

Graph Neural Networks based Fraud Detection: from Research to Application



Yingtong Dou
ydou5@uic.edu



Kay Liu
zliu234@uic.edu

University of Illinois Chicago



@Wells Fargo 12/06/2022

Outline

- **Background: fraud detection and graph neural networks.**
- **Supervised methods and DGFraud.**
- **Unsupervised methods, PyGOD, and benchmark.**
- **Applied graph ML case study.**
- **Challenges, solutions, insights, guideline, and resources.**

Outline

- **Background: fraud detection and graph neural networks.**
- Supervised methods and DGFraud.
- Unsupervised methods, PyGOD, and benchmark.
- Applied graph ML case study.
- Challenges, solutions, insights, guideline, and resources.

Anomaly vs. Fraud

- **Fraud definition according to U.S. Law:**
 - a misrepresentation of a fact, made from one person to another, with knowledge of its falsity and for the purpose of inducing the other to act.
- **Anomaly definition^[1]**
 - An anomaly is a data point that is significantly different from rest of the data.
- **Fraud vs. Anomaly**
 - Not all frauds are anomalies.
 - Not all anomalies are frauds.



Machine Learning in Fraud Detection



Content-based Detectors

Average content similarity

The ratio of exclamation sentences

Description length based on bigrams



Behavior-based Detectors

The frequency of review

Deceptive review count previous week

Max. number of reviews posted in a day



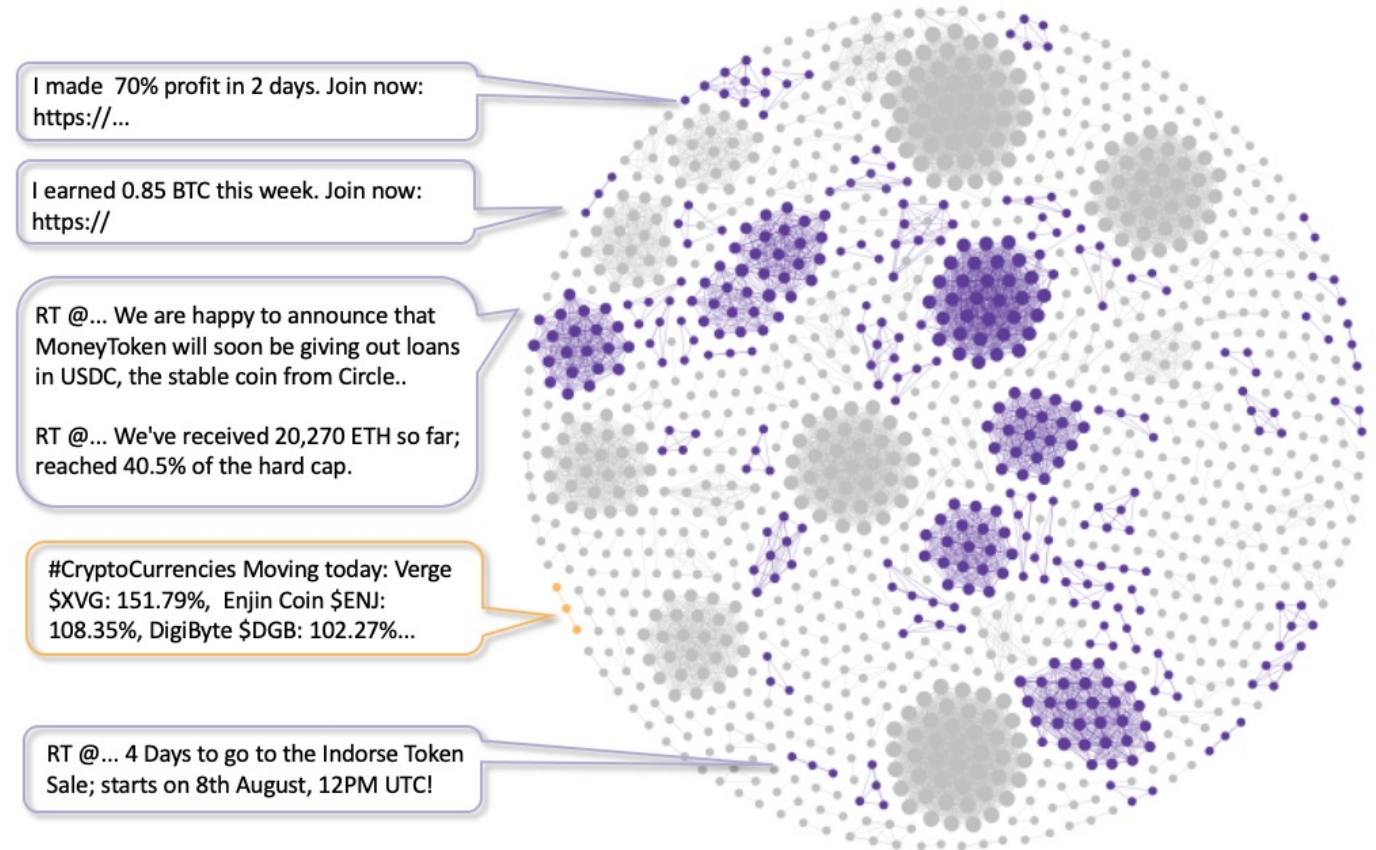
Graph-based Detectors

Discuss in this talk!

Graph-based Fraud Detection

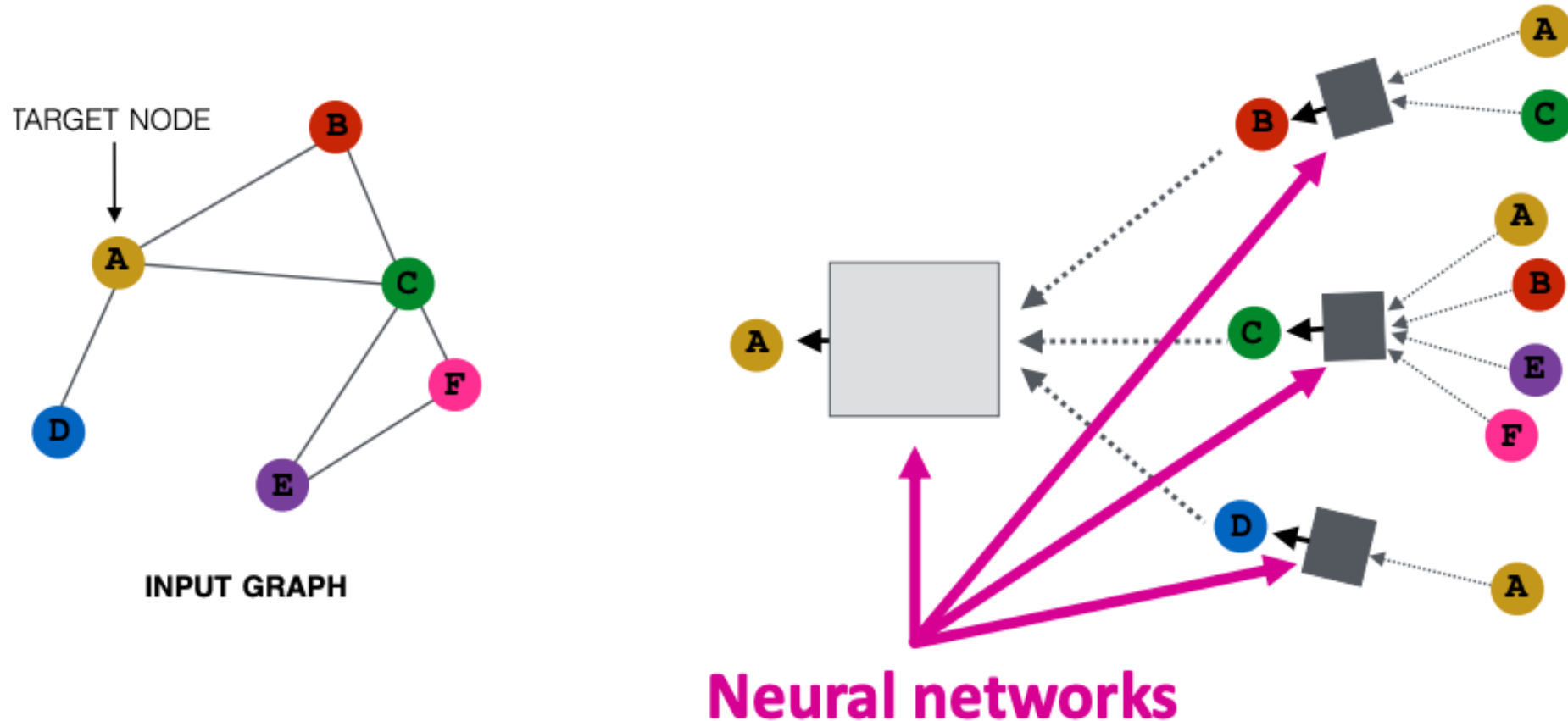


Bot Account on Twitter



Coordinated Accounts on Social Network^[1]

Graph Neural Networks



Key idea: the connected nodes are similar (homophily assumption)

GNN Use Cases in Industry

- [Pinterest](#), [Snapchat](#)
 - Recommender systems
- [Amazon & United Airlines](#)
 - Information extraction
- [AstraZeneca](#)
 - Molecular Generation

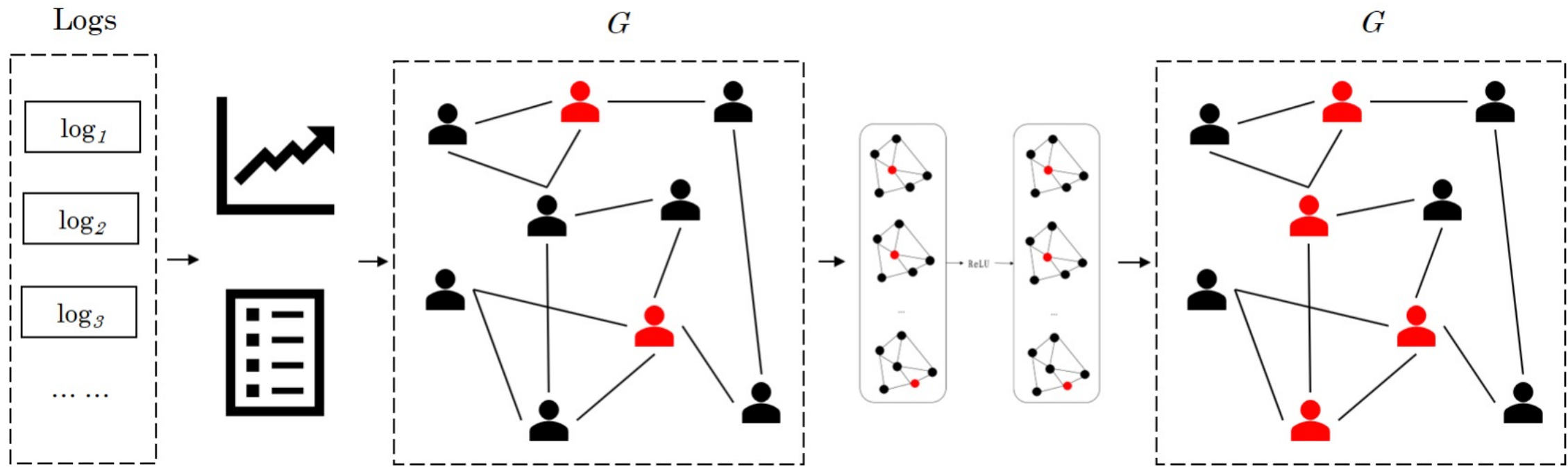
GNN for Fraud Detection

- Insurance Fraud
- Loan Defaulter
- Money Laundering
- Malicious Account
- Transaction Fraud
- Cash-out User
- Bitcoin Fraud

Outline

- Background: fraud detection and graph neural networks.
- **Supervised methods and DGFraud.**
- Unsupervised methods, PyGOD, and benchmark.
- Applied graph ML case study.
- Challenges, solutions, insights, guideline, and resources.

Supervised GNN



(1) Graph Construction.

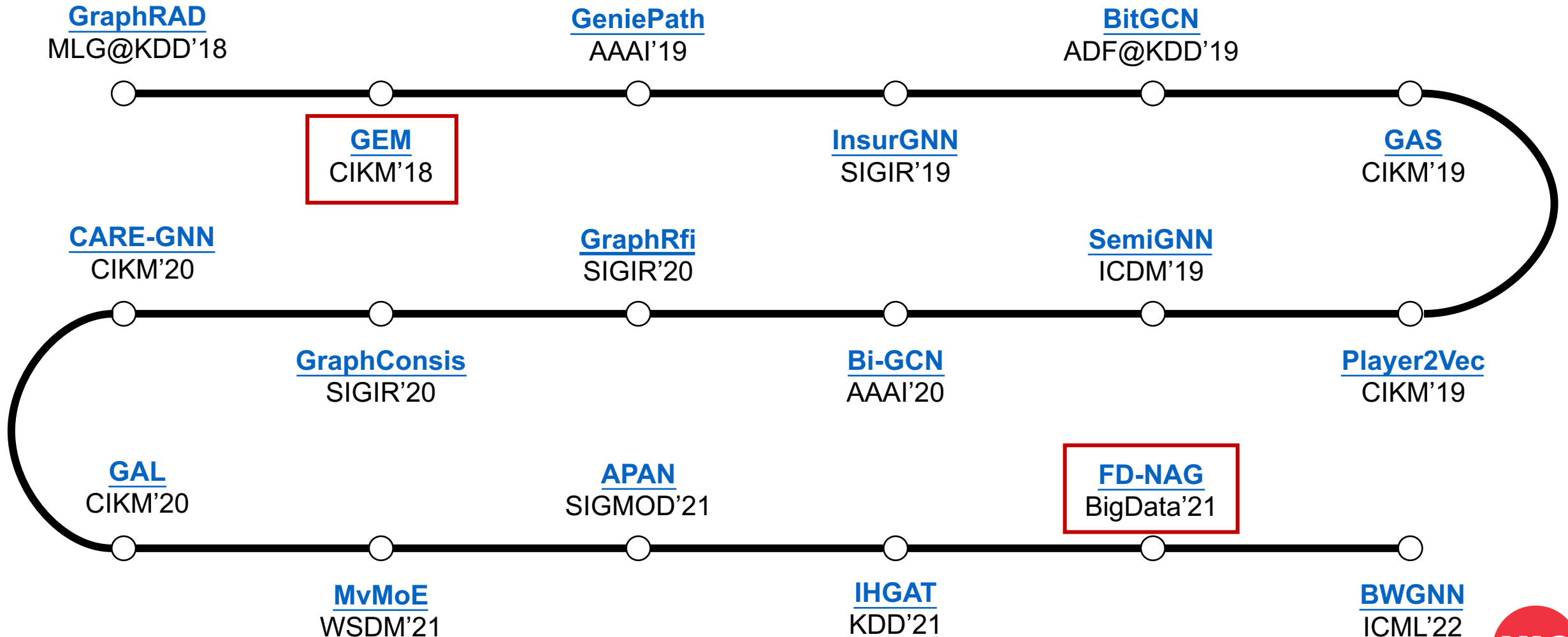
(2) Training GNN on the Graph with labeled nodes.

(3) Classifying Unlabeled Nodes.

Background

Supervised

A Short History of GNN-based Fraud Detection

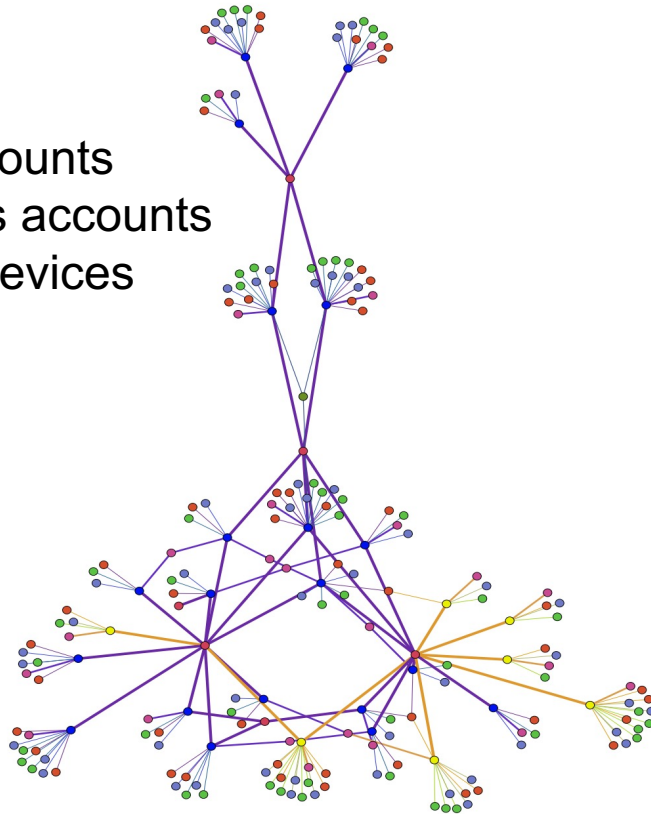


GEM (CIKM'18)

Blue: normal accounts

Yellow: malicious accounts

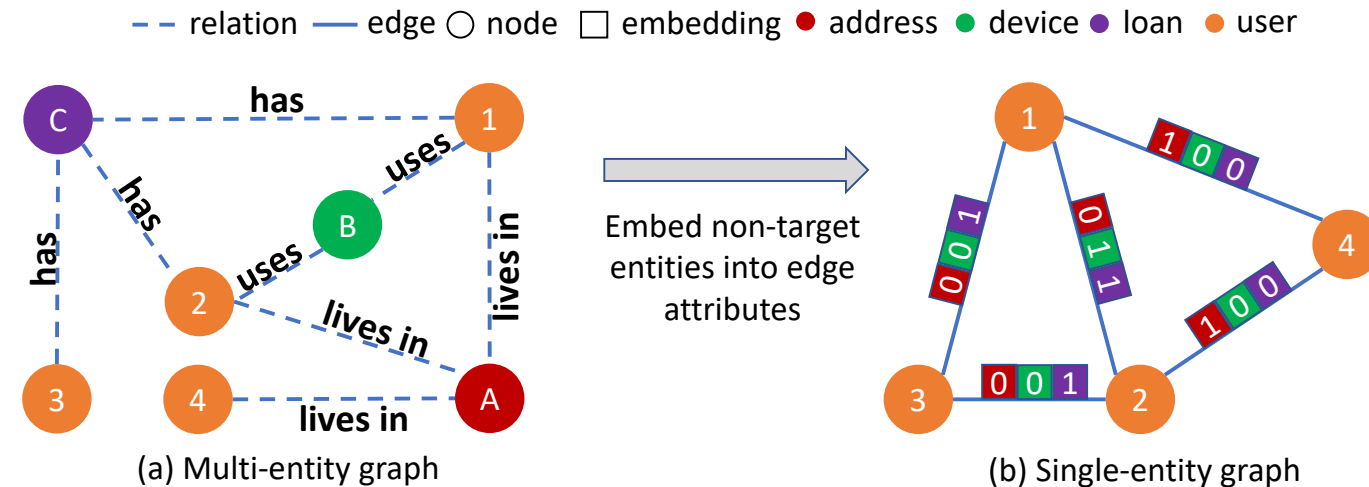
Other: different devices



**Account-Device
Heterogeneous Graph**

- Task: malicious accounts detection in mobile payment service (Alipay).
- **The first paper leveraging the heterogeneous graphs for fraud detection.**
- Device types include UMID, MAC address, IMSI, APDID (Alipay Fingerprint).
- Code is [available](#).

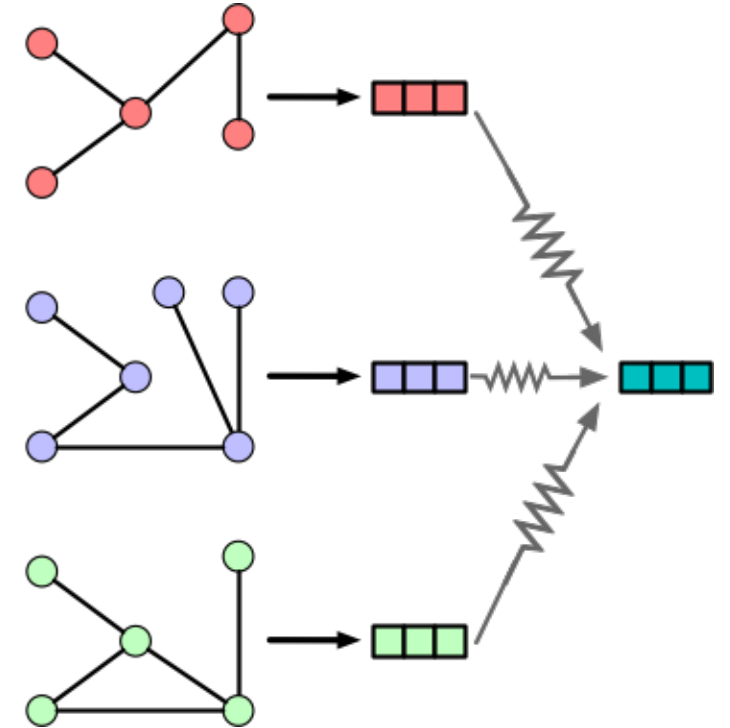
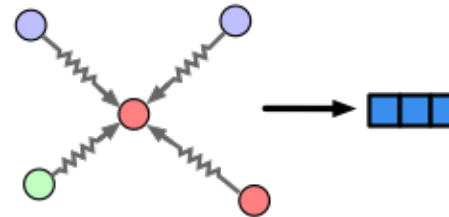
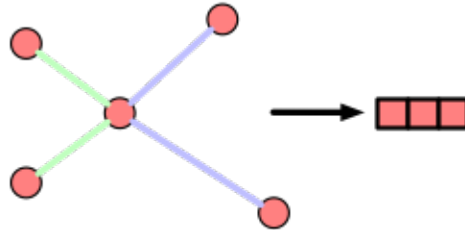
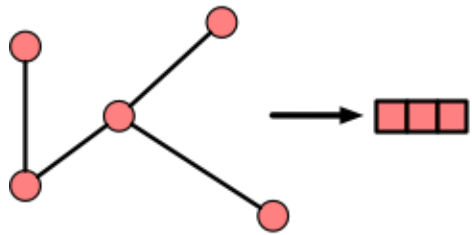
FD-NAG (BigData'21)



Transferring a heterogeneous non-attributed graph to an edge-attributed homogeneous graph

- Task: fraudsters detection in ride sharing services.
- Designing node and edge features for non-attributed graphs.
- Empirically verified the effectiveness of **contrastive learning** in fraud detection.

Graph Schema



Homogeneous

BitGCN
FdGars
GeniePath
FD-NAG

Multi-relation

GraphConsis
CARE-GNN
PC-GNN

Heterogeneous

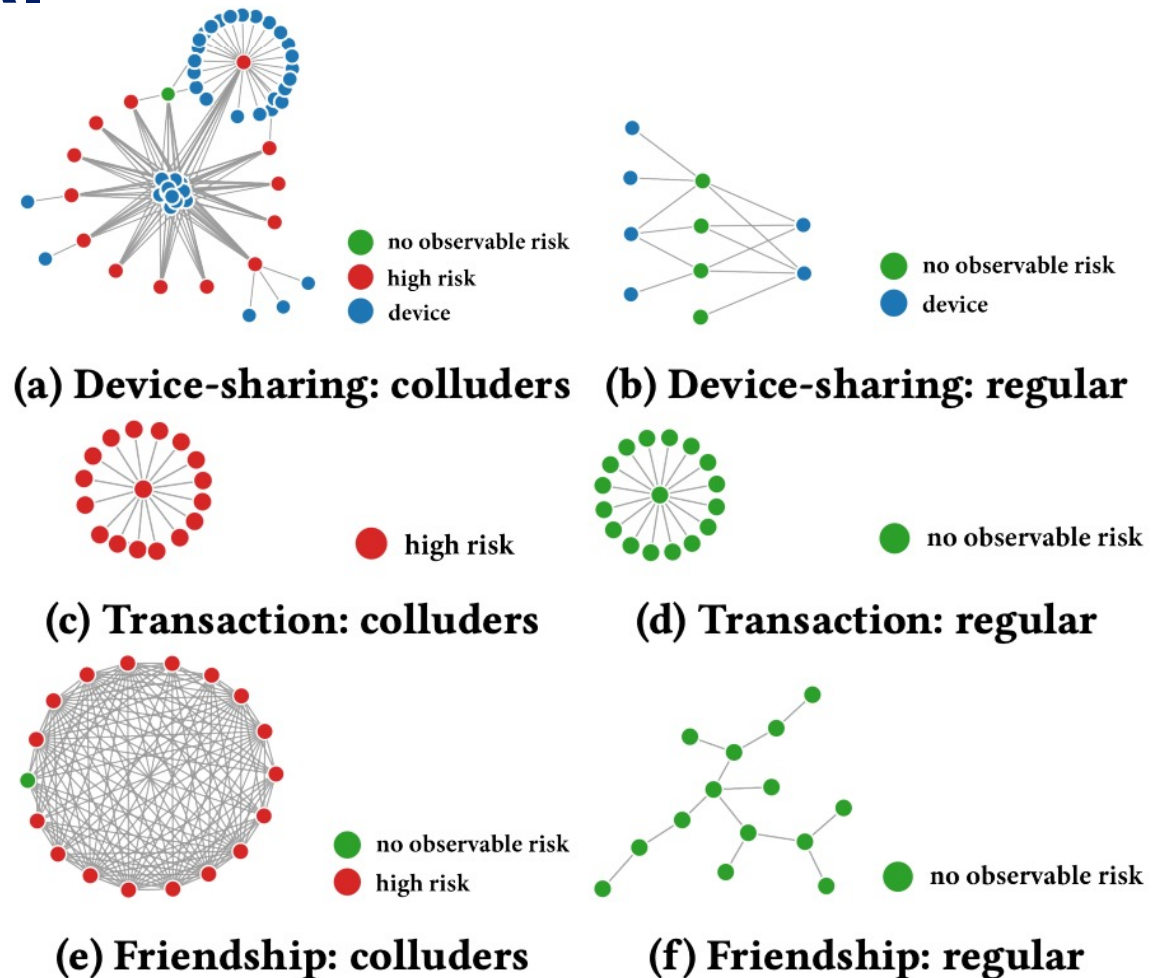
GAS
mHGNN
IHGAT

Hierarchical

GEM
SemiGNN
Player2Vec
AA-HGNN

Graph Schema is Crucial

- Task: finding fraud colluders on an online insurance platform.
- The suspicious signal can only be visible under certain graph schemas.
- Graph schema design is the key step for applied graph machine learning.



DGFraud

- [DGFraud](#) – A Deep Graph-based Toolbox for Fraud Detection.



build passing license Apache-2.0 release v0.1.0 PRs welcome



build passing tensorflow 2.X python 3.6 | 3.7 | 3.8 | 3.9 PRs welcome release v0.1.0

Combined 600+ Stars on GitHub

Model	Application	Graph Type	Base Model
SemiGNN	Financial Fraud	Heterogeneous	GAT, LINE, DeepWalk
Player2Vec	Cyber Criminal	Heterogeneous	GAT, GCN
GAS	Opinion Fraud	Heterogeneous	GCN, GAT
FdGars	Opinion Fraud	Homogeneous	GCN
GeniePath	Financial Fraud	Homogeneous	GAT
GEM	Financial Fraud	Heterogeneous	GCN
GraphSAGE	Opinion Fraud	Homogeneous	GraphSAGE
GraphConsis	Opinion Fraud	Heterogeneous	GraphSAGE
HACUD	Financial Fraud	Heterogeneous	GAT

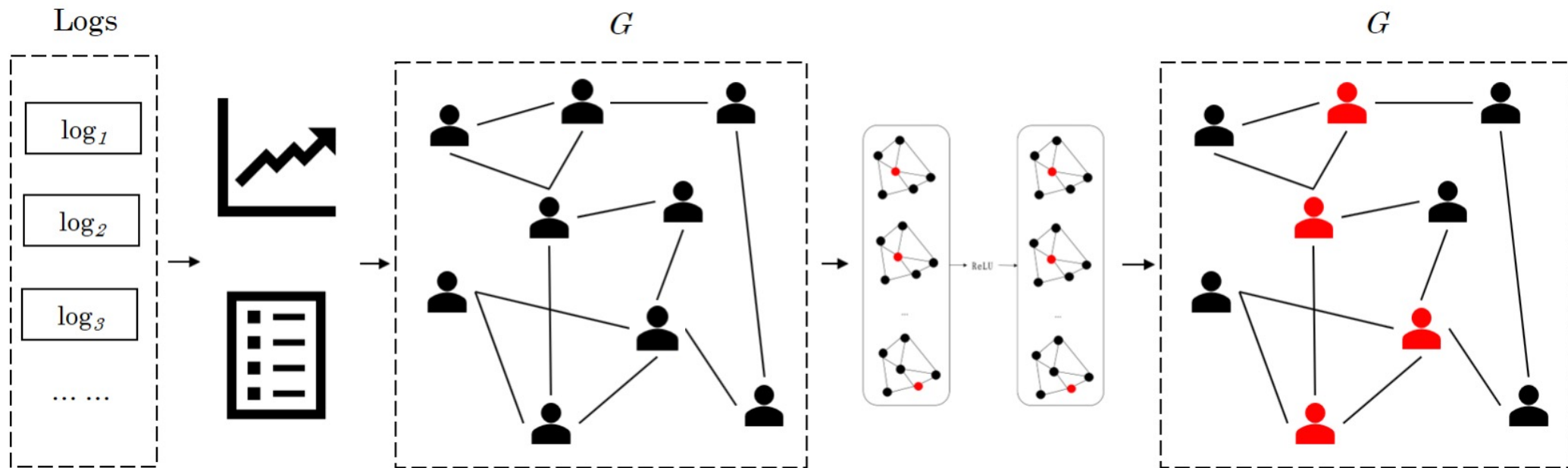
Outline

- Background: fraud detection and graph neural networks.
- Supervised methods and DGFraud.
- **Unsupervised methods, PyGOD, and benchmark.**
- Applied graph ML case study.
- Challenges, solutions, insights, guideline, and resources.

Unsupervised Anomaly Detection with Graphs

- **Label scarcity**
 - Ground truth labels can be expensive, even impossible to obtain.
- **Novelty detection**
 - Unsupervised learning does not rely on existing labeled data.
- **Preprocessing for downstream tasks**
 - E.g., Outlier resistant node classification.

Graph Auto-Encoder (GAE)

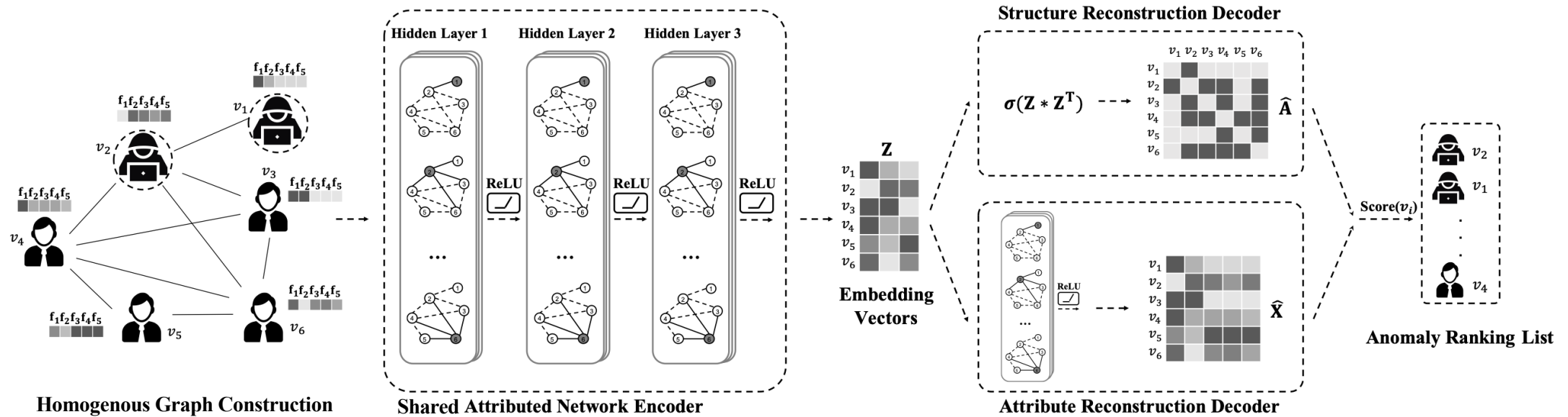


(1) Graph Construction.

(2) Training GAE on the Graph.

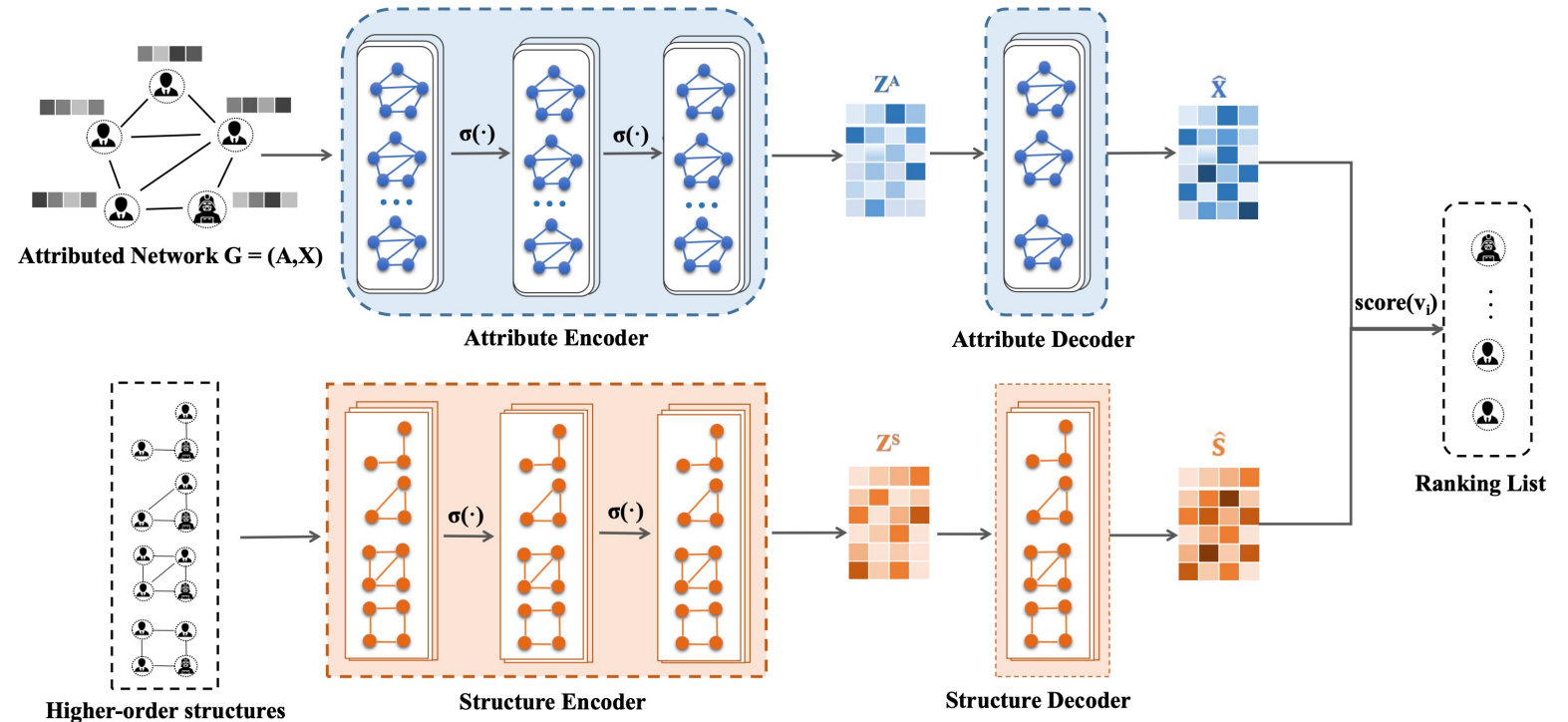
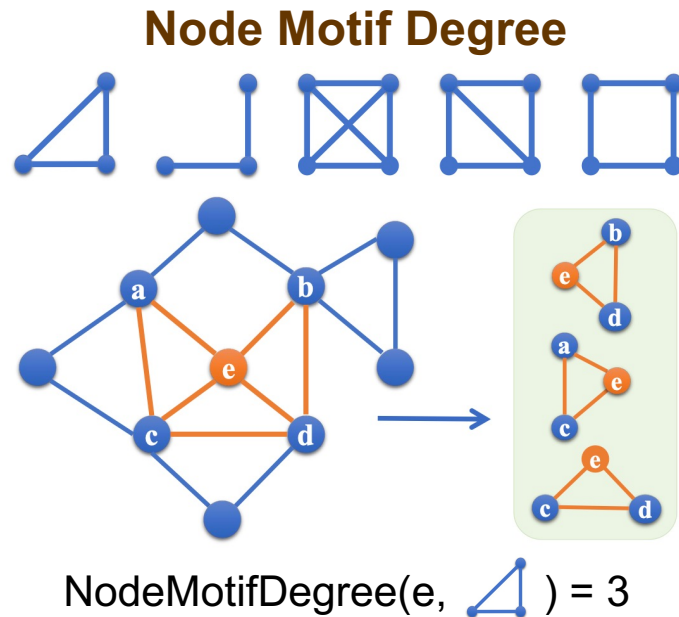
(3) Detecting Outlier Nodes.

DOMINANT (SDM'19)



- The first attempt of **graph auto-encoder** in graph anomaly detection problem.
- Adopted **multi-task learning** framework to jointly detect anomalies from two aspects.
- Using **reconstruction error** of structure and attribute as anomaly score.

GUIDE (Big Data'21)



- Capture higher-order structure information with **node motif degree**.
- Largely improve the scalability for AE, but **huge burden** on node motif degree counting.
- **Imprecise estimate** (e.g., [LRP](#)) can accelerate node motif degree counting.



A Python Library for Graph Outlier/Anomaly Detection

Detecting graph outliers in 5 lines of code

Received 700+ Stars on GitHub

Homepage: <https://pygod.org>

Doc: <https://docs.pygod.org>

Software Paper: <https://arxiv.org/abs/2204.12095>

Email: dev@pygod.org

Backbone	Abbr	Year	Sampling
MLP+AE	MLPAE	2014	Yes
Clustering	SCAN	2007	No
GNN+AE	GCNAE	2016	Yes
MF	Radar	2017	No
MF	ANOMALOUS	2018	No
MF	ONE	2019	No
GNN+AE	DOMINANT	2019	Yes
MLP+AE	DONE	2020	Yes
MLP+AE	AdONE	2020	Yes
GNN+AE	AnomalyDAE	2020	Yes
GAN	GAAN	2020	Yes
GNN+AE	OCGNN	2021	Yes
GNN+AE	CoLA (beta)	2021	In progress
GNN+AE	ANEMONE (beta)	2021	In progress
GNN+AE	GUIDE	2021	Yes
GNN+AE	CONAD	2022	Yes

BOND Benchmark (NeurIPS'22)

- The first comprehensive unsupervised node outlier detection benchmark.
- Provides synthetic, injected, and organic outlier detection dataset.

Data repo: <https://github.com/pygod-team/data>

Benchmark Paper: <https://arxiv.org/abs/2206.10071>

Email: benchmark@pygod.org

Dataset	Type	#Nodes	#Edges	#Feat
'weibo'	organic	8,405	407,963	400
'reddit'	organic	10,984	168,016	64
'disney'	organic	124	335	28
'books'	organic	1,418	3,695	21
'enron'	organic	13,533	176,987	18
'inj_cora'	injected	2,708	11,060	1,433
'inj_amazon'	injected	13,752	515,042	767
'inj_flickr'	injected	89,250	933,804	500
'gen_time'	generated	1,000	5,746	64
'gen_100'	generated	100	618	64
'gen_500'	generated	500	2,662	64
'gen_1000'	generated	1,000	4,936	64
'gen_5000'	generated	5,000	24,938	64
'gen_10000'	generated	10,000	49,614	64

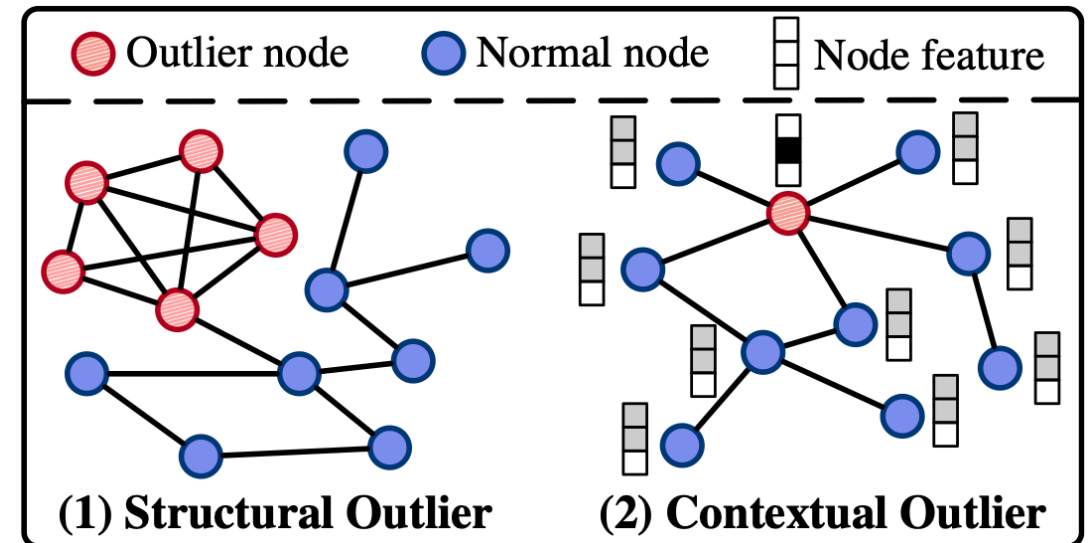
Benchmark on Performance

- No algorithm outperforms on all datasets in expectation.
- Performance on synthetic outliers may not generalize to organic outliers.
- Most Graph OD methods and SGD may be sub-optimal on small graphs.
- Trade-off between algorithm stability and potential in deep graph methods.

Benchmark on Outlier Types

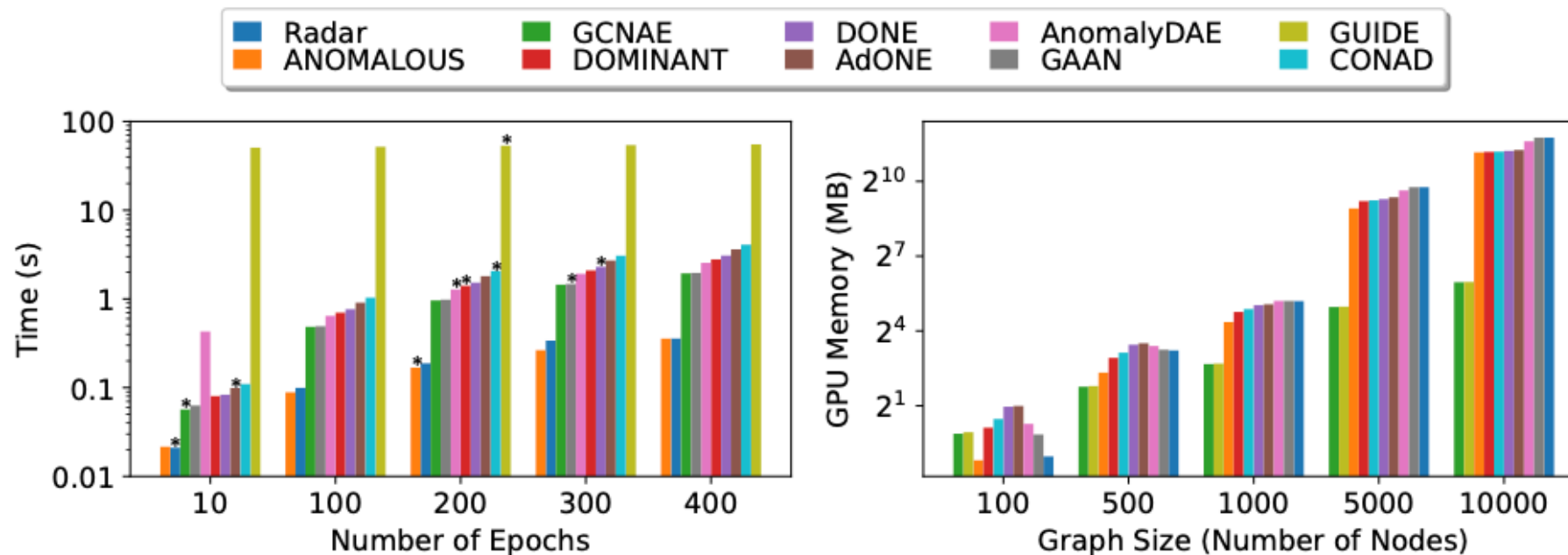
- The **reconstruction** instead of neighbor aggregation detects structural outlier.
- **Low-order structure** is sufficient for detecting structural outlier.
- None of the methods **balances** multiple types of outliers well.

Graph Outlier Taxonomy



Benchmark on Efficiency and Scalability

- Conventional methods are more efficient than deep methods.
- GUIDE improves scalability at an expense of efficiency.



Outline

- Background: fraud detection and graph neural networks.
- Supervised methods and DGFraud.
- Unsupervised methods, PyGOD, and benchmark.
- **Applied graph ML case study.**
- Challenges, solutions, insights, guideline, and resources.

Commercial Graph Database

- Store the data with graph structure.
- Fast update and query.
- Fraud detection capabilities
 - Community/Cycle Detection
 - Link Analysis
 - Graph feature extraction
 - Visualization
- Most of them do not have deep learning/GNN capability.
- TigerGraph ML Workbench has the API to PyG/DGL.



Non-commercial Graph Database

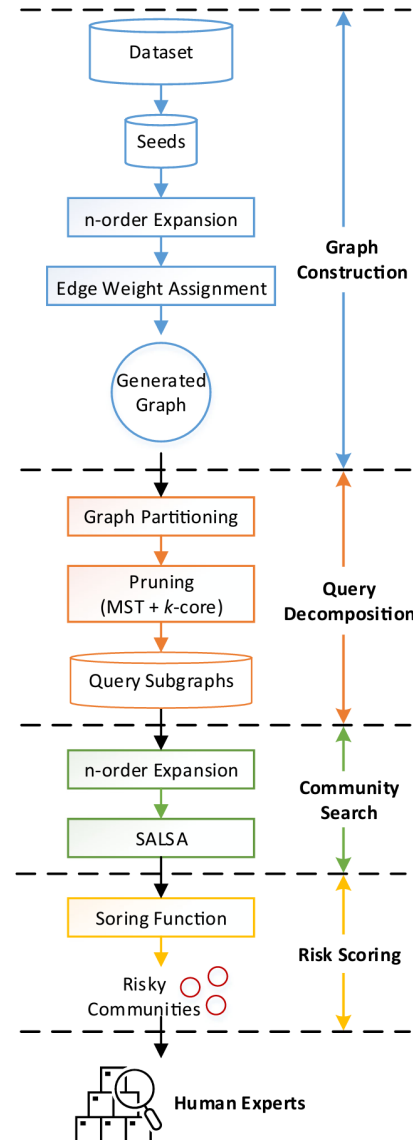
eRiskCom

VLDB Journal 2022

Alibaba

Macquarie University

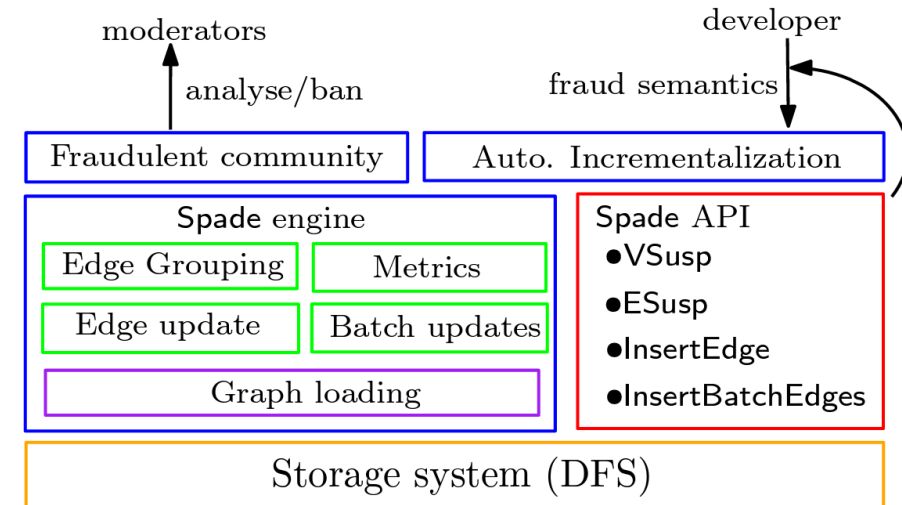
For a suspicious user, find the other key members within the same community of the suspicious users



Spade VLDB 2023

Grab and National University of Singapore

Dynamically maintaining a large graph and fast computing the dense subgraphs

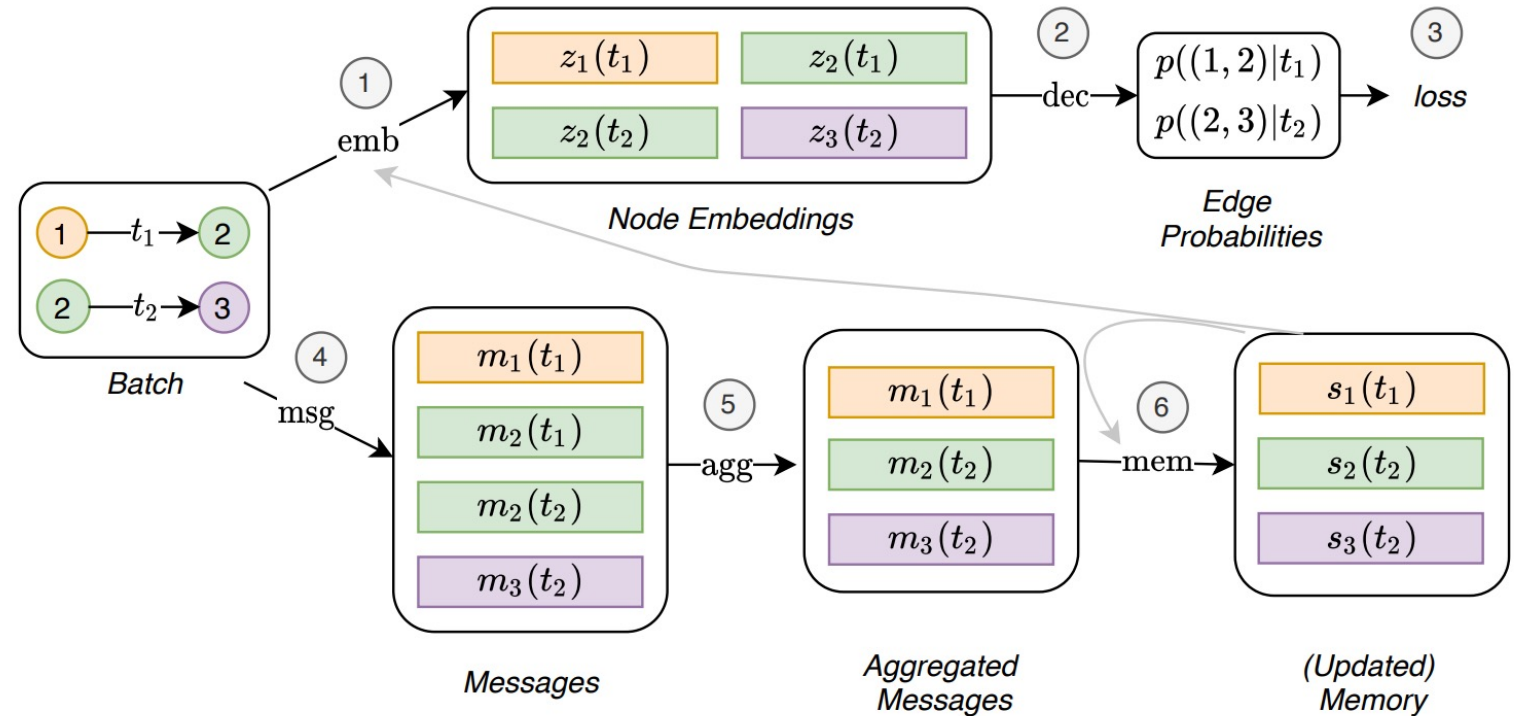


Open Source/Commercial GNN Solution

- [DGL + AWS](#): an end-to-end fraud detection solution (support all DL frameworks).
- [PyG + NVIDIA](#): GNN acceleration on CPU and GPU, no distributed solution.
- [Jraph from DeepMind](#): GNN in JAX, support distributed training.
- [Kumo.ai](#): a startup from PyG team, provide GNN-based fraud detection solution.

Dynamic GNN: A Low-latency Inference Solution

- Each node has a mailbox (memory) to store the up-to-date neighbor information.
- The inference can be done by aggregating information from the memory.
- The memory can be updated asynchronously.
- [APAN](#): e-commerce transaction fraud detection.



Outline

- Background: fraud detection and graph neural networks.
- Supervised methods and DGFraud.
- Unsupervised methods, PyGOD, and benchmark.
- Applied graph ML case study.
- **Challenges, solutions, insights, guideline, and resources.**

Graph-based Financial Fraud Detection Papers

- **Online/Mobile Payment Fraud**
 - [CIKM'18](#), [ICDM'19](#), [AAAI'19](#), [TDSC'20](#).
- **Insurance Fraud**
 - [SIGIR'19](#), [ICME'22](#).
- **Blockchain/Crypto Fraud**
 - [MLF@KDD'19](#), [KDD'21](#), [WWW'22](#), [ADMA'22](#), [KDD'22](#).
- **Loan Defaulting, Loan Fraud, Credit Limit Prediction**
 - [CIKM'19\(1\)](#), [CIKM'19\(2\)](#), [CIKM'20\(1\)](#), [CIKM'20\(2\)](#), [WWW'20](#);
 - [AAAI'21](#), [WSDM'21](#), [WWW'21](#), [SDM'21](#), [arXiv'22](#), [NeurIPS'22](#).
- **Transaction Fraud (e-commerce and credit card)**
 - [ICDE'18](#), [VLDB'19](#), [ICDE'21](#), [arXiv'20](#), [ICAIF'20](#), [SIGMOD'21](#);
 - [KDD'21](#), [WISE'21](#), [TOIS'21](#), [ESA'22](#), [ICML'22](#), [TCSS'22](#).
- **Money Laundering**
 - [MLF@KDD'19](#), [AAAI'20](#).

Key Challenges and Solutions

- **Camouflage**

- Neighboring filtering: [SIGIR'20](#), [CIKM'20](#), [WWW'21](#).
- Aware of adversarial behavior: [IJCAI'20](#), [WWW'20](#).
- Active generative learning: [ACL'20](#).
- Bayesian edge weight inference: [ACL'21](#).

- **Scalability**

- GNN scalability: [MLF@KDD'20](#), [WWW'22\(1\)](#), [WWW'22\(2\)](#), [NeurIPS'22\(1\)](#), [NeurIPS'22\(2\)](#).
- Shallow graph models are more scalable: [MLG@KDD'18](#), [WWW'20](#).

- **Class imbalance**

- Down/Over-sampling: [CIKM'20](#), [WWW'22](#).
- Neighbor selection: [WWW'21](#).
- Data augmentation: [CIKM'20](#).

Key Challenges and Solutions (Cont'd)

- **Label scarcity**

- Active learning: [ICDM'20](#), [TNNLS'21](#), [WWW'22](#).
- Ensemble learning: [CIKM'20](#).
- Meta learning: [WSDM'21](#).

- **Label fidelity**

- Active learning: [TNNLS'21](#).
- Human-in-the-loop: [AAAI'20](#).

- **Data scarcity**

- Data augmentation: [CIKM'20](#), [CIKM'21](#), [ACL'20](#).

Novel Practices

- **Graph Pretraining (Contrastive Learning)**

- Anomalous node is distinguishable from its structural pattern.
- [TNNLS'21](#), [SIGIR'21](#), [arXiv'21\(1\)](#), [arXiv'21\(2\)](#), [ICDM'22](#), [NeurIPS'22](#).

- **Dynamic/Temporal/Streaming Graph**

- Historical information is useful for identifying anomalous.
- Efficiency and cost are bottlenecks.
- [CIKM'21](#), [KDD'21\(1\)](#), [KDD'21\(2\)](#), [SIGMOD'21](#).
- [arXiv'21](#), [SDM'21](#), [ICDM'20](#), [KDD'20](#).
- [ROLAND](#) ([KDD'22](#)), [arXiv'22](#), [ADMA'22](#), [VLDB'23](#).

Novel Practices (Cont'd)

- **Multi-task Learning**

- Credit limit forecasting and credit risk predicting: [WSDM'21](#).
- Fraud detection and recommender system: [SIGIR'20](#).

- **Explainable Anomaly Detection**

- Explainable fraud transaction detection: [arXiv'20](#), [KDD'21](#).
- Explainable fake news detection: [ACL'20](#).
- Explainable fraudulent account detection: [CIKM'22](#).

- **Graph Design**

- Transforming tabular data to graph data for anomaly detection: [IJCAI'22](#).
- Graph architecture search for GNNs: [WWW'22](#).

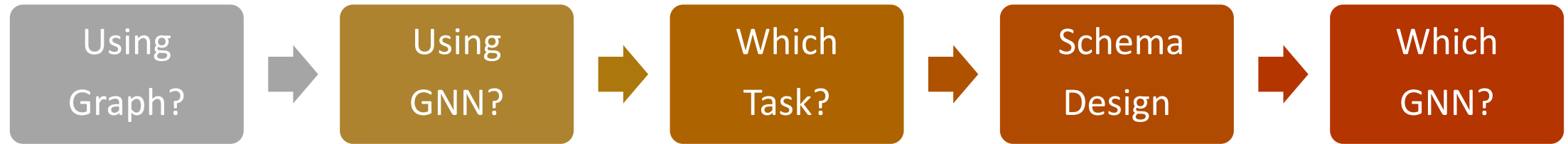
Industry Cases

- **Facebook**
 - [WWW'20](#), [KDD'20](#), [Security'21](#).
- **Amazon**
 - [MLG@KDD'18](#), [KDD'21](#), [WWW'22](#).
- **Alibaba & Ant Group**
 - [CIKM'18](#), [AAAI'19](#), [SIGIR'19](#), [CIKM'19](#), [ICDM'19](#), [IJCAI'20](#), [ACL'20](#), [CIKM'20\(1\)](#);
 - [CIKM'20\(2\)](#), [SIGMOD'21](#), [WSDM'21](#), [WWW'21](#), [AAAI'21](#), [KDD'21\(1\)](#), [KDD'21\(2\)](#), [VLDB'22](#).
- **eBay**
 - [Workshop@AAAI'21](#), [arXiv'20](#), [MLF@KDD'20](#).
- **Others**
 - [App Market](#), [Money Laundering](#), [Fake Invitation \(iQIYI\)](#), [Blockchain\(1\)](#), [Blockchain\(2\)](#);
 - [Blockchain\(3\)](#), [Grab\(1\)](#), [Grab\(2\)](#), [Custom Fraud](#), [Finvolution](#).
- **Tencent**
 - [WWW'19](#), [WWW'20](#), [KDD'21](#).

Insights and Discussions

- **Early detection, prevention vs. detection.**
- **Unsupervised model selection/hyper-parameter tuning.**
- **The longtail distribution of the anomaly types.**
- **The concept drift and continual learning.**
- **The gap between academia and industry.**

How to Apply GNNs in Anomaly Detection



- **Using Graph?**

- The anomalous entities share common properties.
- The anomalous entities have clustering behavior.
- The trade off between cost and effectiveness.

- **Using GNN?**

- The infrastructure.
- The feature availability and feature types.
- Integrating with other modules and tasks.

- **Which Task?**

- Supervised:
 - Node/edge/graph/subgraph classification.
- Unsupervised:
 - Community detection; anomaly detection.

- **Schema Design**

- Node/edge type and node/edge feature.
- Graph schema, node sampling.
- Graph structure is flexible: [SIGIR'19](#), [ICDM'20](#).

- **Which GNN?**

- GNN is chosen based on task and schema.
- Simple GNN model is enough.
- Dynamic GNNs need more efforts.
- GAT and Graph-SAGE are commonly used.

Resources

- **Paper List**

- [Graph-based Fraud Detection Papers and Resources.](#)

- **Tutorial**

- [KDD'20 Tutorial: Deep Learning for Anomaly Detection](#)

- **Libraries**

- [PyG Temporal](#), [PyOD](#), [PyODDS](#), [TODS](#), [UGFraud](#).

- **Recent Surveys**

- [Graph Learning for Anomaly Analytics: Algorithms, Applications, and Challenges.](#)
 - [Graph Computing for Financial Crime and Fraud Detection: Trends, Challenges and Outlook.](#)

Thank you!

Q & A

Yingtong Dou

[@dozee_sim](#)

Kay Liu

[@kayzliu](#)

University of Illinois Chicago

@Wells Fargo

12/06/2022

