



香港中文大學  
The Chinese University of Hong Kong

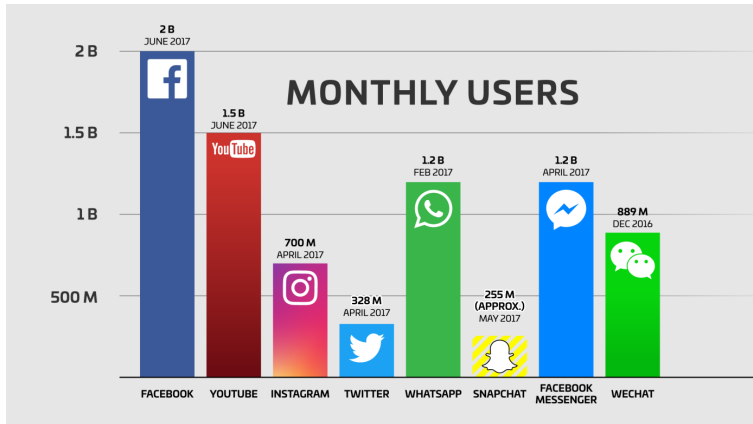


IEEE ICNP 2018

# Sybil Detection in Social-Activity Networks: Modelling, Algorithm and Evaluations

**Xiaoying Zhang**, Xie Hong, John C.S. Lui  
Dept. of Computer Science & Engineering  
The Chinese University of Hong Kong

# Online Social Networks (OSNs)--- Popular & Important



# Millions Fake Users (Sybils) exist on OSNs

FACEBOOK

SOCIAL

## Facebook has disabled almost 1.3 billion fake accounts over the past six months

Facebook will begin publishing more data about how many posts it takes down.

By Kurt Wagner and Rani Molla | May 15, 2018, 10:00am EDT

Technology

BBC



NEWS

## Twitter 'shuts down millions of fake accounts'

🕒 9 July 2018



🔗 Share

# Threats of Sybils

## Fake

Fake Accounts  
Fake Follows  
Fake Likes

Fake  
Social  
Media

TECHNOLOGY

## How Twitter Bots Are Shaping the Election

Between the first two presidential debates, a third of pro-Trump tweets and nearly a fifth of pro-Clinton tweets came from automated accounts.

DOUGLAS GUILBEAULT AND SAMUEL WOOLLEY NOV 1, 2016

## Can social media influence financial markets?



11 Nov 2015

Costas Milas  
Professor, University of Liverpool



This article is published in collaboration with *The Conversation*.

Scottish financial trader James Alan Craig has been charged in the US for allegedly using Twitter to manipulate share prices. According to the US Department of Justice, the 62 year old, from Dunragit in Dumfries and Galloway, caused shareholders to lose more than \$1.6m (£1.1m) after allegedly spreading "fraudulent" information about companies on the social network. According to



**RepublicTV**  
@republictv

Unofficial #RepublicTV satirical website.  
100% #Fake, #Satire, #Parody. All tweets  
are imaginary and fake.

Mumbai, India  
republictv.com

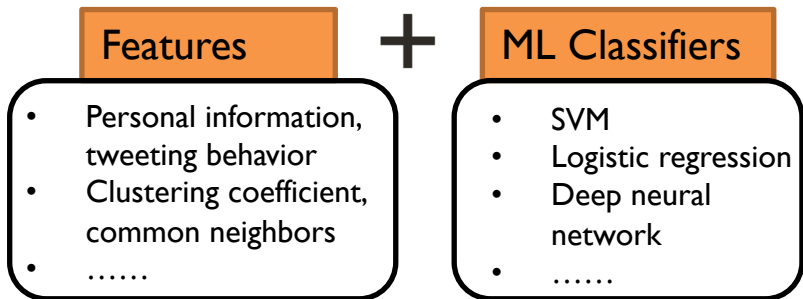


**Republic**  
@republic

Official handle of India's only  
independent news venture. Republic is  
independent. Republic is global. Republic  
is your movement. Join us.

# Existed Sybil detection methods

- **Feature-based methods**

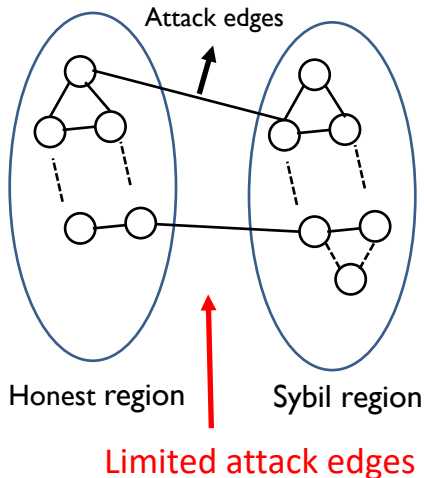


Detect only sybils with  
known patterns !

# Existed Sybil detection methods

- Graph-based methods

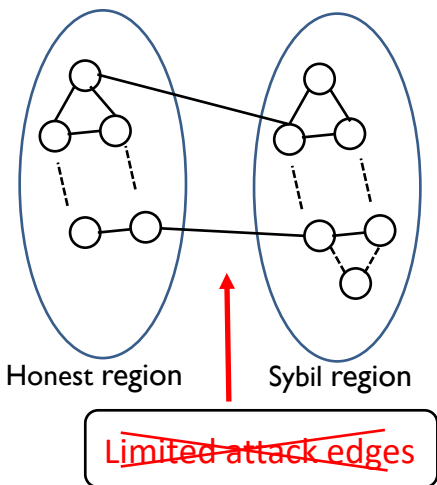
**Limited-attack-edge assumption:** honest users seldom make friends with sybils.



- Techniques

- Random walk  
([SybilRank \[NSDI'12\]](#))
- Fundamental matrix  
([SybilWalk \[DSN'17\]](#))
- Belief propagation  
([SybilBelief \[TIFS'2014\]](#), [SybilScar \[NFCOM'2017\]](#))
- .....

# Inefficiency of Previous Attack Model



- **Limited-attack-edge assumption may not hold !**
  - Sybils can easily befriend with honest users, thus large attack edges exist ([Twitter games](#)[ACSAC'14], [LinkFarm](#)[WWW'12]).
- Breaking the assumption leads to low accuracy.

# Our objective



- What's the **realistic** attack model?
- How to design **efficient algorithms** to detect sybils under realistic attack model?



# Contributions



Realistic Attack  
Model

- The Social-Activity Attack Model



Efficient Detection  
Algorithm

- Sybil\_SAN



Theoretical  
Analysis

- Convergence & Sensitivity analysis



Extensive  
Experiments

- Synthetic & real datasets from Twitter

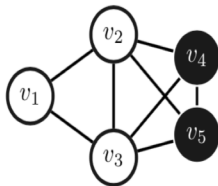
# Social-Activity Network Model: Main intuition

Even honest users befriend with sybils, they **seldom initiate activities to sybils**.

- mention sybils in their own tweets
- reply/retweet sybils' tweets
- ....

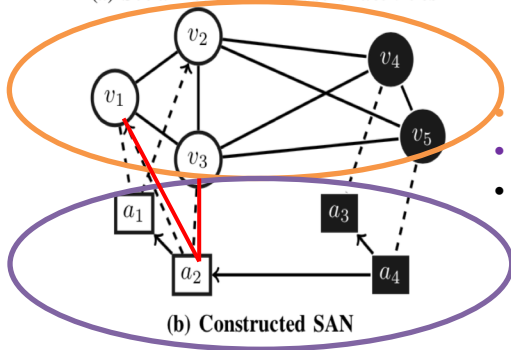
It has been **verified** by the analysis of a dataset containing thousands of real sybils in Twitter by Zhang et.al[TON'2016].

# The Social and Activity Network Model



$(v_1, a_1)$ : user  $v_1$  creates tweet  $a_1$   
 $(v_3, a_2)$ : user  $v_3$  creates tweet  $a_2$   
 $(v_4, a_3)$ : user  $v_4$  creates tweet  $a_3$   
 $(v_5, a_4)$ : user  $v_5$  creates tweet  $a_4$   
 $(a_2, a_1)$ : tweet  $a_2$  retweets tweet  $a_1$   
 $(a_4, a_2)$ : tweet  $a_4$  retweets tweet  $a_2$   
 $(a_4, a_3)$ : tweet  $a_4$  replies tweet  $a_3$   
 $(a_2, v_1)$ : tweet  $a_2$  mentions user  $v_1$   
 $(a_1, v_2)$ : tweet  $a_1$  mentions user  $v_2$

(a) Social network and users' activities



(b) Constructed SAN

Layer 1: friendship graph  $G$

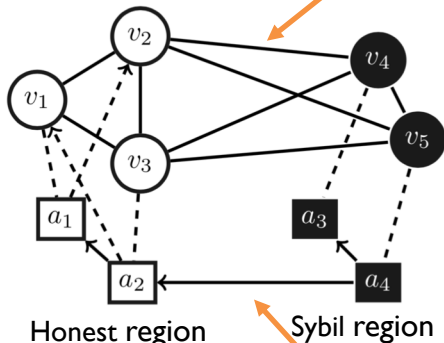
• Layer 2: Activity graph  $\tilde{G}$ .

• Between layers: user-activity mapping graph

# A More Realistic Attack Model

## Friendship attacks ( $N_A$ ):

- Property 1.  $N_A$  can take any value in  $\{0, 1, \dots, |V_h| \times |V_s|\}$ .



## Incoming interactions attacks ( $\alpha W_h$ ):

- Activities initiated from honest users to sybils
- $\alpha \approx 10^{-5}$  ([TON'2016]).
- $W_h$ : # of activities among honest users

## Outgoing interaction attacks ( $\beta W_h$ ):

- Property 2.  $\beta$  can be arbitrary large.

# Contributions



Realistic Attack  
Model

- The Social-Activity Network Model



Efficient Detection  
Algorithm

- Sybil\_SAN



Theoretical  
Analysis

- Convergence & Sensitivity analysis



Extensive  
Experiments

- Synthetic & real datasets from Twitter

# Sybil\_SAN

## Input

Social-Activity-Network

A small set of labelled users

- $S_s$ : sybils
- $S_h$ : honest users)

## Process

Sybil\_SAN

- Initialize nodes' trust/distrust score  $s/s_{dis}$
- Trust/distrust distribution on SAN
- Rank nodes according to  $s + s_{dis}$

## Output

A rank of nodes

nodes with low rank -> sybils

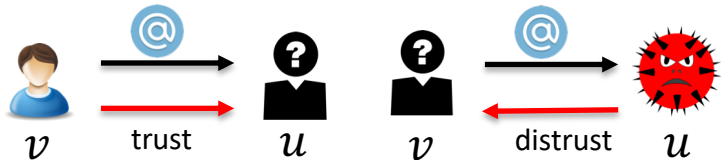
# Sybil\_SAN

- Initialization ( $\mathbf{s}, \mathbf{s}_{dis}$ )

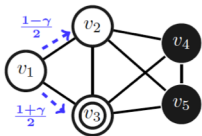
$$\mathbf{s}_i = \begin{cases} \frac{1}{|S_h|}, & i \in S_h \\ \mathbf{0}, & o.w \end{cases} \quad (\mathbf{s}_{dis})_i = \begin{cases} -\frac{1}{|S_s|}, & i \in S_s \\ \mathbf{0}, & o.w \end{cases}$$

# Sybil\_SAN

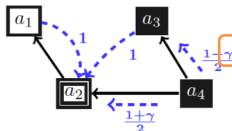
- Trust/distrust Distribution



- Distribution in each layer:
  - personalized pagerank ( $\gamma$ )



(a) friendship graph



(b) activity-following graph

Tradeoff between exploitation and exploration



# Sybil\_SAN: Coupled random walks

- **Mutual reinforcement relationship between users and activities.**
  - The activities of a trusted user can be trusted.
  - An activity with high trust score can certify the trustiness of its creator.

---

**Algorithm 2:** Coupling random walk

---

```
1 The walker starts with node  $v_i$  with probability  $s_i$ .
2 repeat
3   if The walker is at a user node  $v_i \in \mathcal{V}$  then
4     With probability  $\lambda_i$  it walks one step on the
      friendship graph according to the
      WalkOnFriendGraph algorithm.
5     With probability  $(1 - \lambda_i)$ , it walks  $2k + 1$  steps
      on the user-activity graph according to the
      WalkOnUserActivityGraph algorithm.
6   if The walker is at an activity node  $a_i \in \mathcal{A}$  then
7     With probability  $\lambda_{|\mathcal{V}|+i}$ , the walker walks  $n$  steps
      on the activity-following graph according to the
      WalkOnActivityFollowingGraph algorithm.
8     With probability  $(1 - \lambda_{|\mathcal{V}|+i})$  it takes  $2k + 1$ 
      steps on the user-activity graph according to the
      WalkOnUserActivityGraph algorithm.
9 until converge
```

---

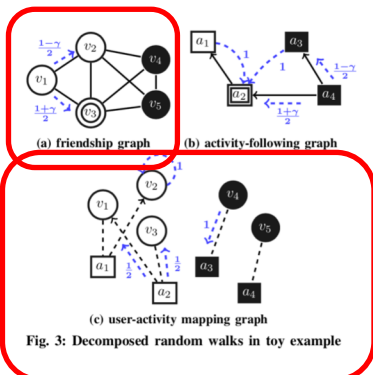


Fig. 3: Decomposed random walks in toy example

# Contributions



Realistic Attack  
Model

- The Social-Activity Network Model



Efficient Detection  
Algorithm

- Sybil\_SAN



Theoretical  
Analysis

- Convergence & Sensitivity analysis



Extensive  
Experiments

- Synthetic & real datasets from Twitter

# Theoretical Analysis: Convergence

## Theorem 1:

**Sufficient conditions** to guarantee Sybil\_SAN **converge** to **unique** trust score

- friendship graph  $\mathcal{G}$  is connected
- $0 < \gamma < 1$
- $0 < \lambda_i < 1, \forall i$  satisfying  $i > |\mathcal{V}|$  or  $v_i$  has activities

## Theorem 2:

Suppose  $\|s^{t+1} - s\| \leq \epsilon$  in Sybil\_SAN is measured by 1-norm, then Sybil\_SAN **stops in at most**  $1 + \frac{1}{v} \ln\left(\frac{4}{\epsilon s_{min}^*}\right)$  rounds, where

- $s_{min}^* = \min_i s_i^*$
- $v$ : spectral gap of the Markov chain

# Theoretical Analysis: Sensitivity analysis

Trust score under graph  
without attack edges

Trust score estimated  
by Sybil\_SAN

**Theorem 3:**

$$\frac{||\mathbf{s}^* - \tilde{\mathbf{s}}^*||}{||\tilde{\mathbf{s}}^*||} \leq \epsilon_{sd}$$

where  $\epsilon_{sd}$  is defined as :

$$\epsilon_{sd} \triangleq ||[(\mathbf{P}_{cr} - \mathbf{I})(\mathbf{P}_{cr}^T - \mathbf{I}) + \mathbf{e}^T \mathbf{e}]^{-1} \times (\mathbf{E}(\mathbf{P}_{cr} - \mathbf{I} + \mathbf{E}^T) + (\mathbf{P}_{cr} - \mathbf{I})\mathbf{E}^T)||$$

# Contributions



Realistic Attack  
Model

- The Social-Activity Network Model



Efficient Detection  
Algorithm

- Sybil\_SAN



Theoretical  
Analysis

- Convergence & Sensitivity analysis



Extensive  
Experiments

- Synthetic & real datasets from Twitter

# Synthetic datasets

- **Honest region:**

- a public Twitter dataset([Weng et.al 2013])

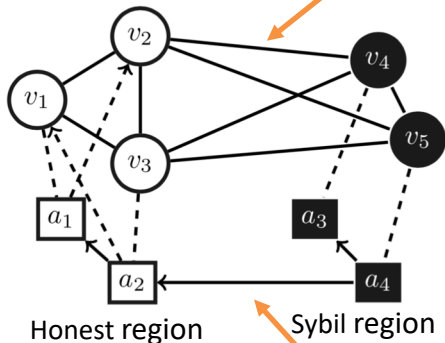
# of nodes	# of edges	# of activities
543,785	28,397,413	21,426,709

- Two type of activities:
    - user  $v_i$  retweets user  $v_j$ 's tweets
    - user  $v_i$  mentions user  $v_j$
- **Sybil region** ( $N_S, M$ ):
  - $M$  disconnected clusters, all together  $N_S$  sybils.
  - the Preferential Attachment (PA) model
- **Attack** ( $N_A, \alpha, \beta$ )

# A More Realistic Attack Model

## Friendship attacks ( $N_A$ ):

- Property 1.  $N_A$  can take any value in  $\{0, 1, \dots, |V_h| \times |V_s|\}$ .



## Incoming interactions attacks ( $\alpha W_h$ ):

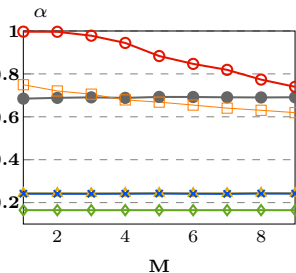
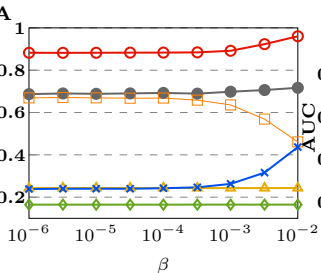
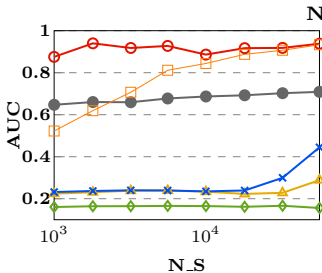
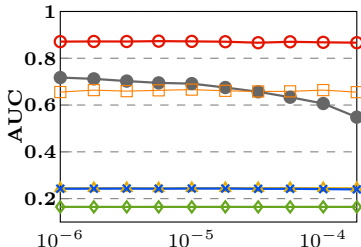
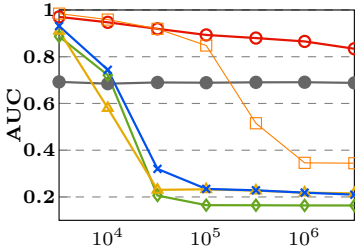
- Activities initiated from honest users to sybils
- $\alpha \approx 10^{-5}$  ([TON'2016]).
- $W_h$ : # of activities among honest users

## Outgoing interaction attacks ( $\beta W_h$ ):

- Property 2.  $\beta$  can be arbitrary large.

# Experimental results

—○— SAN —□— SWALK —◇— SSCAR —▲— SR-U —\*— SR-W —●— INTER



Sybil\_SAN outperforms compared algorithms under various attacks.



## Experiments on real dataset

- A crawled subnetwork from Twitter starting from public 991 sybils.
  - Sybils: blocked users
  - Honest users: unblocked users

# of honest users	# of sybils	# of edges	# of interactions
409, 694	40, 548	222,944,310	102,693,769

## Experiments on real dataset

- A crawled subnetwork from Twitter starting from public 991 sybils.
  - Sybils: blocked users
  - Honest users: unblocked users (**noisy**).
- **Results**

	Inter	SR_W	SR_U	SScar	SWalk	SAN
AUC	0.62	0.48	0.52	0.15	0.44	<b>0.73</b>
<b>Improved ratio</b>	<b>17.7%</b>	<b>52.1%</b>	<b>40.4%</b>	<b>386%</b>	<b>66%</b>	

# Contributions



Realistic Attack  
Model

- The Social-Activity Network Model



Efficient Detection  
Algorithm

- Sybil\_SAN



Theoretical  
Analysis

- Convergence & Sensitivity analysis



Extensive  
Experiments

- Synthetic & real datasets from Twitter



thank you!