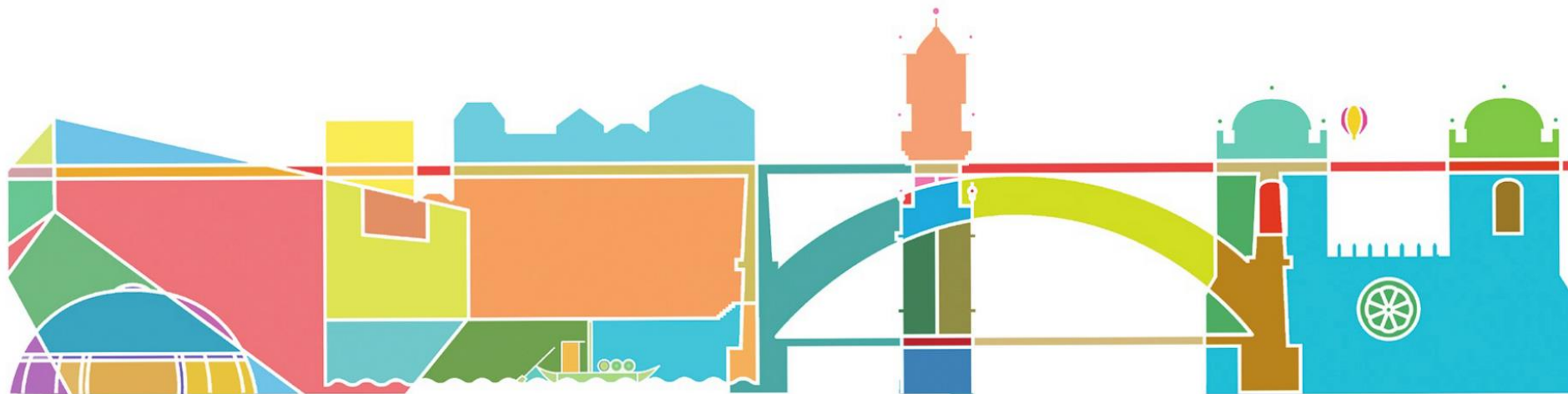


# Two Step Graph-based Semi-supervised Learning for Online Auction Fraud Detection

---

Phiradet Bangcharoensap<sup>1</sup>, Hayato Kobayashi<sup>2</sup>, Nobuyuki Shimizu<sup>2</sup>,  
Satoshi Yamauchi<sup>2</sup>, and Tsuyoshi Murata<sup>1</sup>

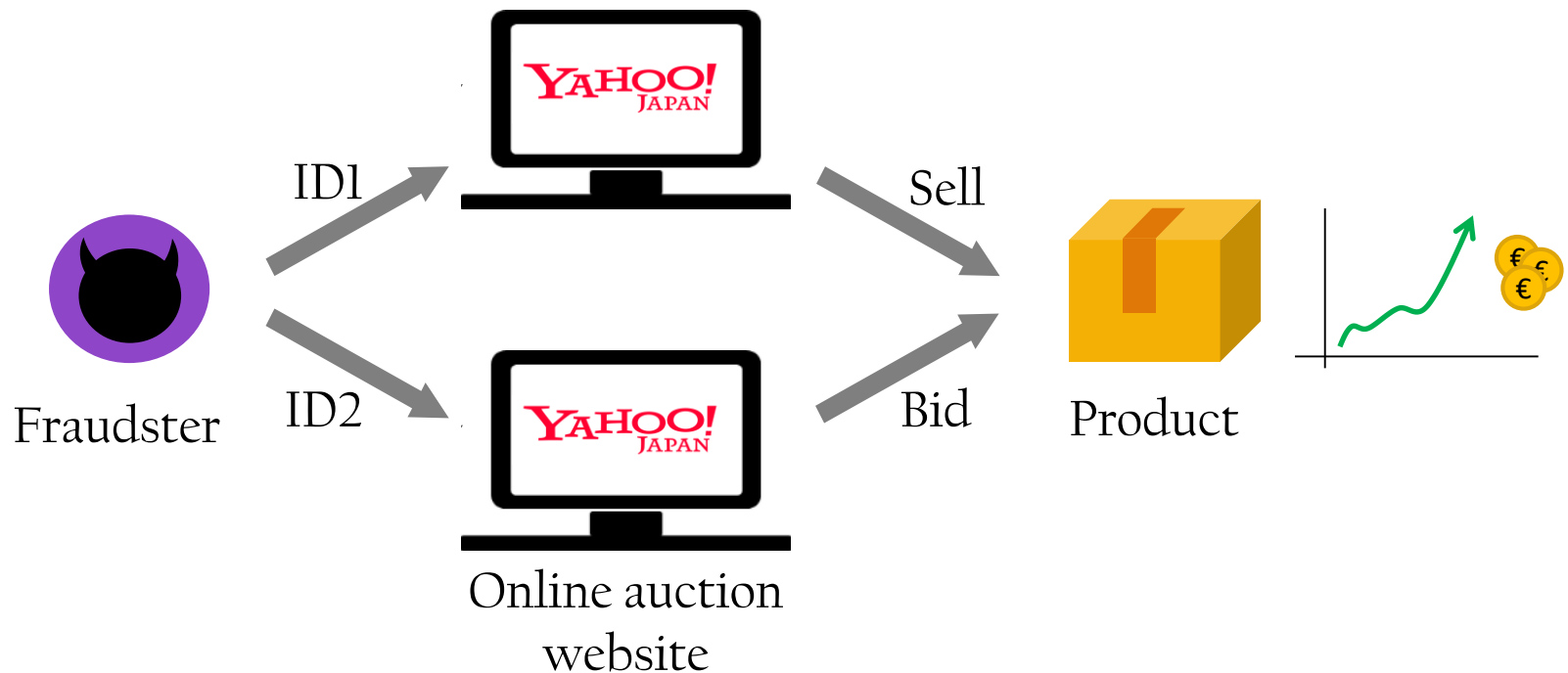
<sup>1</sup>Tokyo Institute of Technology, <sup>2</sup>Yahoo Japan Corporation



# Definition of Fraudster

## Competitive Shilling

auction users who bid on their product, as other user IDs, in order to drive up the final price.



# Key Ideas

---



## Fraudsters

- frequently participate in auctions hosted by fraudulent sellers working in a same group



## Innocents

- rarely interact with fraudsters
  - frequently interact with famous sellers
- or
- uniformly interact with various sellers

# Key Ideas

---



## Fraudsters

- frequently participate in auctions hosted by fraudulent sellers working in a same group

**H**omophily



## Innocents

- rarely interact with fraudsters
- frequently interact with famous sellers

*or*

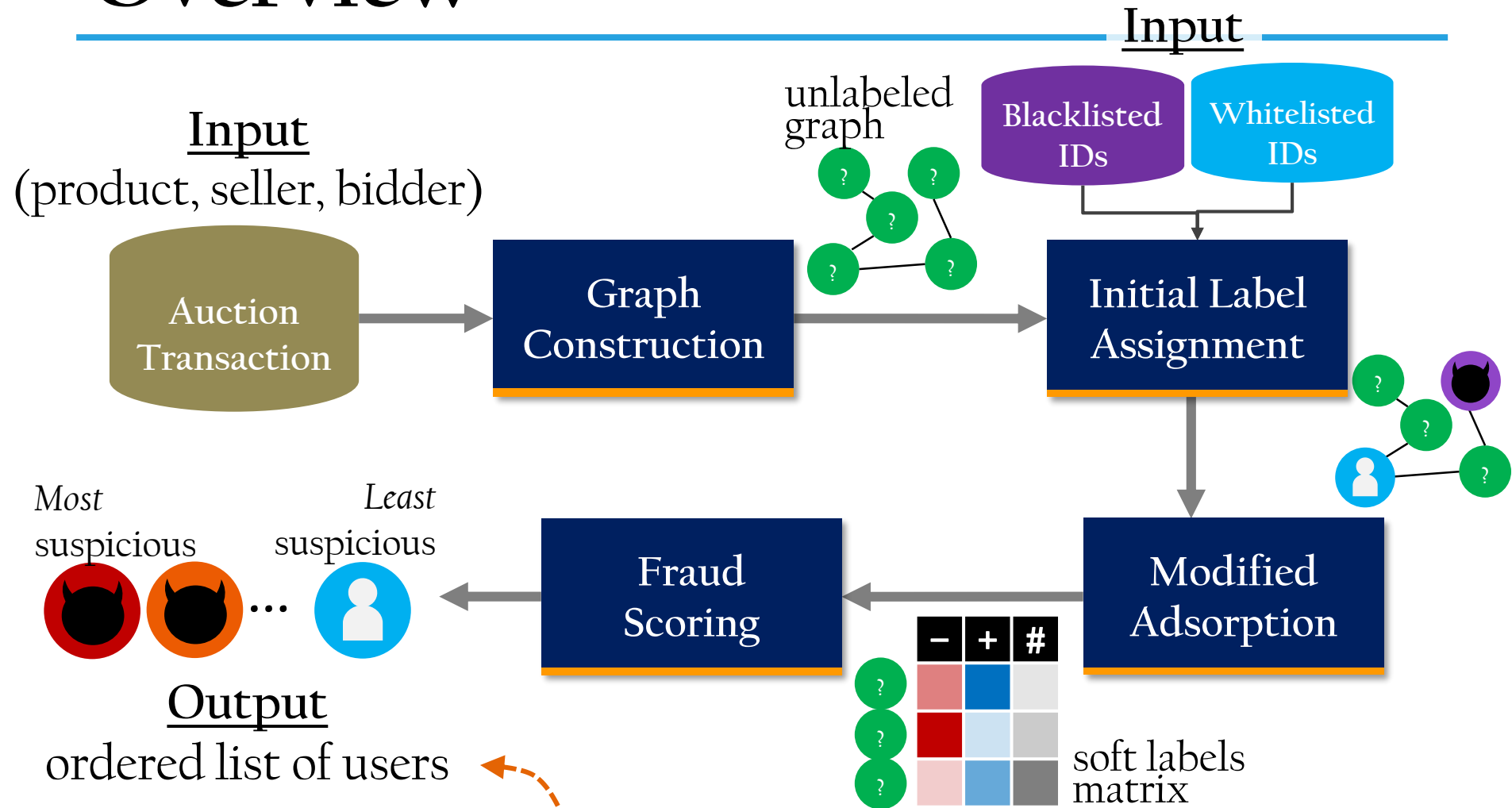
**U**niformly interact with various sellers

# Contributions

---

1. Novel application of Modified Adsoption (MAD) [Talukdar & Crammer, ECMLPKDD'09]
  - Have been previously used in NLP
  - **H**omophily: smoothness constraint
  - **U**niformity of innocents: dummy label
2. Incorporate weighted degree centrality
  - Fraudsters tend to form very strong ties.
  - Help us to yield better results

# Overview

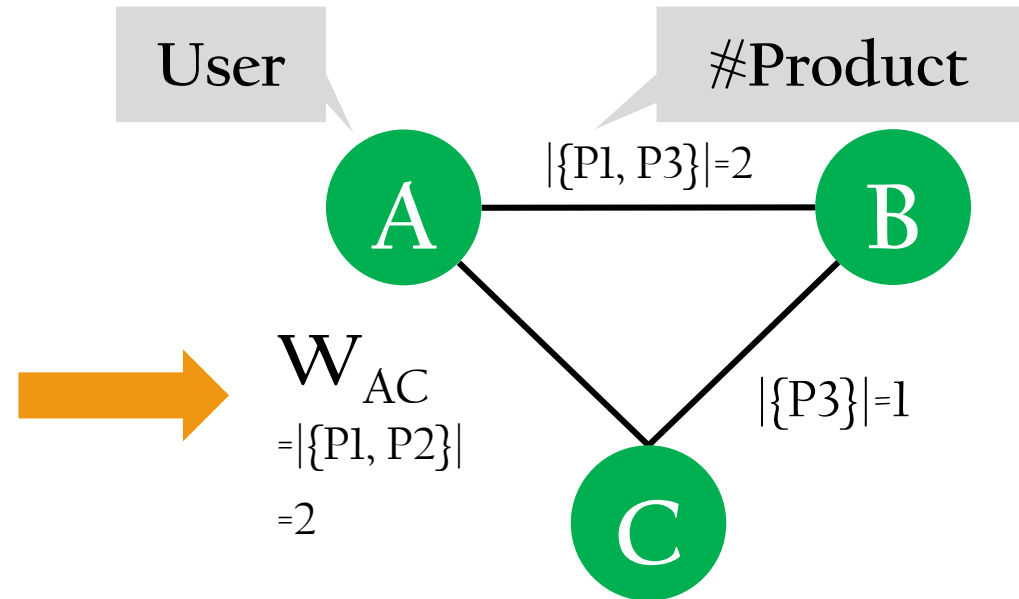


**Objective:** Fraudsters working in the same collusion with the blacklisted users are ranked at the top.

# Graph Construction

Product	Seller	Bidder
P1	A	B
P1	A	C
P2	A	C
P3	B	A
P3	B	C
P3	B	C
P3	B	C

Online auction  
transaction

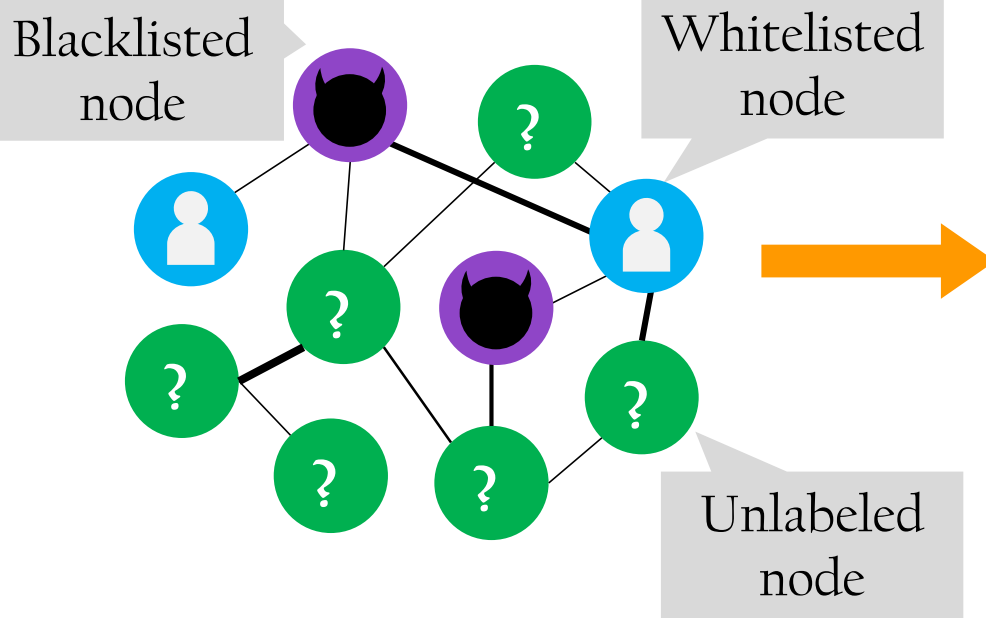


Weighted undirected  
graph

# Graph-based SSL

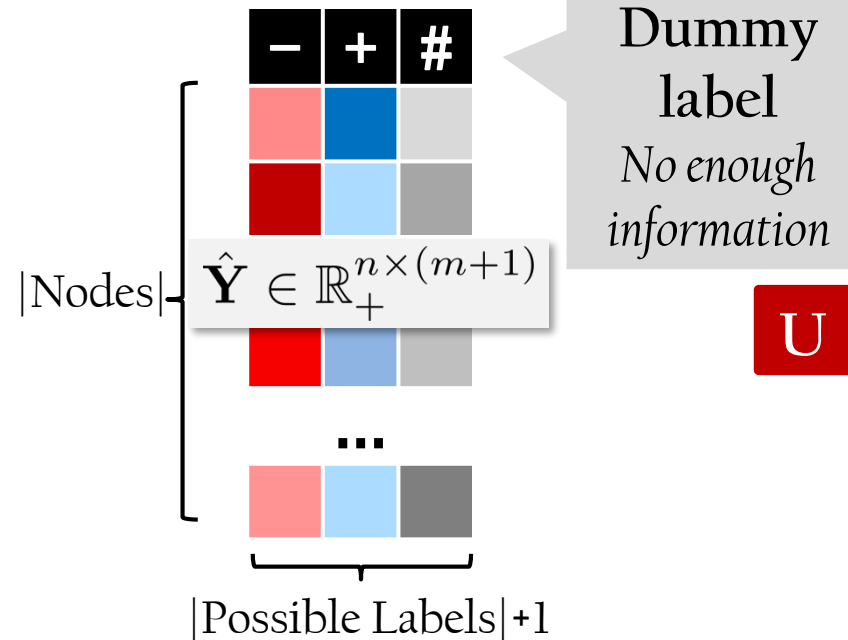
Modified Adsorption (MAD) [Talukdar & Crammer,'09] is used.

Input: partially labeled  
weighted undirected graph



Node: instance that want to classify  
Edge: similarity between instances

Output: soft label matrix



Assign a score indicating likelihood of  
being each label (soft labels)



# Dummy Label

- Exceptional case of all other labels

Entropy  
Amount of uncertainty

$$\text{dummy}_v \sim - \sum_{u \in N(v)} \frac{\mathbf{W}_{uv}}{k_w(v)} \log \frac{\mathbf{W}_{uv}}{k_w(v)}$$

Neighbors of vertex  $v$       Weighted degree of vertex  $v$

$$k_w(v) = \sum_{u \in N(v)} \mathbf{W}_{uv}$$

**U**

The score of dummy is high when the vertex uniformly interacts with its neighbors.

# Modified Absorption (MAD)

Tradeoff between fitting and smoothness constraints

- **Fitting**: retain initial labels of seed nodes
- **Smoothness**: assign same labels to adjacent nodes

H

Solving the convex optimization problem

$$\min_{\hat{\mathbf{Y}}} \sum_{l \in \mathcal{L}} \left[ \underbrace{\mu_1 (\mathbf{Y}_l - \hat{\mathbf{Y}}_l)^\top \mathbf{S} (\mathbf{Y}_l - \hat{\mathbf{Y}}_l)}_{\text{Fitting}} + \underbrace{\mu_2 \hat{\mathbf{Y}}_l^\top \mathbf{L} \hat{\mathbf{Y}}_l}_{\text{Smoothness}} + \underbrace{\mu_3 \left\| \hat{\mathbf{Y}}_l - \mathbf{R}_l \right\|^2}_{\text{Regularization}} \right]$$

where  $\hat{\mathbf{Y}}$  is a matrix storing scores of labels (soft label matrix)

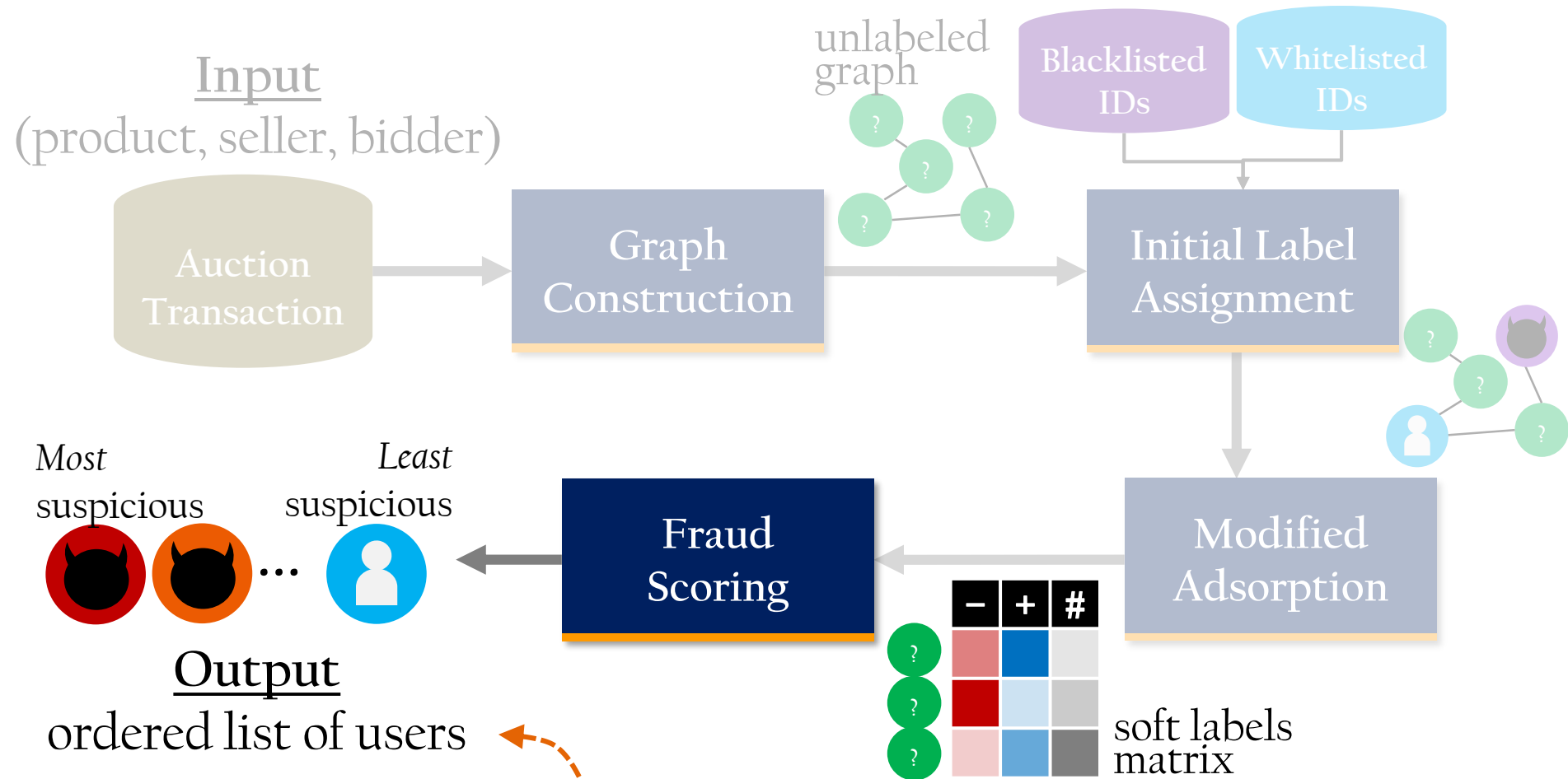
$\mathbf{Y}$  stores seed information

$\mathbf{S}$  indicates positions of seed vertices

$\mathbf{L}$  is the Laplacian matrix

$\mathbf{R}$  encodes scores of the dummy label and  $L^2$  regularization.

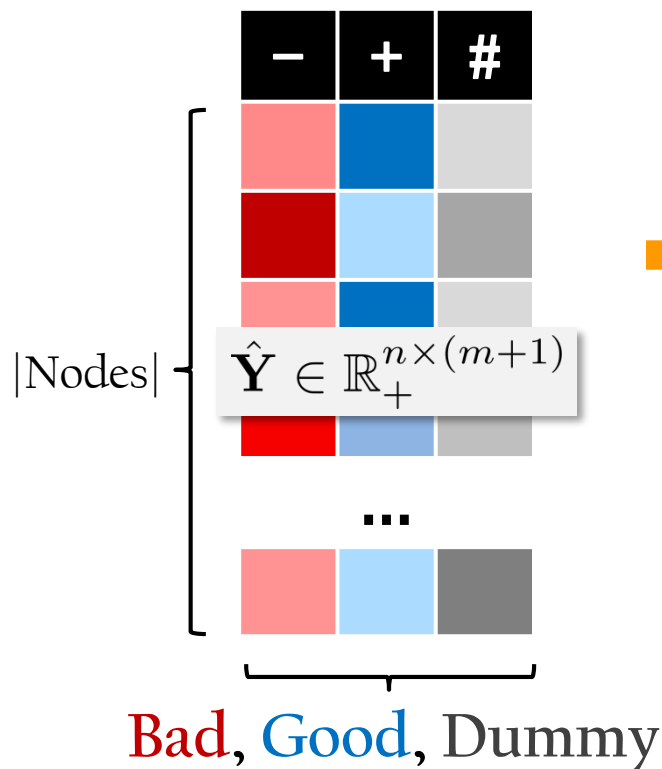
# Overview (2)



**Objective:** Fraudsters working in the same collusion with the blacklisted users are ranked at the top.

# Fraud Scoring

Input: soft label matrix




Output: fraud score of nodes

$$\varphi(v, \hat{\mathbf{Y}}) = \frac{\hat{\mathbf{Y}}_{v1}}{\sum_{l=1}^{m+1} \hat{\mathbf{Y}}_{vl}} \quad \text{MAD}$$

The ratio of **Bad**'s score to total scores

# Contributions

---

1. Novel application of Modified Adsorption (MAD) [Talukdar & Crammer, ECMLPKDD'09]
  - **H**omophily: smoothness constraint
  - **U**niform interaction of innocents: dummy label
2. Incorporate weighted degree centrality (WDC)
  -  Fraudsters form very strong ties.

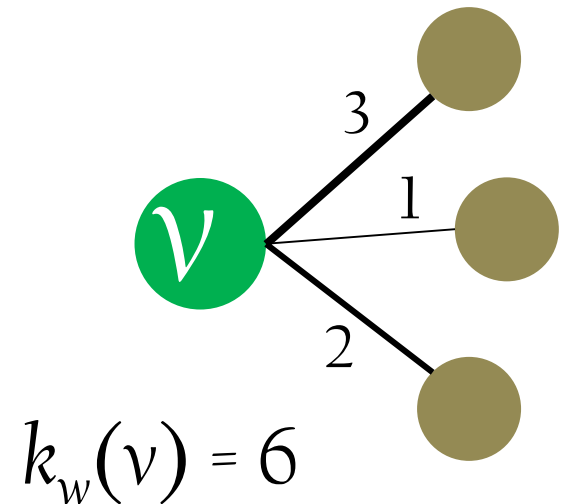
# Weighted Degree Centrality (WDC)

Weighted degree centrality of vertex  $v$  is the total weights of edges originating from  $v$

$$k_w(v) = \sum_{u \in N(v)} w_{uv}$$

Neighbors of  $v$

Weight of an edge  $(u,v)$



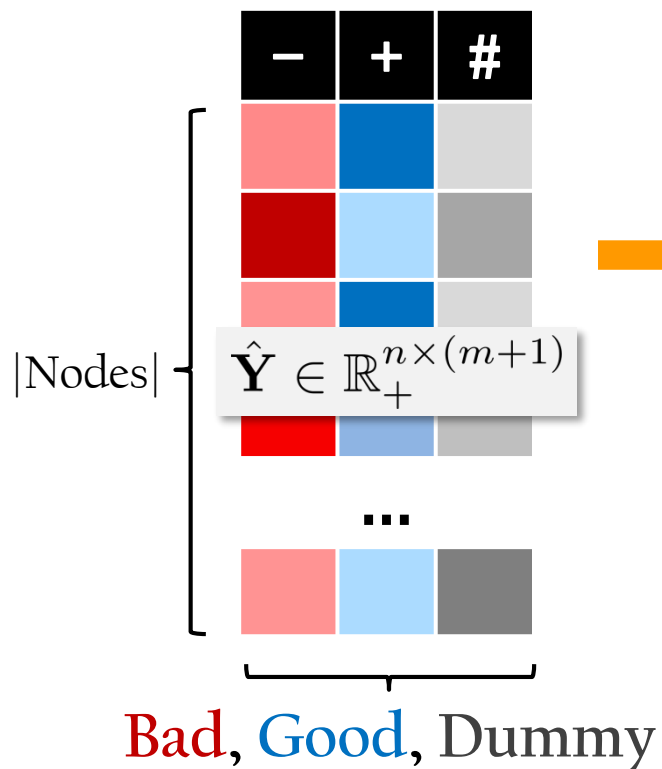
Fraudsters tend to have higher *weighted degree centralities* because of stronger ties.

H

# Fraud Scoring + WDC

Input: soft label matrix

Output: fraud score of nodes



$$\rho(v, \hat{\mathbf{Y}}) = \varphi(v, \hat{\mathbf{Y}})$$

2-STEP

$$+ \frac{\gamma}{|N(v)|} \sum_{u \in N(v)} \mathbf{W}_{uv} \varphi(u, \hat{\mathbf{Y}})$$

Neighbors of  
vertex  $v$

Weight of an  
edge  $(u,v)$

$$\varphi(v, \hat{\mathbf{Y}}) = \frac{\hat{\mathbf{Y}}_{v1}}{\sum_{l=1}^{m+1} \hat{\mathbf{Y}}_{vl}}$$

MAD

# Experiments

---

- Questions
  1. Does the dummy label help?
  2. Comparison with unsupervised methods
  3. Comparison with a state-of-the-art Sybil defense method
- Evaluation metric

Used normalized discounted cumulative gain (NDCG) to compare results with the blacklisted users

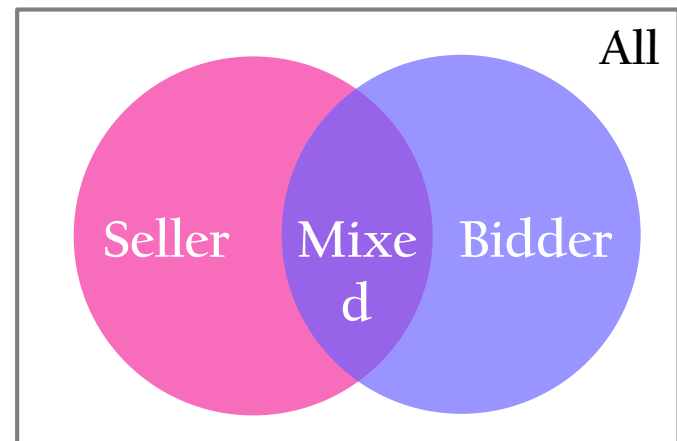
Higher NDCG is better.



# Dataset

---

- Real-world dataset from YAHUOKU<sup>1</sup>
  - The largest online auction site in Japan
  - Operated by Yahoo! Japan
- Auction transaction
  - ≈ 16 million transactions
  - ≈ 2 million users
  - ≈ 550 blacklisted users
  - ≈ 10,000 whitelisted users



<sup>1</sup>[auctions.yahoo.co.jp/](http://auctions.yahoo.co.jp/)

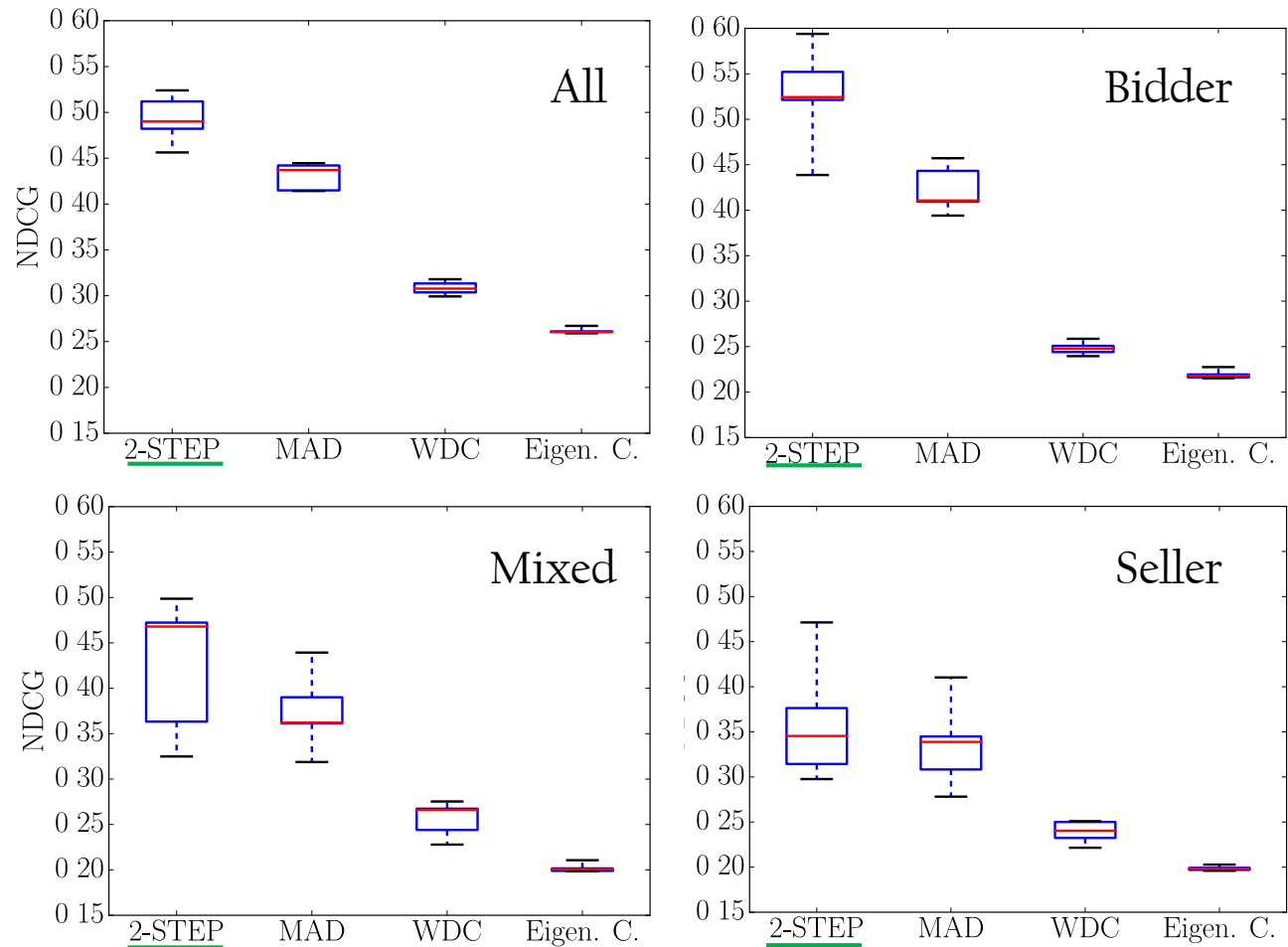
# With VS Without Dummy Label

Node type	with dummy		w/o dummy	
	$\langle \text{NDCG} \rangle$	SD	$\langle \text{NDCG} \rangle$	SD
All	<u>0.431</u>	0.015	0.406	0.019
Bidder	<u>0.423</u>	0.026	0.397	0.035
Seller	<u>0.336</u>	0.049	0.284	0.029
Mixed	<u>0.374</u>	0.044	0.319	0.024

- Dummy label has a true advantage.
- Support the key idea that innocents tend to interact with neighbors uniformly


 U

# Proposed VS Unsupervised



Compare with

- 1) Weighted degree centrality (WDC)
- 2) Eigenvector centrality (Eigen. C.)

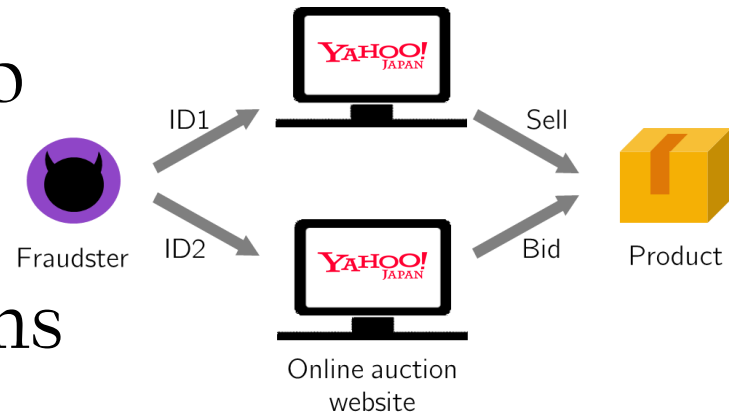
2-STEP method  
outperforms MAD.

Unsupervised methods  
yield poor results.

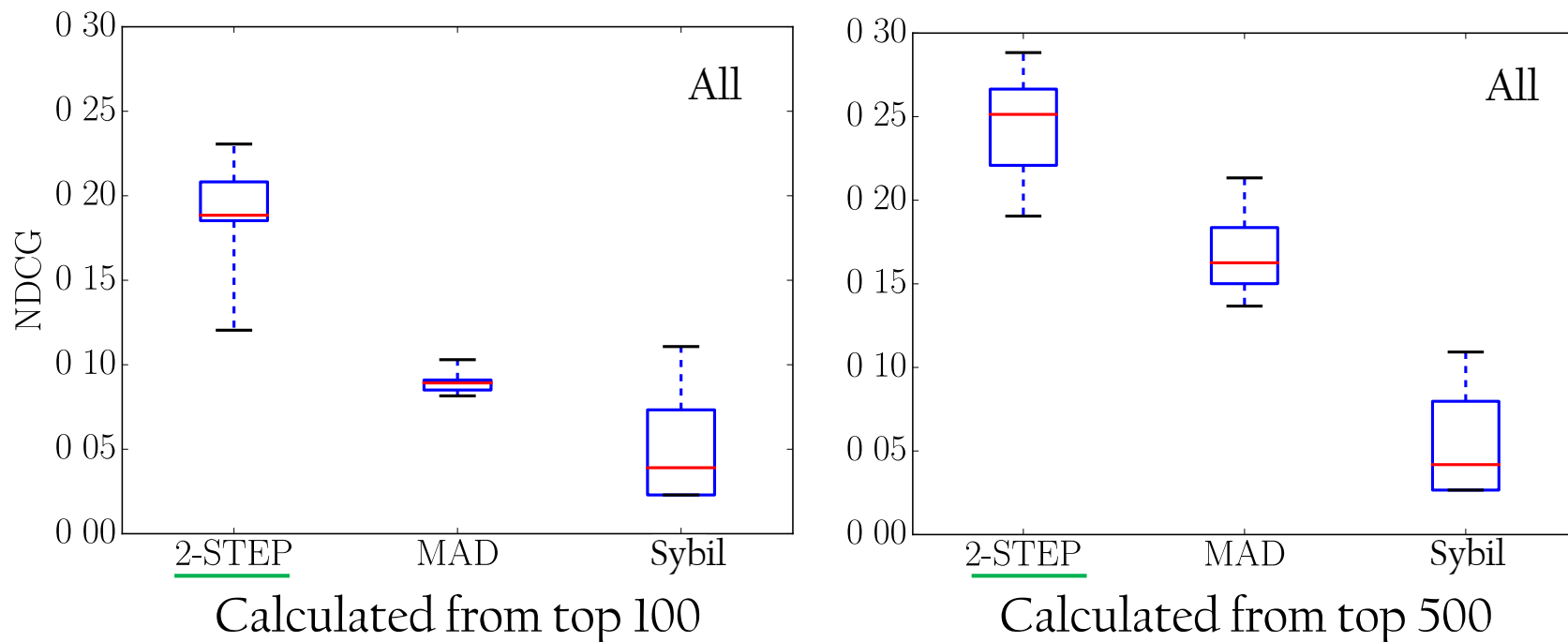
Fraudulent sellers are  
more difficult.

# Sybil Defense Method

- Sybil: malicious attackers who
  - create multiple identities
  - influence working of systems
- Shill bidders are one type of Sybil
- We compared our method with a state-of-the-art Sybil defense method [Viswanath et al., SIGCOMM'10]
  - On basis of community detection



# Proposed VS Sybil



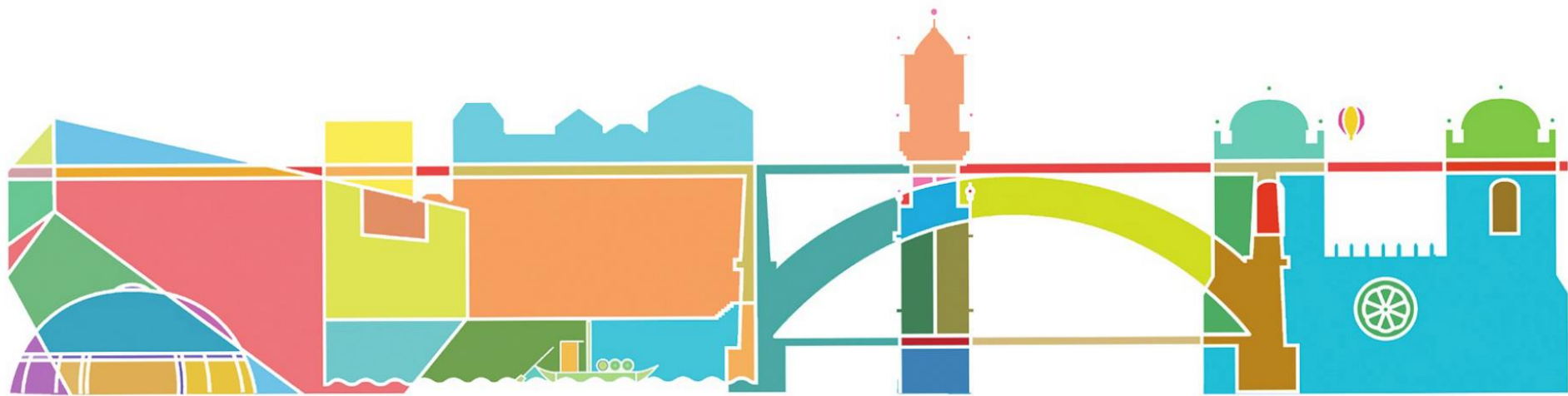
- Our method outperforms the state-of-the-art Sybil defense method.
- Fraudsters and innocents may not form well-established communities.

# Conclusion

---

- Proposed an online auction fraud detection approach
- Motivated by two main ideas
  - **U**niformity of innocents
  - **H**omophily
    - Fraudsters tend to have higher WDCs.
- Incorporated WDC to the method
- Our extended method yields better results.

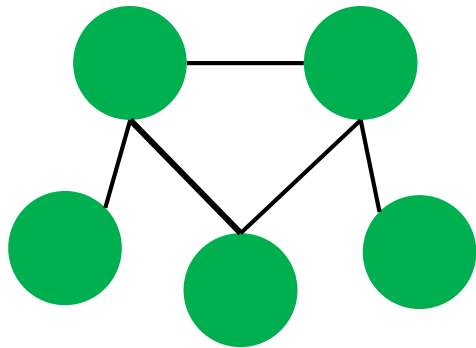
Thank you



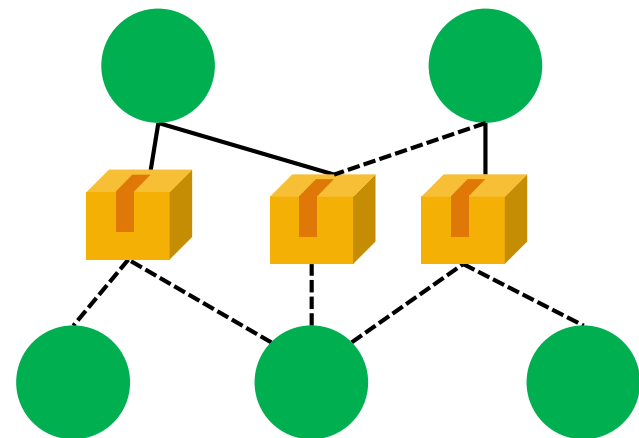
# Future Works

---

- Study limitation of the method
- Incorporate other heuristics
  - Bidding strategy
  - Value of products
- Extend the method to heterogeneous network



Homogeneous network



Heterogeneous network



# Scalability

---

- The optimization process of MAD can be parallelized in MapReduce framework.
  - Map: sends its current label to neighbors
  - Reduce: update its label information
- Hadoop-based implementation is available.
  - Junto Label Propagation Toolkit:  
<https://github.com/parthatalukdar/junto/>