bots | 88% | 12% |
| 🔴 Influential Bots | 100% | 0% |

🟣 Normal Accounts 87.07%
🟠 Normal Bots 12.36%
🔴 Influential Bots 0.14%
🟢 Verified Accounts 0.43%

@sean_spicier: not the real Spicer!
(word is misspelled)

# #ReleaseTheMemo – Memoday

MichaelSteele
TeamPelosi
MarkWarner
RepSwalwell
tedlieu
krassenstein
RanttMedia
ProudResister
activist360

⬤ Normal Accounts 87.07%

ClintonM614
AmericanVoterUS
Trumperland
southern4MAGA
mike_Zollo
DonnaWR8
1776Stonewall
bocavista2016
steph93065
StockMonsterVIP
mflynnJR
he_Trump_Train
Education4Libs
RepStevenSmith
ScottPresler
Makada_
FoxNews
TheRealJulian
DaveNYviii
DonaldJTrumpJr
Jim_Jordan
MarkDice
RealJamesWoods
PrisonPlanet
thecjpearson
RandPaul

# #ReleaseTheMemo – Memoday



Normal Accounts 85.43%
Normal Bots     14.11%
Influential Bots  0.07%

# #ReleaseTheMemo – Memoday



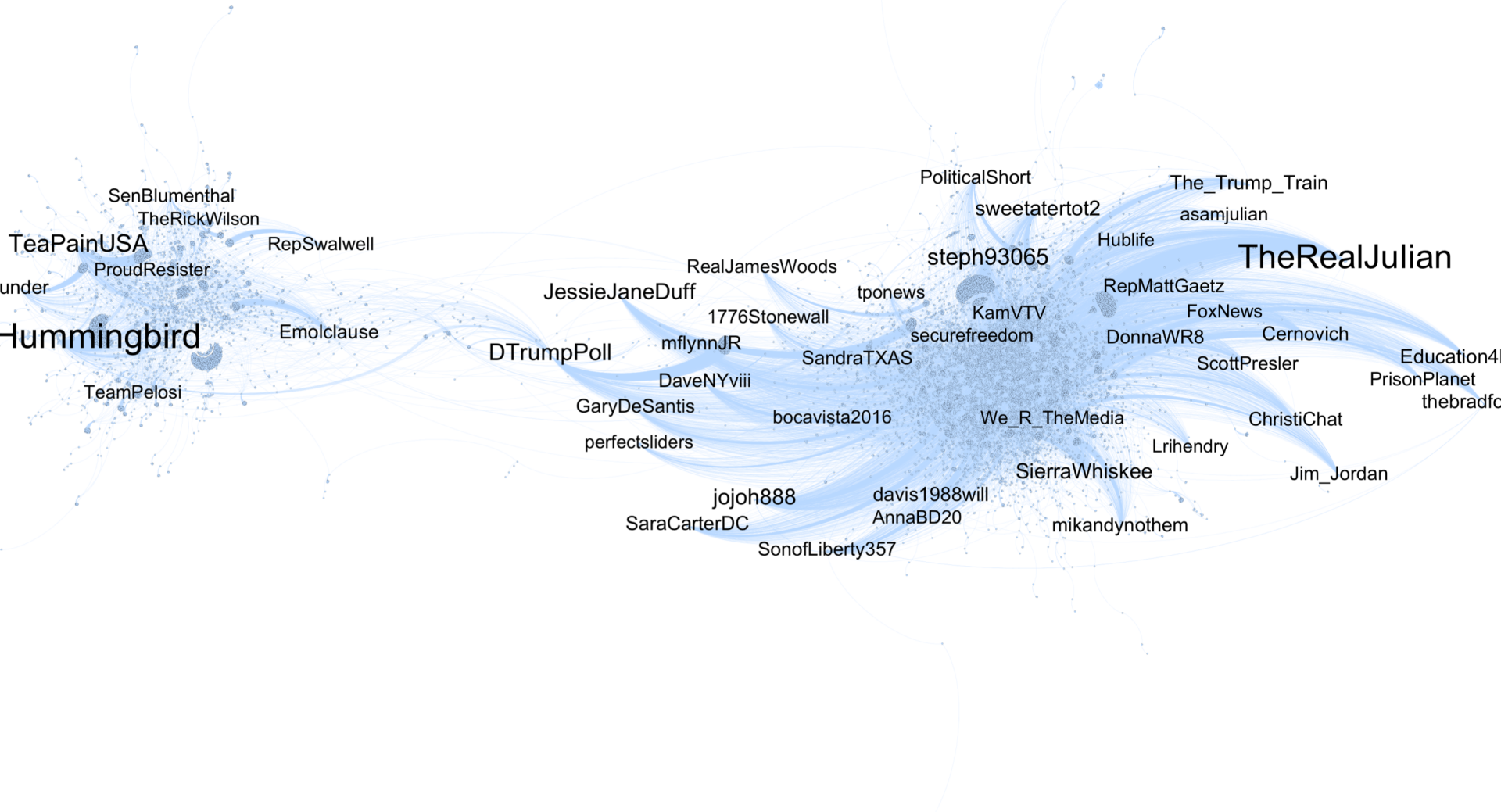|                  | Pro-Trump | Anti-Trump |
|------------------|-----------|------------|
| 🟠 Normal Bots     | 68%       | 32%        |
| 🔴 Influential Bots | 100%      | 0%         |

🟣 Normal Accounts 85.43%
🟠 Normal Bots 14.11%
🔴 Influential Bots 0.07%
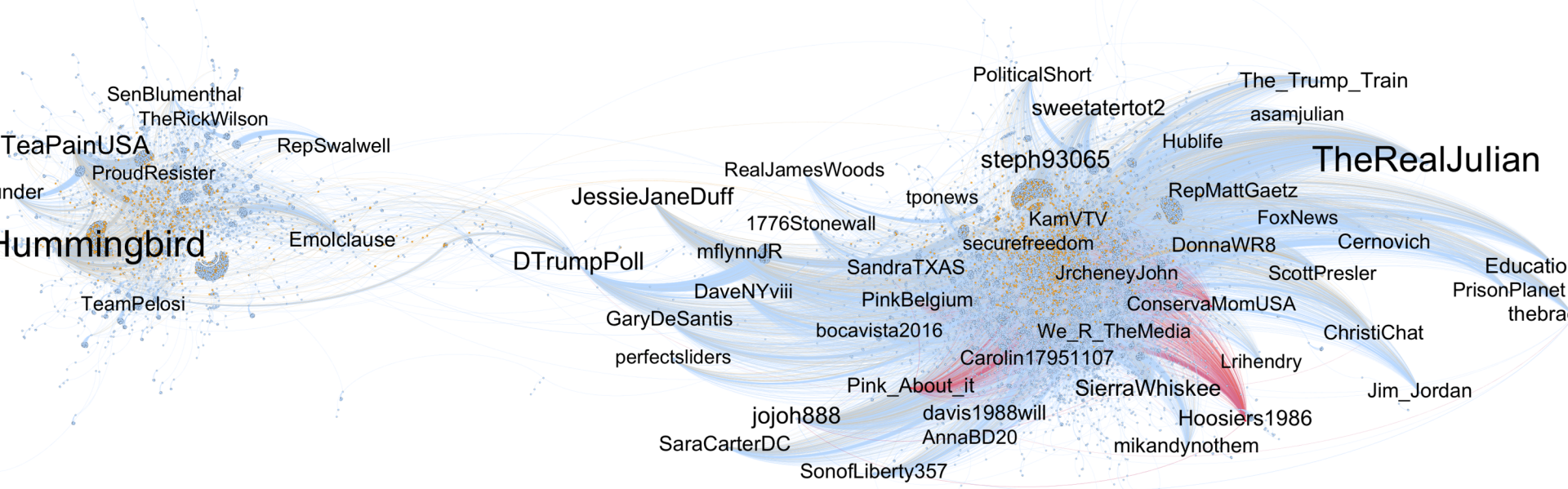🟢 Verified Accounts 0.38%

@TheRealJulian: not the real Julian!

# #ReleaseTheMemo – Post memoday
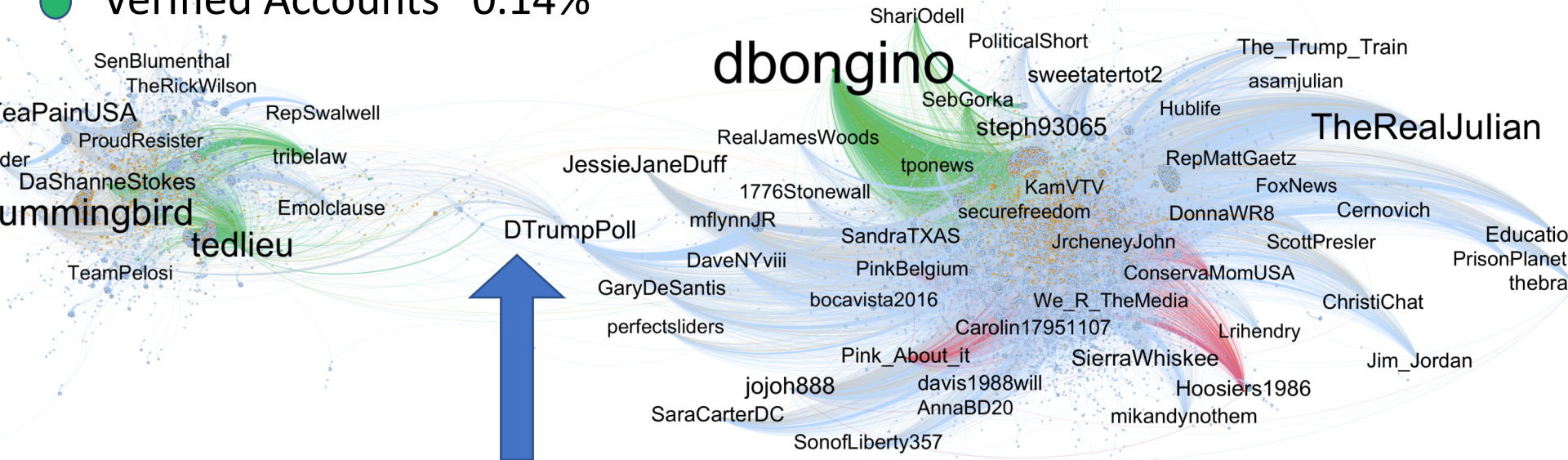


Normal Accounts 87.10%

# #ReleaseTheMemo – Post memoday

- Normal Accounts 87.10%
- Normal Bots         12.73%
- Influential Bots     0.02%

# #ReleaseTheMemo – Post memoday

Normal Accounts 87.10%
Normal Bots 12.73%
Influential Bots 0.02%
Verified Accounts 0.14%

|  | Pro-Trump | Anti-Trump |
|---|---|---|
| Normal Bots | 55% | 45% |
| Influential Bots | 100% | 0% |



@DTrumpPoll: Impartial polls about Trump

# Case study 2: #ReleaseTheMemo

## What do we observe?

- 15% of pro-Trump cluster are bots and 13% of anti-Trump cluster are bots.

- Bots are artificially making trending hashtags
  21% #ReleaseTheMemo, 24% #MemoDay,
  92% # SecretSociety, 34% #IAmNOTaRussianBot
  tweets are generated by bots.

- There are still bots in the dataset that we do not identify. Having access to account registration information would be helpful.

# Application, Action and Impact of our bot detection work

# Bot Detection (Impact)

1. Application: Twitter Bot Monitor
   - Backend: Bot Detection, Spam Campaign Detection, API
   - Frontend: Bot Visualization, Bot Trend Monitor, Trending URL Monitor
   - To date, our Twitter Bot Monitor is still tracking and collecting suspicious accounts (http://water.clear.rice.edu:18000/)

2. Publications
   - Paper published on 2017 International Conference on Social Informatics
   - Presented our work at Oxford University, UK
   - Another paper submitted to IEEE transactions on intelligent systems is under review
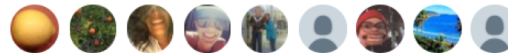
# Bot Detection (Impact)



Online discussion on Twitter

# Bot Detection (Impact)

**MOTHERBOARD**

*VICE*

# Why Twitter Is the Best Social Media Platform for Disinformation

estimated that **up to 15 percent** of all Twitter accounts are bots. In September, **another study** from Rice University put the number at up to 23 percent, out of a global active user base of **approximately 330 million**.

Media coverage, paper cited by *Vice* (November 1, 2017)
*https://motherboard.vice.com/en_us/article/bj7vam/why-twitter-is-the-best-social-media-platform-for-disinformation*

# Contact and response

I reached out to:

- URL Shortening Services (bitly, tiny url, hootsuite, tiny cc, dlvr, ifttt)
- Domain registrars (Namecheap, GoDaddy)
- Domain hosting services (Tiggee, Liquid Web, Digital Ocean)
- Google Network Abuse Team, Google Safe Browsing
- Social Media Company (Twitter, last December)

# Contact and response

**URL Shortening Services**

**Domain Registrars**

**Web Hosting Services**

**Browser Services**

**Twitter**

# Contact and response (Four responded)

# Namecheap replied, but said cannot take action

Thank you for the detailed explanation of the issue.

Unfortunately, we were unable to validate your claim(s), since in this situation, Namecheap acts as the registrar only. Our ability to investigate the matter is limited since the content transmitted via the website is not located on our server.

Considering the aforementioned points, we recommend that you contact the hosting provider, who would be in a better position to validate your claim(s) and take the appropriate action. For your convenience, here are the contact details of the company that owns the IP-Addresses assigned to the subject domains: http://whois.domaintools.com/192.241.145.46

Additionally, if you believe you are aware of an attempted crime, you can file a complaint through Internet Crime Complaint Center at https://complaint.ic3.gov , who are in the best position to fully investigate any such issue across any/all service providers.

Please let us know should you have any further questions.

# Tiny.cc replied and took down reported URLs

Hi Zhouhan,

Thanks for contacting me and thanks for your work on this project. If you have more suspicious domains, please share them. We have our own internal blacklist of domains that can be added to.

Abuse is probably the biggest challenge to running a free URL shortening service as abusers put a great amount of effort and ingenuity into new methods.

We have a large stack of filters to test each URL before it is allowed to be shortened. This "validation" stack includes checking against 3rd party blacklists (Google Safe Browsing API, DNS queries to SURBL, Spamhaus, etc.)  Mostly checking and filtering at the domain level. But there are other patterns of abuse that we have recognized over the years and try to detect at the front and stop it before URL is shortened.
As you know, spam, phish and abuse is largely a reactive game, as there is no way to proactively know for sure how a URL will be used or abused.

Best regards,

# Ow.ly replied but said Twitter should take action

Hi Zhouhan,

Thanks for reaching out.

We appreciate your efforts on reporting this kind of spammy behaviour. As for performing any action on the links, I'm afraid we can't just remove these links.

We do remove links that violate copyright or contain phishing/malware, but this kind of content is not agains our ToS. While we make efforts to detect bots on our system, we rely on the end social network (Twitter in this case) to be the one flagging and terminating the social media accounts.

Please forward us the whole list of suspicious domains and ow.ly links, so we can correlate with our users and monitor their activity.

Thanks again for your help.

# Bit.ly did not reply, but took down reported URLs



**STOP** - there might be a problem with the requested link

The link you requested has been identified by bitly as being potentially problematic. This could be because a bitly user has reported a problem, a black-list service reported a problem, because the link has been shortened more than once, or because we have detected potentially malicious content. This may be a problem because:

- Some URL-shorteners re-use their links, so bitly can't guarantee the validity of this link.
- Some URL-shorteners allow their links to be edited, so bitly can't tell where this link will lead you.
- Spam and malware is very often propagated by exploiting these loopholes, neither of which bitly allows for.

The link you requested may contain inappropriate content, or even spam or malicious code that could be downloaded to your computer without your consent, or may be a forgery or imitation of another website, designed to trick users into sharing personal or financial information.

**bitly suggests that you**
- Change the original link, and re-shorten with bitly
- Close your browser window
- Notify the sender of the URL

Or, continue at your own risk to
http://www.sexyarb.com/vcdKNsycK.html

# Reaching out to Twitter

- On December 7, 2017, we gave an internal presentation to Twitter Content Quality team and Data Science team

- Twitter thanked us for our work and presentation, and introduced us to data scientists and engineers working on anti-spam topics

# Reaching out to Twitter

- On February 21, 2018, Twitter rolled out an update of its anti-spam policy[1]

- The policy explicitly tells Twitter service providers "***Do not (and do not allow your users to) simultaneously post identical or substantially similar content to multiple accounts***."

- This is exactly the criteria of bots defined in our work.

[1] https://blog.twitter.com/developer/en_us/topics/tips/2018/automation-and-the-use-of-multiple-accounts.html

# Who is more responsible?

- URL shortening services are responsive and **willing to cooperate**.

- Domain registrars **cannot take action** if the website is hosted on another IP.

- Domain hosting services are **unresponsive**. If they don't take action, spammers will keep abusing other services.

# Future work

- Investigate new types of malicious activities.

- Recently we found bots tweeting cryptojacking links[1]

- They are websites secretly running cryptocurrency mining script in one's browser, consuming CPU power.



[1] example malicious link: http://technimum.com/blog/tehlukesizlik/6554.html

# Future work

- Update detection algorithm to catch new types of bots.
- Recently found bots truncating texts[1] to evade detection

**9 Embarrassing** Times When Selena Gomez Faced Wardrobe Malfunction

**Embarrassing** Times When Selena Gomez Faced Wardrobe Malfunction

**mbarrassing** Times When Selena Gomez Faced Wardrobe Malfunction

**barrassing** Times When Selena Gomez Faced Wardrobe Malfunction

**rassing** Times When Selena Gomez Faced Wardrobe Malfunction

…

[1] Example final landing URL: http://loveforsomething.com/s1onnq-9-malfunction-when-faced-embarrassing-gomez-times-sd24b

# Conclusions

- Our unsupervised detection system detects malicious accounts and spam campaigns 24/7 without human intervention.

- Attackers and spammers are evolving and getting more sophisticated.

- Academia and Industry have to work together to develop better algorithms and to implement stricter policies.