

Unveiling the Secrets without Data: Can Graph Neural Networks Be Exploited through Data-Free Model Extraction Attacks?

Yuanxin Zhuang

Chuan Shi *

Mengmei Zhang

Beijing University of Posts and Telecommunications

Key Laboratory of Trustworthy Distributed Computing and Service (BUPT), Ministry of Education

Jinghui Chen

The Pennsylvania State University

Lingjuan Lyu

SONY AI

Pan Zhou

Huazhong University of Science and Technology

Lichao Sun

Lehigh University

Yuanxin ZHUANG

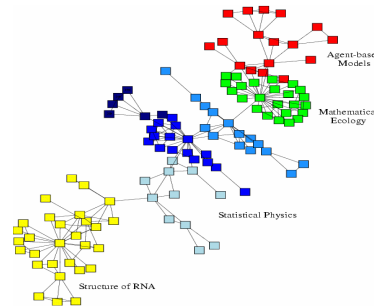
yzhuang436@connect.hkust-gz.edu.cn

Why Do Network Representation?

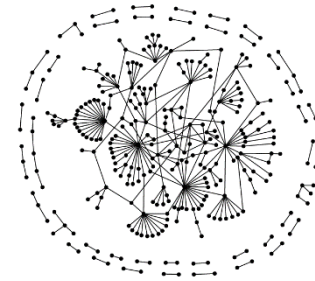
Networks are a general language for describing and modeling complex systems



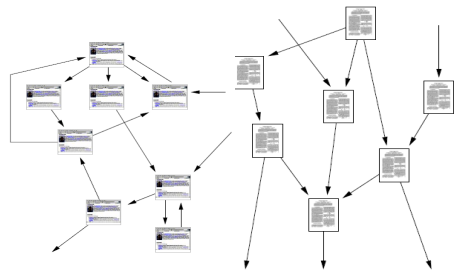
Social networks



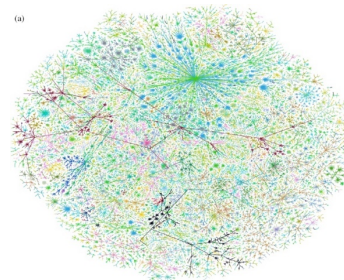
Economic networks



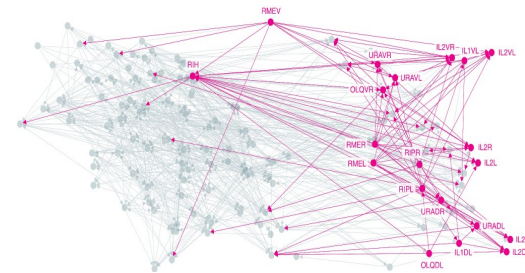
Biomedical networks



Information networks



Internet



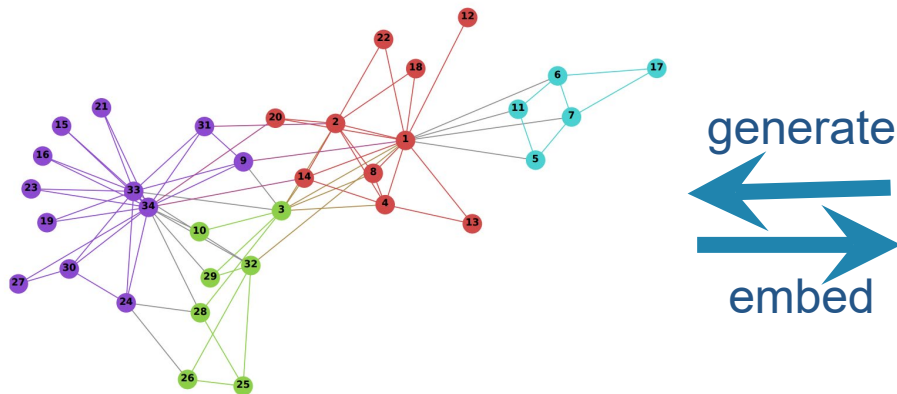
Networks of neurons

Network Representation

Network Representation

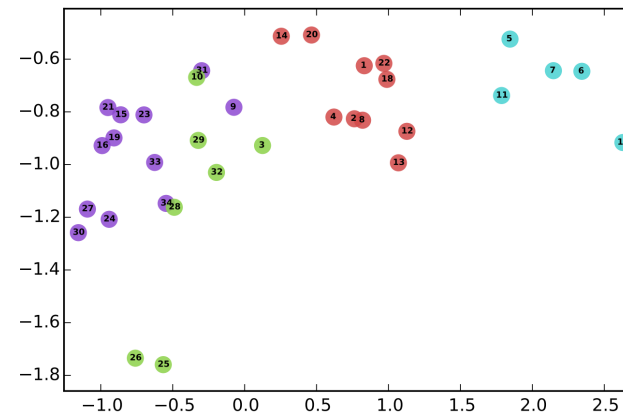
Embed each node of a network into a low-dimensional vector space

- Easy to parallel
- Can apply classical ML methods



Applications

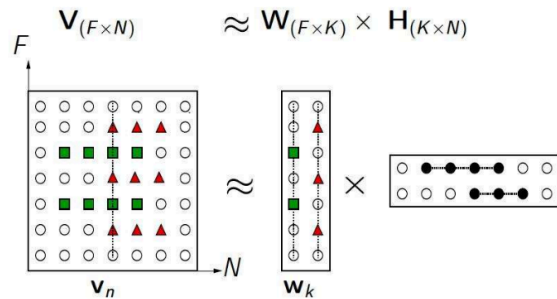
- Node classification
- Link predication
- Community detection
- Network evolution
- ...



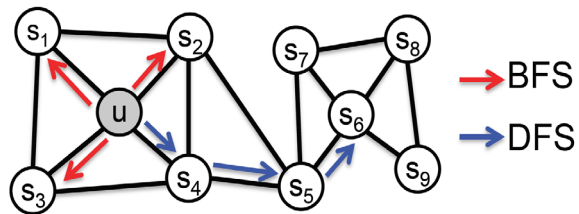
Network Representation

Shallow model

- Factorization based
 - e.g., Laplacian eigenmaps

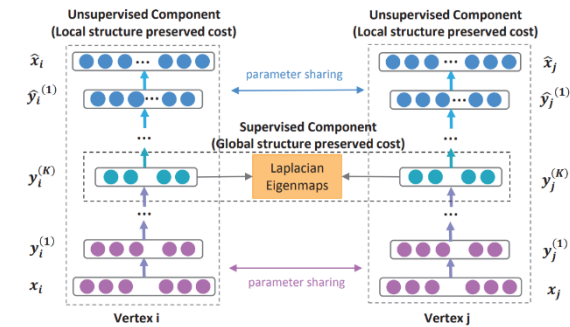


- Random walk based
 - e.g., DeepWalk, node2vec

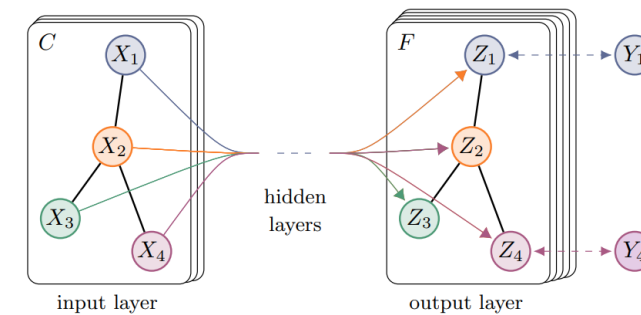


Deep model

- Autoencoder based
 - e.g., DNGR and SDNE

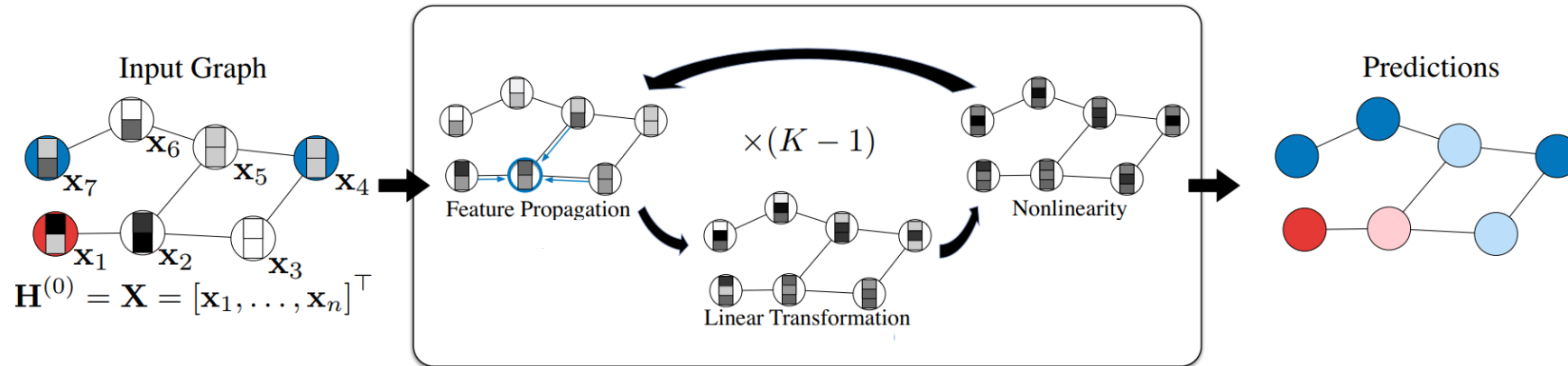


- GNN based
 - e.g., GCN, GraphSage, GAT



Graph Neural Networks

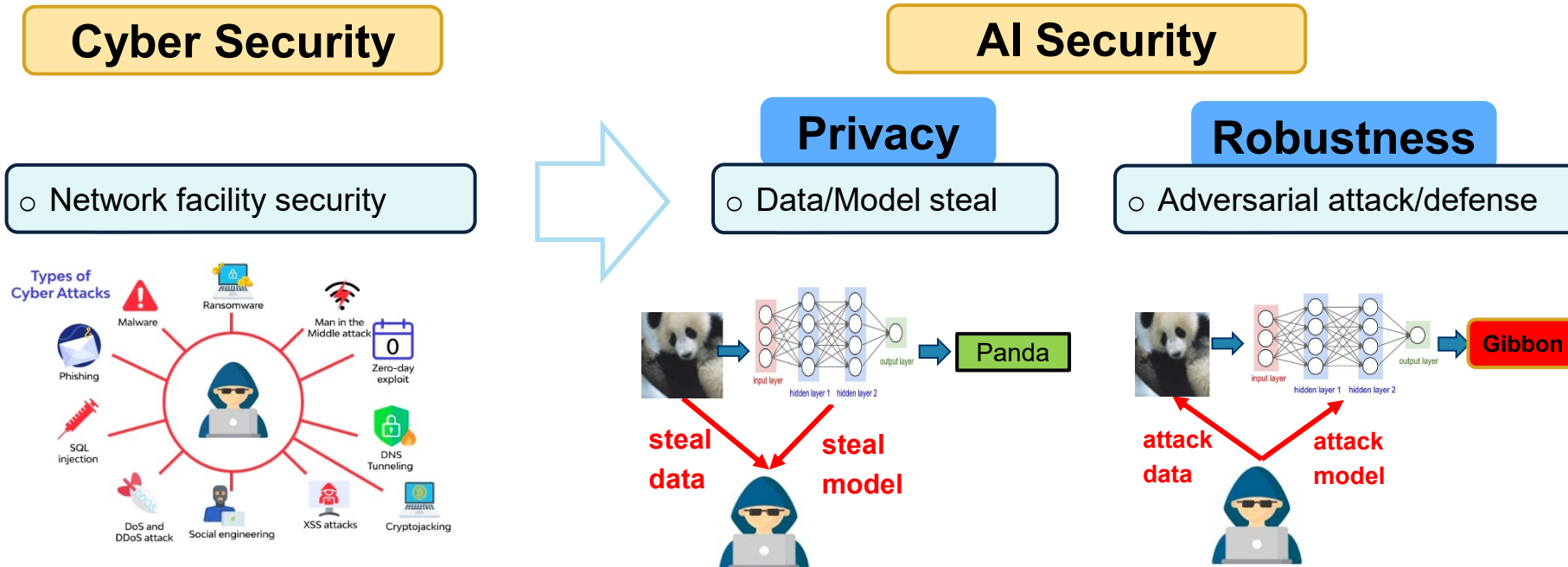
Learning Process of GNN



1. Input graph \mathbf{A} and node feature \mathbf{X}
2. Initialization: $h_i^{(0)} = \sigma(W_0 X_i), \forall i$
3. k -th iteration: $h_i^{(k)} \leftarrow \sigma \left(W_1 h_i^{(k-1)} + W_2 \sum_{j \in \mathcal{N}(i)} h_j^{(k-1)} \right), \forall i$
4. Prediction: $\hat{\mathbf{Y}} = \text{softmax}(\mathbf{A}\mathbf{H}^{(K-1)}\mathbf{W})$

AI Security

- Focus on self security of AI algorithms
- Guarantee the integration and privacy of AI model and data
 - Privacy: The ability of protecting model and data from leakage
 - Robustness: Model adaptability to uncertainty, noise, and attacks



Examples of AI Security

Privacy

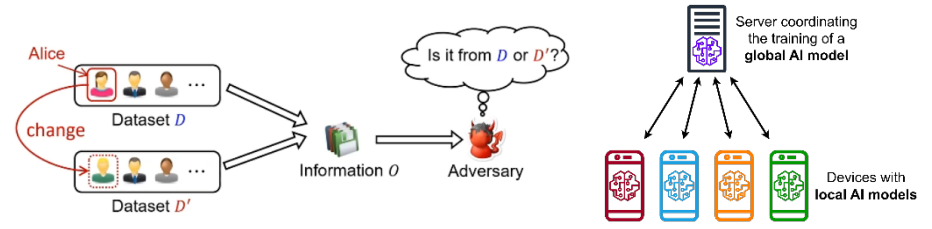
Privacy Attack



Membership inference

Model extraction

Privacy Defense

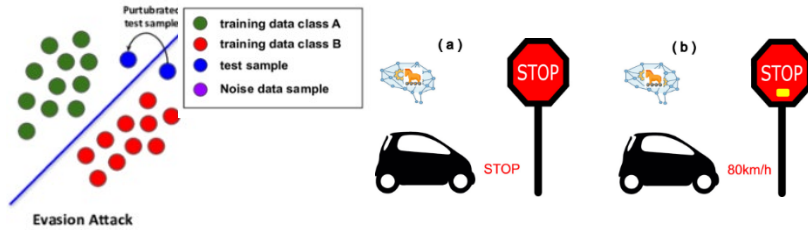


Differential privacy

Federated learning

Robustness

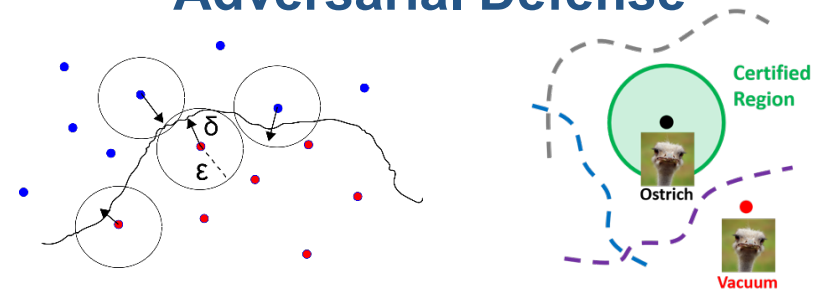
Adversarial Attack



Evasion attack

Backdoor attack

Adversarial Defense



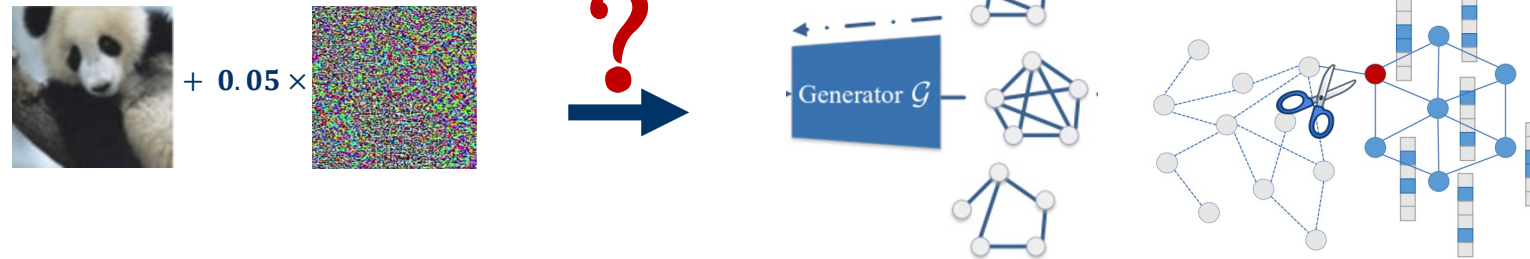
Adversarial training

Certifiable robustness

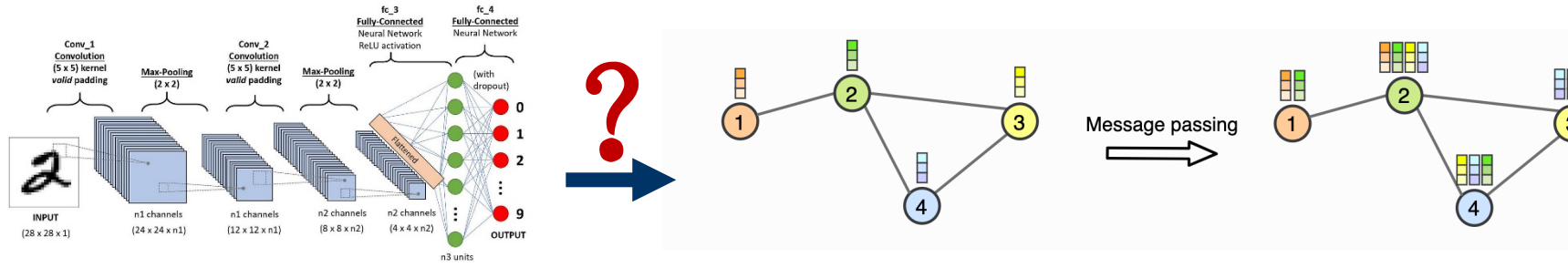
Security of GNNs

Challenges

- Unique data type
 - Discreteness graph structure



- Unique model design
 - Message-passing mechanism



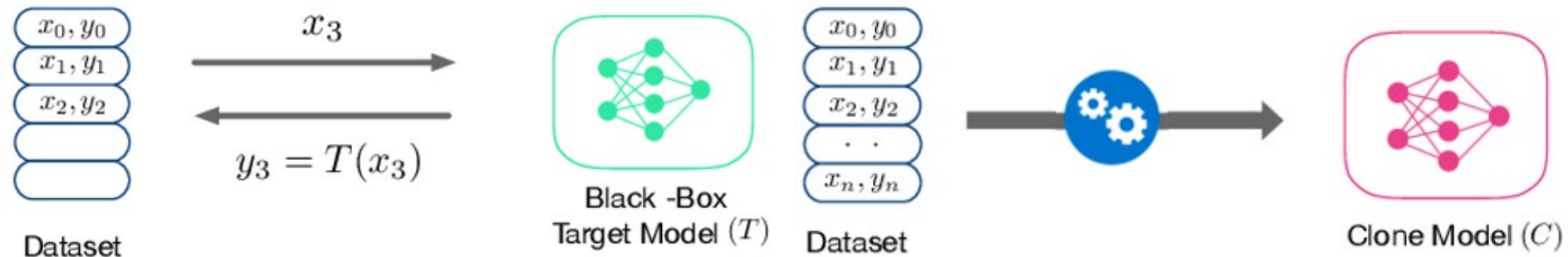
Model Extraction

□ Model Extraction

- Extracting model's parameters and structure for analysis, replication, and security assessment.

□ Steps of Model Extraction

- Generate data to capture the behavior of the target model.
- Train a new clone model using the constructed dataset.



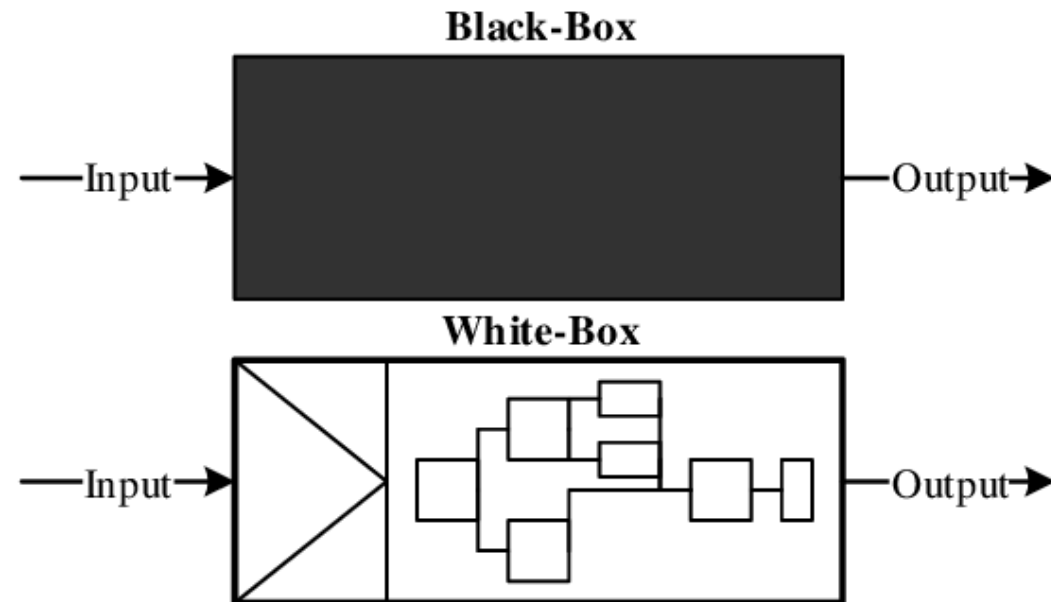
Step1: Construct Training Dataset

Step2: Train the Clone Model

Model Extraction

□ Categories of Model Extraction

- **White-box** involves direct access to the internal parameters and structure of a model
- **Black-box** relies on observing the model's input-output behavior to infer its internal workings



Data-Free Model Extraction Attacks

□ Motivation

- Data is **unobtainable** and characterized by **privacy concerns**
- In real-world, pre-trained models are proprietary and **black-box**.

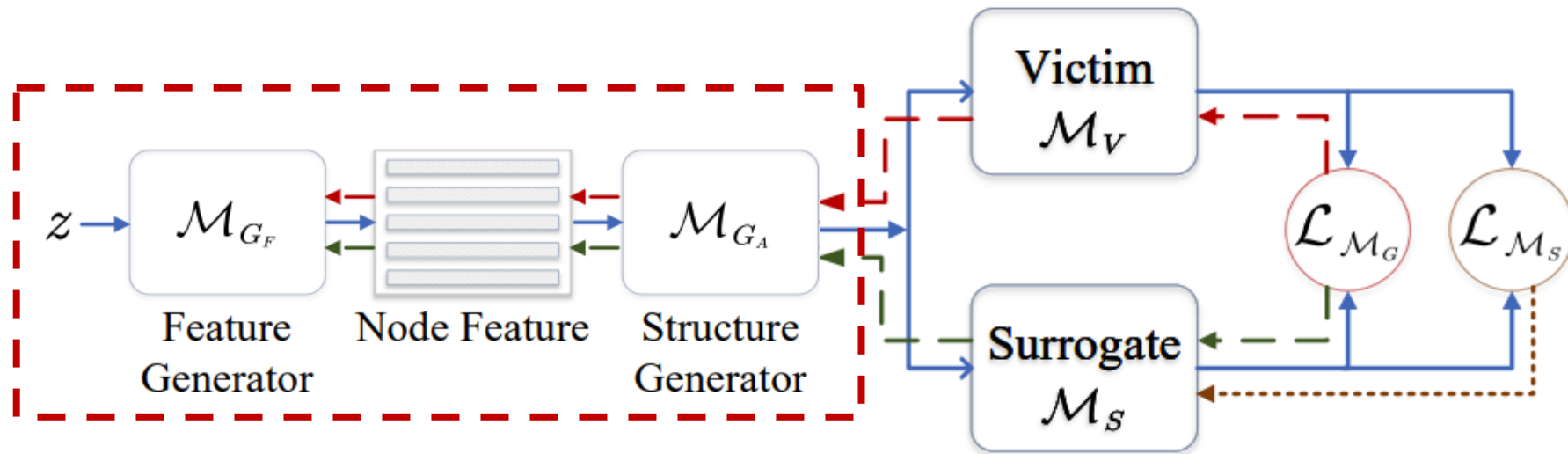
□ Comparison with Existing Work

Method	Data-Free	Node Classification		Link Prediction
		Transductive	Inductive	
DeFazio et al.'s work [6]		✓		
Wu et al.'s work [42]		✓		
Shen et al.'s work [36]			✓	
STEALGNN	✓	✓	✓	✓

Framework of StealGNN

Data-free Model Extraction Attack against GNN (StealGNN)

□ Model Framework



Step 1:

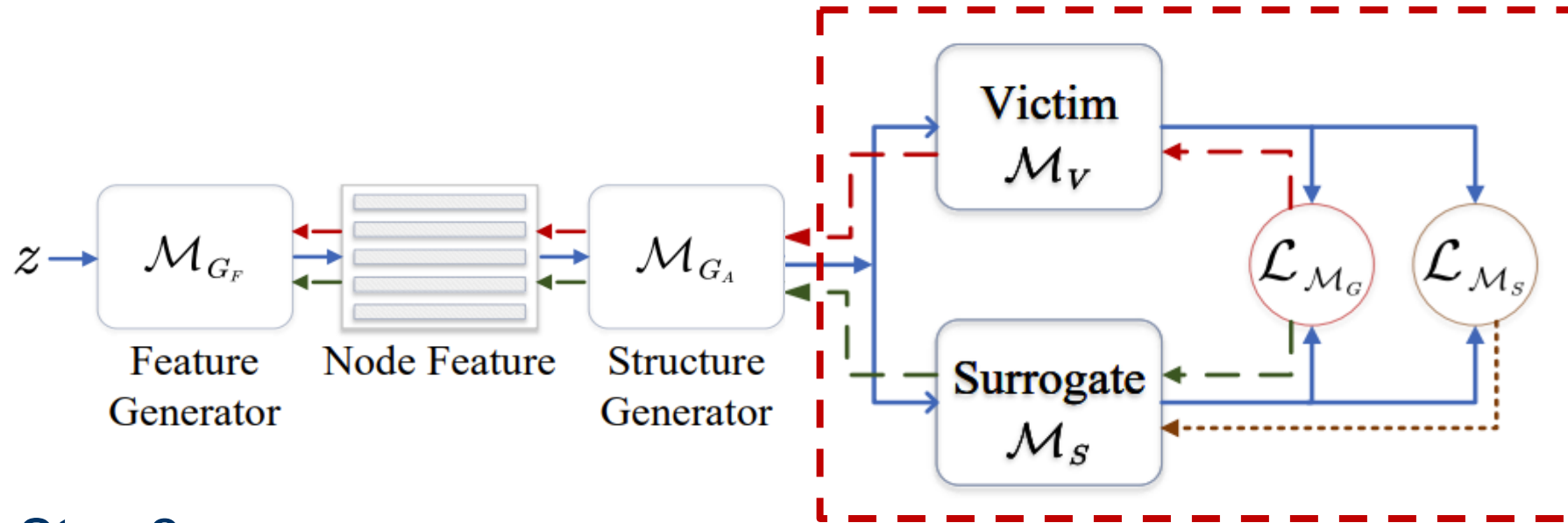
Generate graph from noise:

$$\mathbf{G} = \mathcal{M}_G(z; \theta_G); \quad z \sim \mathcal{N}(0, I),$$

Framework of StealGNN

Data-free Model Extraction Attack against GNN (StealGNN)

□ Model Framework



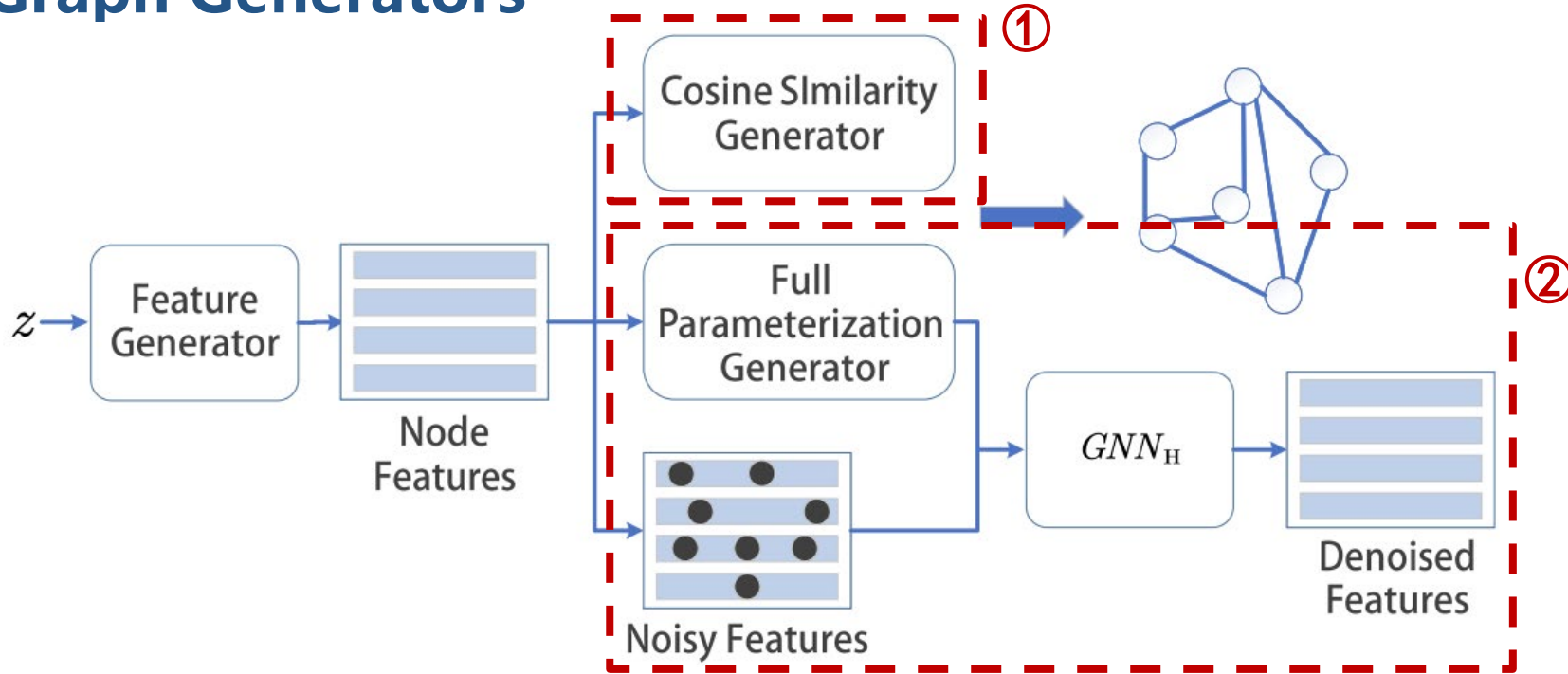
Step 2:

Extract model using generated graph :

$$\vec{Y}_V = \mathcal{M}_V(\mathbf{G}, \theta_V); \quad \vec{P}_S = \mathcal{M}_S(\mathbf{G}; \theta_S)$$

Framework of StealGNN

□ Graph Generators



Two structure generation methods:

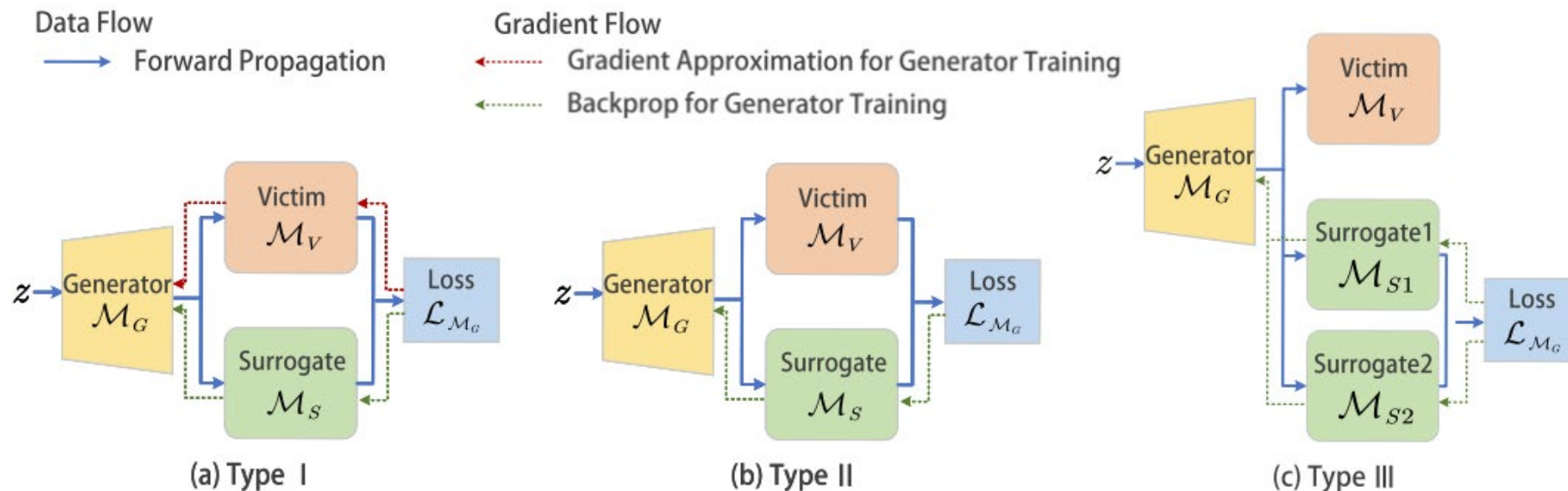
① **Cosine Similarity:** $\mathbf{A}_{cosij} = \text{Cosine}(\mathbf{F}_i, \mathbf{F}_j) = \frac{\mathbf{F}_i \cdot \mathbf{F}_j}{\max(\|\mathbf{F}_i\|_2, \|\mathbf{F}_j\|_2)}$

② **Full Parameterization:**

$$\mathcal{L}_{\text{denoise}} = \frac{1}{|\mathbf{F}_{idx}|} \sum_{i \in \mathbf{F}_{idx}} \mathcal{L}_{\text{reconstruction}}(\mathbf{F}_i, \mathbf{H}(\tilde{\mathbf{F}}, \mathbf{A}_{fp}; \boldsymbol{\theta}_H)_i)$$

Framework of StealGNN

□ Generator Parameter Update Strategy



Type I : Gradient Approximation + Surrogate. **Type II** : One Surrogate.

$$\mathcal{L}_{\mathcal{M}_G} = -\mathcal{L}_{\mathcal{M}_S}$$

Type III: Using the inconsistency between predictions of two surrogate models to generate informative graphs.

$$\mathcal{L}_{\mathcal{M}_G} = -\frac{1}{N} \sum_{i=1}^N [\text{Std}(\mathbf{P}_{S1_i}, \mathbf{P}_{S2_i})]$$

Experiments

Model Performance (node classification)

\mathcal{M}_V	Dataset	Cora		Pubmed		A-Computers		OGB-Arxiv	
	$\mathcal{M}_S(\text{GCN})$	Accuracy	Fidelity	Accuracy	Fidelity	Accuracy	Fidelity	Accuracy	Fidelity
GAT	Real Data	80.09 \pm 0.89 (+0.53)	92.96 \pm 0.51	76.94 \pm 0.31 (-0.12)	89.96 \pm 2.51	79.20 \pm 1.67 (-0.89)	90.51 \pm 1.12	53.89 \pm 1.02 (-0.77)	88.73 \pm 1.94
	Random Graph	68.26 \pm 3.84 (-11.30)	73.70 \pm 2.95	59.35 \pm 2.95 (-17.71)	71.79 \pm 2.49	55.63 \pm 4.61 (-24.46)	52.29 \pm 3.62	36.78 \pm 5.25 (-17.88)	63.52 \pm 4.39
	Attack I-E	81.11 \pm 0.50 (+1.55)	93.09 \pm 0.26	77.57 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack I-A _{cos}	80.79 \pm 0.75 (+1.23)	93.90 \pm 0.42	77.15 \pm 0.93 (+0.09)	90.34 \pm 2.94	70.73 \pm 1.90 (-9.36)	81.18 \pm 1.18	51.72 \pm 2.31 (-2.94)	78.77 \pm 1.68
	Attack I-A _{fp}	81.23 \pm 0.68 (+1.67)	94.10 \pm 0.38	77.40 \pm 0.87 (+0.34)	90.90 \pm 2.76	71.12 \pm 1.78 (-8.54)	81.45 \pm 1.10	52.15 \pm 2.23 (-2.51)	79.22 \pm 1.58
	Attack II-E	80.98 \pm 0.57 (+1.42)	92.78 \pm 0.39	77.15 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack II-A _{cos}	81.17 \pm 0.71 (+1.61)	93.19 \pm 0.47	77.40 \pm 0.87 (+0.34)	90.90 \pm 2.76	71.12 \pm 1.78 (-8.54)	81.45 \pm 1.10	52.15 \pm 2.23 (-2.51)	79.22 \pm 1.58
	Attack II-A _{fp}	81.35 \pm 0.66 (+1.79)	93.50 \pm 0.43	77.57 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack III-E	81.12 \pm 0.52 (+1.56)	93.02 \pm 0.26	77.31 \pm 0.78 (+0.51)	91.91 \pm 2.76	70.73 \pm 1.90 (-9.36)	81.18 \pm 1.18	51.72 \pm 2.31 (-2.94)	78.77 \pm 1.68
GCN	Real Data	81.09 \pm 0.69 (+0.49)	92.84 \pm 0.29	77.06 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Random Graph	69.83 \pm 2.61 (-10.78)	81.79 \pm 2.28	62.78 \pm 2.95	73.70 \pm 2.95	59.35 \pm 2.95 (-17.71)	71.79 \pm 2.49	55.63 \pm 4.61 (-24.46)	52.29 \pm 3.62
	Attack I-E	80.67 \pm 0.45 (+0.06)	93.57 \pm 0.96	75.97 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack I-A _{cos}	80.68 \pm 0.29 (+0.07)	93.77 \pm 1.37	76.77 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack I-A _{fp}	80.91 \pm 0.24 (+0.30)	93.96 \pm 1.24	76.79 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack II-E	79.79 \pm 0.23 (-0.82)	92.81 \pm 0.79	74.11 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack II-A _{cos}	80.07 \pm 0.59 (-0.54)	92.84 \pm 0.31	76.41 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack II-A _{fp}	81.12 \pm 0.28 (+0.51)	93.65 \pm 0.12	76.71 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack III-E	80.78 \pm 0.31 (+0.17)	93.39 \pm 0.15	76.41 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
SAGE	Real Data	79.84 \pm 1.25 (+0.51)	92.57 \pm 0.28	76.71 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Random Graph	60.39 \pm 2.95 (-18.94)	78.39 \pm 3.38	58.75 \pm 2.95	73.70 \pm 2.95	59.35 \pm 2.95 (-17.71)	71.79 \pm 2.49	55.63 \pm 4.61 (-24.46)	52.29 \pm 3.62
	Attack I-E	81.32 \pm 0.65 (+1.99)	93.01 \pm 0.12	77.06 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack I-A _{cos}	80.89 \pm 0.47 (+1.56)	93.97 \pm 0.49	77.22 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack I-A _{fp}	81.10 \pm 0.35 (+1.77)	93.85 \pm 0.33	77.40 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack II-E	81.28 \pm 0.67 (+1.95)	93.92 \pm 0.70	76.81 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack II-A _{cos}	81.01 \pm 0.51 (+1.68)	95.14 \pm 0.17	77.15 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack II-A _{fp}	80.98 \pm 0.41 (+1.65)	93.68 \pm 0.39	77.22 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack III-E	81.56 \pm 0.45 (+2.23)	93.28 \pm 0.09	77.41 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
GAT	Real Data	81.45 \pm 0.63 (+2.12)	93.10 \pm 0.10	77.21 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Random Graph	81.55 \pm 0.61 (+2.22)	93.15 \pm 0.09	77.25 \pm 0.14 (+0.02)	92.90 \pm 0.41	70.55 \pm 2.43 (-8.69)	86.65 \pm 1.32	58.55 \pm 3.78 (-10.63)	78.25 \pm 1.97
	Attack I-E	81.35 \pm 0.66 (+1.79)	93.50 \pm 0.43	77.57 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack I-A _{cos}	81.45 \pm 0.62 (+1.89)	93.80 \pm 0.41	77.40 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack I-A _{fp}	81.68 \pm 0.58 (+2.12)	94.05 \pm 0.39	77.80 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack II-E	81.12 \pm 0.52 (+1.56)	93.02 \pm 0.26	77.31 \pm 0.78 (+0.51)	91.91 \pm 2.76	70.73 \pm 1.90 (-9.36)	81.18 \pm 1.18	51.72 \pm 2.31 (-2.94)	78.77 \pm 1.68
	Attack II-A _{cos}	81.17 \pm 0.71 (+1.61)	93.19 \pm 0.47	77.40 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack II-A _{fp}	81.35 \pm 0.66 (+1.79)	93.50 \pm 0.43	77.57 \pm 0.78 (+0.51)	92.09 \pm 2.56	65.68 \pm 1.91 (-14.41)	73.48 \pm 2.13	52.17 \pm 2.94 (-2.49)	80.26 \pm 2.23
	Attack III-E	81.12 \pm 0.52 (+1.56)	93.02 \pm 0.26	77.31 \pm 0.78 (+0.51)	91.91 \pm 2.76	70.73 \pm 1.90 (-9.36)	81.18 \pm 1.18	51.72 \pm 2.31 (-2.94)	78.77 \pm 1.68

Experiments

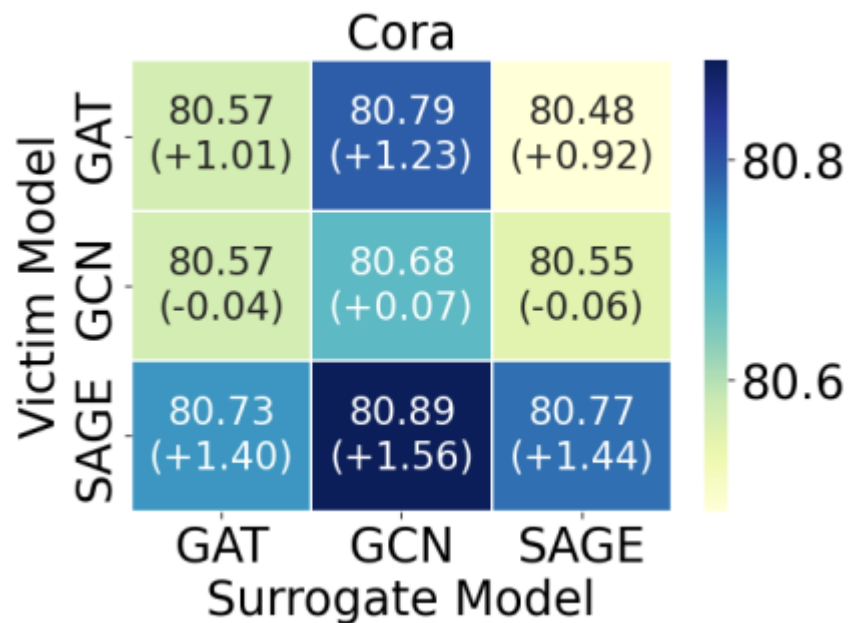
Model Performance (link prediction)

\mathcal{M}_V	Dataset	Cora	Pubmed	A-computers	\mathcal{M}_V	Dataset	Cora	Pubmed	A-computers
GAT	Real Data	89.30 \pm 0.59 (-0.38)	65.70 \pm 1.58 (-4.67)	84.64 \pm 1.78 (+1.89)	GAT	Real Data	87.70 \pm 0.47	68.50 \pm 2.20	85.11 \pm 1.38
	Random Graph	75.32 \pm 2.12 (-14.36)	50.21 \pm 3.89 (-20.16)	66.79 \pm 2.78 (-15.96)		Random Graph	78.25 \pm 0.35	52.10 \pm 1.95	72.80 \pm 1.05
	Attack I- \bar{E}	89.01 \pm 0.56 (-0.67)	64.78 \pm 1.61 (-5.59)	80.29 \pm 1.48 (-2.46)		Attack I- \bar{E}	85.91 \pm 0.86	68.19 \pm 1.20	79.99 \pm 1.05
	Attack I- A_{cos}	88.71 \pm 0.32 (-0.97)	64.98 \pm 1.92 (-5.39)	80.62 \pm 1.19 (-2.13)		Attack I- A_{cos}	86.79 \pm 0.62	67.08 \pm 0.99	81.73 \pm 1.49
	Attack I- A_{fp}	89.12 \pm 0.22 (-0.56)	64.25 \pm 1.38 (-6.12)	79.92 \pm 0.92 (-2.83)		Attack I- A_{fp}	85.45 \pm 0.92	68.60 \pm 1.05	79.75 \pm 0.98
	Attack II- \bar{E}	88.96 \pm 0.63 (-0.72)	64.58 \pm 1.82 (-5.79)	80.19 \pm 1.19 (-2.56)		Attack II- \bar{E}	89.96 \pm 0.83	66.29 \pm 1.52	79.61 \pm 1.38
	Attack II- A_{cos}	89.38 \pm 0.47 (-0.30)	64.39 \pm 1.25 (-5.98)	79.89 \pm 0.92 (-2.86)		Attack II- A_{cos}	90.72 \pm 0.67	65.42 \pm 1.65	80.59 \pm 1.17
	Attack II- A_{fp}	88.95 \pm 0.62 (-0.73)	66.12 \pm 1.61 (-4.25)	78.98 \pm 1.28 (-3.77)		Attack II- A_{fp}	90.40 \pm 0.72	65.75 \pm 1.58	80.85 \pm 1.10
	Attack III- \bar{E}	88.95 \pm 0.50 (-0.73)	66.78 \pm 1.32 (-3.59)	79.47 \pm 0.98 (-3.28)		Attack III- \bar{E}	91.02 \pm 0.64	66.02 \pm 1.53	81.13 \pm 1.04
	Attack III- A_{cos}	89.25 \pm 0.53 (-0.43)	65.78 \pm 1.42 (-4.59)	80.35 \pm 1.07 (-2.90)		Attack III- A_{cos}	90.56 \pm 0.64	66.10 \pm 1.52	81.08 \pm 1.08
Attack III- A_{fp}	89.45 \pm 0.35 (-0.23)	66.25 \pm 1.15 (-4.12)	80.52 \pm 0.98 (-2.98)	Attack III- A_{fp}	90.73 \pm 0.64	68.82 \pm 0.36	81.03 \pm 1.12		
GCN	Real data	89.59 \pm 0.23 (-0.09)	87.72 \pm 0.18 (-0.38)	89.59 \pm 0.23 (-0.09)	GCN	Real data	89.59 \pm 0.23 (-0.09)	87.72 \pm 0.18 (-0.38)	89.59 \pm 0.23 (-0.09)
	Random Graph	80.12 \pm 2.18 (-9.56)	72.85 \pm 3.29 (-15.25)	80.12 \pm 2.18 (-9.56)		Random Graph	80.12 \pm 2.18 (-9.56)	72.85 \pm 3.29 (-15.25)	80.12 \pm 2.18 (-9.56)
	Attack I- \bar{E}	90.78 \pm 0.73 (+1.10)	85.54 \pm 1.09 (-2.36)	90.78 \pm 0.73 (+1.10)		Attack I- \bar{E}	90.78 \pm 0.73 (+1.10)	85.54 \pm 1.09 (-2.36)	90.78 \pm 0.73 (+1.10)
	Attack I- A_{cos}	90.83 \pm 0.41 (+1.15)	84.20 \pm 1.67 (-2.90)	90.83 \pm 0.41 (+1.15)		Attack I- A_{cos}	90.83 \pm 0.41 (+1.15)	84.20 \pm 1.67 (-2.90)	90.83 \pm 0.41 (+1.15)
	Attack I- A_{fp}	89.70 \pm 0.60 (-0.02)	87.45 \pm 0.85 (-0.76)	89.70 \pm 0.60 (-0.02)		Attack I- A_{fp}	89.70 \pm 0.60 (-0.02)	87.45 \pm 0.85 (-0.76)	89.70 \pm 0.60 (-0.02)
	Attack II- \bar{E}	89.95 \pm 0.52 (+0.27)	87.69 \pm 0.74 (-0.41)	89.95 \pm 0.52 (+0.27)		Attack II- \bar{E}	89.95 \pm 0.52 (+0.27)	87.69 \pm 0.74 (-0.41)	89.95 \pm 0.52 (+0.27)
	Attack II- A_{cos}	89.09 \pm 0.25 (-0.59)	88.05 \pm 0.28 (-0.05)	89.09 \pm 0.25 (-0.59)		Attack II- A_{cos}	89.09 \pm 0.25 (-0.59)	88.05 \pm 0.28 (-0.05)	89.09 \pm 0.25 (-0.59)
	Attack II- A_{fp}	89.30 \pm 0.22 (-0.38)	88.25 \pm 0.25 (+0.15)	89.30 \pm 0.22 (-0.38)		Attack II- A_{fp}	89.30 \pm 0.22 (-0.38)	88.25 \pm 0.25 (+0.15)	89.30 \pm 0.22 (-0.38)
	Attack III- \bar{E}	90.25 \pm 0.45 (+0.57)	87.85 \pm 0.62 (-0.25)	90.25 \pm 0.45 (+0.57)		Attack III- \bar{E}	90.25 \pm 0.45 (+0.57)	87.85 \pm 0.62 (-0.25)	90.25 \pm 0.45 (+0.57)
	Attack III- A_{cos}	90.40 \pm 0.42 (+0.72)	88.09 \pm 0.55 (-0.01)	90.40 \pm 0.42 (+0.72)		Attack III- A_{cos}	90.40 \pm 0.42 (+0.72)	88.09 \pm 0.55 (-0.01)	90.40 \pm 0.42 (+0.72)
Attack III- A_{fp}	90.65 \pm 0.40 (+0.97)	88.12 \pm 0.52 (+0.02)	90.65 \pm 0.40 (+0.97)	Attack III- A_{fp}	90.65 \pm 0.40 (+0.97)	88.12 \pm 0.52 (+0.02)	90.65 \pm 0.40 (+0.97)		
SAGE	Real data	91.15 \pm 0.20 (-1.32)	87.90 \pm 0.09 (+1.06)	91.15 \pm 0.20 (-1.32)	SAGE	Real data	91.15 \pm 0.20 (-1.32)	87.90 \pm 0.09 (+1.06)	91.15 \pm 0.20 (-1.32)
	Random Graph	77.25 \pm 1.35 (-10.54)	75.60 \pm 1.75 (-11.34)	77.25 \pm 1.35 (-10.54)		Random Graph	77.25 \pm 1.35 (-10.54)	75.60 \pm 1.75 (-11.34)	77.25 \pm 1.35 (-10.54)
	Attack I- \bar{E}	88.68 \pm 0.54 (-3.79)	85.79 \pm 0.48 (-1.05)	88.68 \pm 0.54 (-3.79)		Attack I- \bar{E}	88.68 \pm 0.54 (-3.79)	85.79 \pm 0.48 (-1.05)	88.68 \pm 0.54 (-3.79)
	Attack I- A_{cos}	87.79 \pm 0.96 (-4.68)	86.94 \pm 0.51 (+0.00)	87.79 \pm 0.96 (-4.68)		Attack I- A_{cos}	87.79 \pm 0.96 (-4.68)	86.94 \pm 0.51 (+0.00)	87.79 \pm 0.96 (-4.68)
	Attack I- A_{fp}	88.05 \pm 0.45 (-4.42)	87.30 \pm 0.55 (+0.36)	88.05 \pm 0.45 (-4.42)		Attack I- A_{fp}	88.05 \pm 0.45 (-4.42)	87.30 \pm 0.55 (+0.36)	88.05 \pm 0.45 (-4.42)
	Attack II- \bar{E}	89.01 \pm 0.56 (-3.46)	86.77 \pm 0.39 (-0.07)	89.01 \pm 0.56 (-3.46)		Attack II- \bar{E}	89.01 \pm 0.56 (-3.46)	86.77 \pm 0.39 (-0.07)	89.01 \pm 0.56 (-3.46)
	Attack II- A_{cos}	90.04 \pm 0.47 (-2.43)	88.26 \pm 0.08 (+1.42)	90.04 \pm 0.47 (-2.43)		Attack II- A_{cos}	90.04 \pm 0.47 (-2.43)	88.26 \pm 0.08 (+1.42)	90.04 \pm 0.47 (-2.43)
	Attack II- A_{fp}	88.50 \pm 0.35 (-4.97)	87.80 \pm 0.50 (+0.96)	88.50 \pm 0.35 (-4.97)		Attack II- A_{fp}	88.50 \pm 0.35 (-4.97)	87.80 \pm 0.50 (+0.96)	88.50 \pm 0.35 (-4.97)
	Attack III- \bar{E}	90.05 \pm 0.50 (-2.42)	87.32 \pm 0.30 (+0.48)	90.05 \pm 0.50 (-2.42)		Attack III- \bar{E}	90.05 \pm 0.50 (-2.42)	87.32 \pm 0.30 (+0.48)	90.05 \pm 0.50 (-2.42)
	Attack III- A_{cos}	90.31 \pm 0.40 (-2.16)	88.15 \pm 0.10 (+1.31)	90.31 \pm 0.40 (-2.16)		Attack III- A_{cos}	90.31 \pm 0.40 (-2.16)	88.15 \pm 0.10 (+1.31)	90.31 \pm 0.40 (-2.16)
Attack III- A_{fp}	90.59 \pm 0.40 (-1.88)	87.43 \pm 0.25 (+0.59)	90.59 \pm 0.40 (-1.88)	Attack III- A_{fp}	90.59 \pm 0.40 (-1.88)	87.43 \pm 0.25 (+0.59)	90.59 \pm 0.40 (-1.88)		

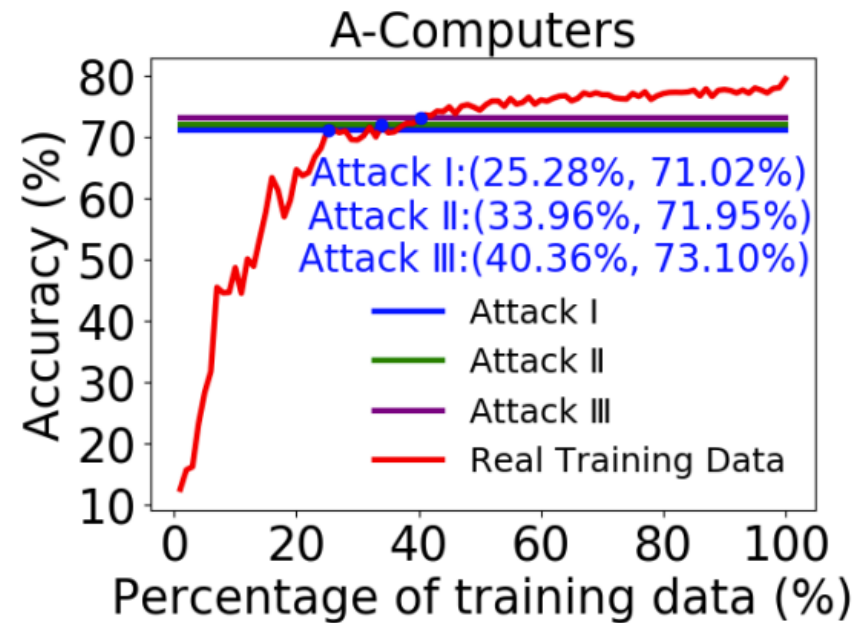
Experiments

□ Model Analysis

- Model Comparison



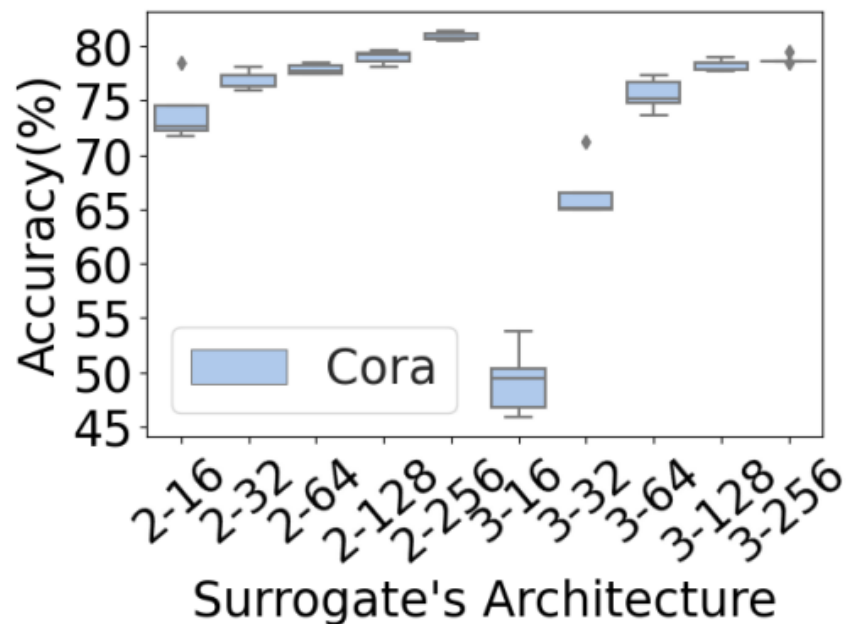
- Quantitative Analysis



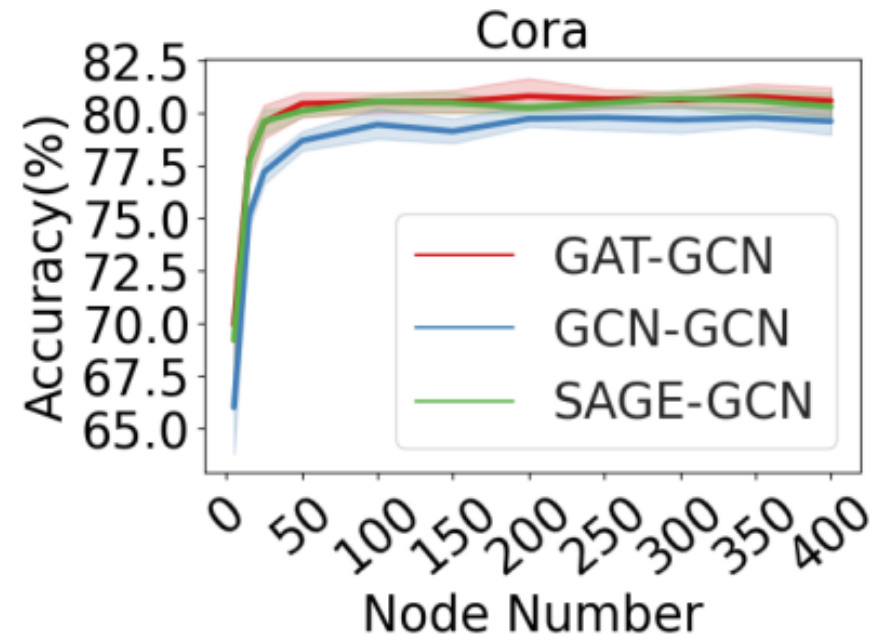
Experiments

□ Model Analysis

- Surrogate Architectures



- Generated Graph Size



Conclusion

- **Effectiveness:** StealGNN successfully performs data-free model extraction, replicating GNNs with high accuracy.
- **Innovation:** Introduced the first framework targeting GNNs, showcasing multiple strategies for generator updates
- **Security Implications:** Revealed significant vulnerabilities in GNN models, highlighting the need for stronger defenses.

Yuanxin ZHUANG

yzhuang436@connect.hkust-gz.edu.cn

Thank you !