# All Your GNN Models And Data Belong To Me *

Yang Zhang and his research group (CISPA Helmholtz Center for Information Security)
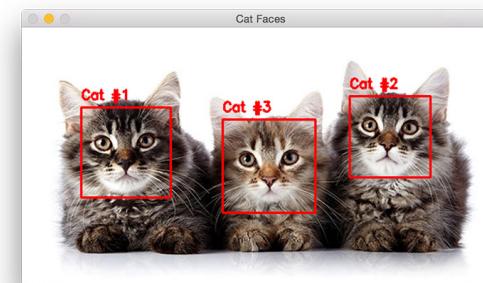Yun Shen (Spot by NetApp)
Azzedine Benameur (Spot by NetApp)

*All attacks discussed in this talk are simulated in the lab environment.

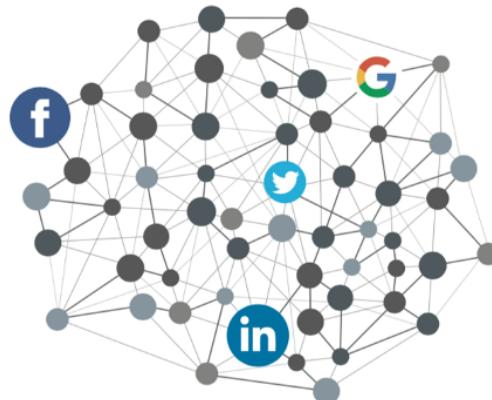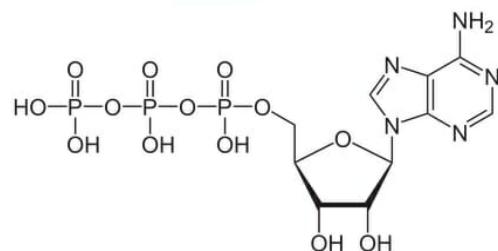# The Age of Machine Learning



Image/Text/Video/Audio

# Graph

# Graphs Are Everywhere

Graphs are combinatorial structures, have arbitrary sizes,
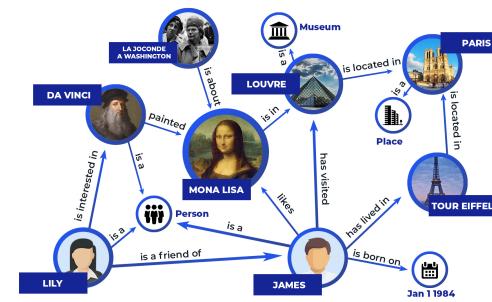and contain multi-modal information

Social Networks
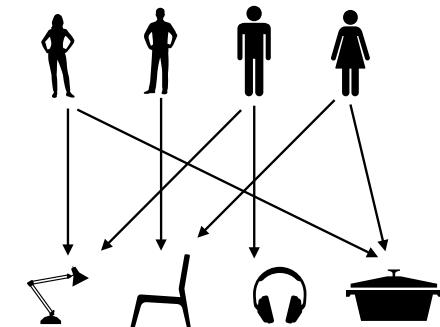
Molecules

Knowledge Graphs

User-Item Graphs

# Graph Applications Are Everywhere

Graph-based applications pervasively exist in our everyday life



Social Networks

Demographic Inference
Age group of Bob

Link Prediction
Do you (Alice) know Bob?

Molecules

Toxicity Prediction

Knowledge Graphs

Knowledge Mining

User-Item Graphs

Recommendation

We found an item you may be interested!

# Graph Neural Netwok (GNN)

- Traditional neural networks are designed for grids (e.g., images) or sequences (e.g., text)

  - CNNs for images

  - RNNs for sequences



Image

Convolved Feature

# Graph Neural Netwok (GNN)



**Input Graph**

**Message Passing**

Aggregator

Graph Convolution Network (GCN)
Graph Sample and Aggregate (GraphSAGE)
Graph Isomorphism Network (GIN)
Graph Attention Network (GAT)

**Euclidean Space**

Node Classification

Link Prediction

Demographic Inference
Age group of Bob

Link Prediction
Do you (Alice) know Bob?

# Graph Neural Netwok (GNN)

A B C D E F

- **Mean pooling**
- **Max pooling**

B A C D E F

**Node Embeddings**

**Graph Embedding**

Graph Classification

Graph Matching

Graph Visualization



Original network · Pooled network at level 1 · Pooled network at level 2 · Pooled network at level 3 · Graph classification

Toxicity Prediction

[1] Hierarchical Graph Representation Learning with Differentiable Pooling. Ying et. al.

# Graph Neural Netwok (GNN)



DeepMind Uses GNNs to Boost Google Maps ETA Accuracy by up to 50%

DeepMind unveiled a partnership with Google Maps that has leveraged advanced GNNs to improve ETA accuracy.



Cornell & Facebook AI Simplified Graph Learning Approach Outperforms SOTA GNNs

Cornell and Facebook AI "Correct and Smooth" high-accuracy graph learning method is fast to train and outperforms big Graph Neural Network models.

Insights

## Graph ML at Twitter

By Michael Bronstein
Wednesday, 2 September 2020

WHAT IS IT?

## Neo4j Graph Data Science

Neo4j Graph Data Science is a connected data analytics and machine learning platform that helps you understand the connections in big data to answer critical questions and improve predictions.

Read 5 Graph Data Science Basics

## Introducing Amazon SageMaker Support for Deep Graph Library (DGL): Build and Train Graph Neural Networks

Posted On: Dec 3, 2019

Amazon SageMaker support for the Deep Graph Library (DGL) is now available. With DGL, you can improve the prediction accuracy of recommendation, fraud detection, and drug discovery systems using Graph Neural Networks (GNNs).

kumo

how it works    about us    request early access

### From siloed tasks to an enterprise graph.

Conventional enterprise AI treats every predictive task separately in a silo. However, enterprise data represents a rich, interconnected web of business relationships, interactions, customers, transactions, and more. By leveraging the connectedness of enterprise data, Kumo enables a technical leap-frog in AI.

## Graph

# The Age of Machine Learning

# The Age of Adversarial Machine Learning

**Bloomberg**

AI Poisoning Is the Next Big Risk in Cybersecurity

25 Apr · Opinion

**IEEE Spectrum**

How Adversarial Attacks Could Destabilize Military AI Systems

**WIRED**

Even Artificial Neural Networks Can Have Exploitable 'Backdoors'

**Air Force Magazine**

Does AI Present a New Attack Surface for Adversaries?

29 Sept 2021

# Overview *

|  | Graph | GNN |
|---|---|---|
| **Security** |  | **Model extraction attack** |
| **Privacy** | **Link re-identification attack**<br><br>**Property inference attack**<br><br>**Subgraph inference attack** |  |

*All attacks discussed in this talk are simulated in the lab environment.

# Link Re-Identification Attack

|  | Graph | GNN |
|---|---|---|
| **Security** | | |
| **Privacy** | **Link re-identification attack** | Identify if two nodes are connected in the **training data** |

# Link Re-Identification Attack (Scenario 1)

**Scenario**

GNN → Posterior Scores

Security Boundary

GNN model:
**Node classification**

Attacker's **capability**:
1. posteriors of nodes (from training data) obtained from the target model

Posterior Scores

Private Information

GPU intensive

panda dog cat

panda dog cat

?

# Link Re-Identification Attack (Scenario 1)

Scenario

GNN → Posterior Scores

Security Boundary

GNN model:
**Node classification**

Attacker's **capability**:
1. posteriors of nodes (from training data) obtained from the target model

Posterior Scores

(0.70, 0.20, 0.1)

panda dog cat

(0.20, 0.70, 0.1)

panda dog cat

Posterior Similarity

>= Threshold
**link exists**

< Threshold
**link does not exist**

# Link Re-Identification Attack (Scenario 2)

**Scenario**

GNN → Posterior Scores

Security Boundary

GNN model:
**Node classification**

Attacker's **capability**:
1. posteriors of nodes (from training data) obtained from the target model
2. **have a shadow dataset**

# Link Re-Identification Attack (Scenario 2)

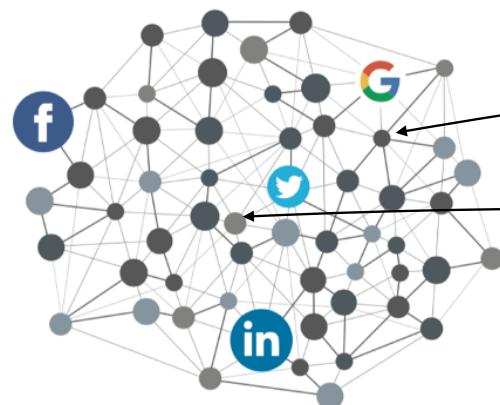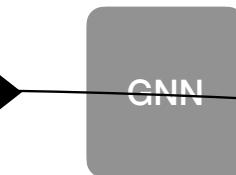**Scenario**

GNN → Posterior Scores

Security Boundary

GNN model:
**Node classification**

Attacker's **capability**:
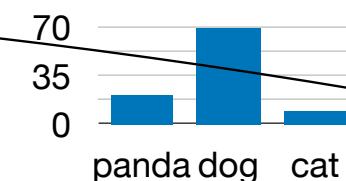1. posteriors of nodes (from training data) obtained from the target model
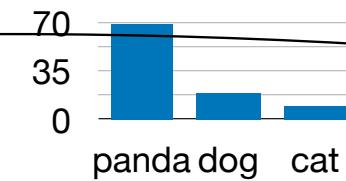2. **have a shadow dataset**

**Training with pos/neg edges**

**Shadow Dataset**

[14, 6, 9]

[10, 5, 8]

[5, 3, 12]

[8, 5, 13]

[12, 7, 8]

GNN

Posterior Scores

50
25
0
Cook  Actor  Barber  Coach

70
35
0
Cook  Actor  Barber  Coach

concat

MLP

# Link Re-Identification Attack (Scenario 2)

Training with pos/neg edges

50
25
0
Cook  Actor  Barber  Coach

70
35
0
Cook  Actor  Barber  Coach

GNN

**Distance (8)**

**Entropy (4)**

MLP

| Metrics | Definition |
|---|---|
| Cosine | $1 - \dfrac{f(u) \cdot f(v)}{\|f(u)\|_2 \, \|f(v)\|_2}$ |
| Euclidean | $\|f(u) - f(v)\|_2$ |
| Correlation | $1 - \dfrac{(f(u) - \overline{f(u)}) \cdot (f(v) - \overline{f(v)})}{\|(f(u) - \overline{f(u)})\|_2 \|(f(v) - \overline{f(v)})\|_2}$ |
| Chebyshev | $\max_i |f_i(u) - f_i(v)|$ |
| Braycurtis | $\dfrac{\sum |f_i(u) - f_i(v)|}{\sum |f_i(u) + f_i(v)|}$ |
| Manhattan | $\sum_i |f_i(u) - f_i(v)|$ |
| Canberra | $\sum_i \dfrac{|f_i(u) - f_i(v)|}{|f_i(u)| + |f_i(v)|}$ |
| Sqeuclidean | $\|f(u) - f(v)\|_2^2$ |

| Operator | Definition | Operator | Definition |
|---|---|---|---|
| Average | $\dfrac{f_i(u) + f_i(v)}{2}$ | Weighted-L1 | $|f_i(u) - f_i(v)|$ |
| Hadamard | $f_i(u) \cdot f_i(v)$ | Weighted-L2 | $|f_i(u) - f_i(v)|^2$ |

# Link Re-Identification Attack (Scenario 2)

**Shadow**

GNN

Cook Actor Barber Coach

Cook Actor Barber Coach

Training with pos/neg edges

**Distance (8)**

**Entropy (4)**

MLP

**Unified Input**

**Target**

GNN

panda dog cat

panda dog cat

**Distance (8)**

**Entropy (4)**

MLP

Testing

# Link Re-Identification Attack (Scenario 2)

**AUC**

| Target Dataset | AIDS | COX2 | DHFR | Shadow Dataset ENZYMES | PROTEINS_full | Citeseer | Cora | Pubmed |
|---|---|---|---|---|---|---|---|---|
| AIDS | - | $0.720 \pm 0.009$ | $0.690 \pm 0.005$ | $\mathbf{0.730 \pm 0.010}$ | $0.720 \pm 0.005$ | $0.689 \pm 0.019$ | $0.650 \pm 0.025$ | $0.667 \pm 0.014$ |
| COX2 | $0.755 \pm 0.032$ | - | $0.831 \pm 0.005$ | $0.739 \pm 0.116$ | $\mathbf{0.832 \pm 0.009}$ | $0.762 \pm 0.009$ | $0.773 \pm 0.008$ | $0.722 \pm 0.024$ |
| DHFR | $0.689 \pm 0.004$ | $\mathbf{0.771 \pm 0.004}$ | - | $0.577 \pm 0.044$ | $0.701 \pm 0.010$ | $0.736 \pm 0.005$ | $0.740 \pm 0.003$ | $0.663 \pm 0.010$ |
| ENZYMES | $\mathbf{0.747 \pm 0.014}$ | $0.695 \pm 0.023$ | $0.514 \pm 0.041$ | - | $0.691 \pm 0.030$ | $0.680 \pm 0.012$ | $0.663 \pm 0.009$ | $0.637 \pm 0.018$ |
| PROTEINS_full | $0.775 \pm 0.020$ | $0.821 \pm 0.016$ | $0.528 \pm 0.038$ | $0.822 \pm 0.020$ | - | $\mathbf{0.823 \pm 0.004}$ | $0.809 \pm 0.015$ | $0.809 \pm 0.013$ |
| Citeseer | $0.801 \pm 0.040$ | $0.920 \pm 0.006$ | $0.842 \pm 0.036$ | $0.846 \pm 0.042$ | $0.848 \pm 0.015$ | - | $\mathbf{0.965 \pm 0.001}$ | $0.942 \pm 0.003$ |
| Cora | $0.791 \pm 0.019$ | $0.884 \pm 0.005$ | $0.811 \pm 0.024$ | $0.804 \pm 0.048$ | $0.869 \pm 0.012$ | $\mathbf{0.942 \pm 0.001}$ | - | $0.917 \pm 0.002$ |
| Pubmed | $0.705 \pm 0.039$ | $0.796 \pm 0.007$ | $0.704 \pm 0.042$ | $0.708 \pm 0.067$ | $0.752 \pm 0.014$ | $0.883 \pm 0.006$ | $\mathbf{0.885 \pm 0.005}$ | - |

# Property/Subgraph Inference Attack

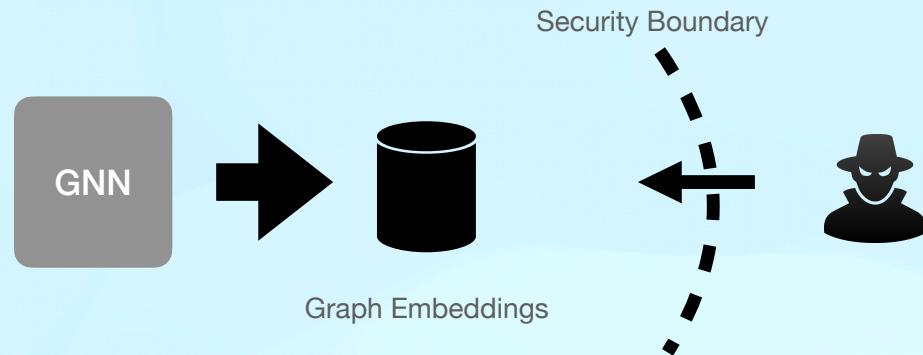|  | Graph | GNN |
|---|---|---|
| **Security** | | |
| **Privacy** | **Property inference attack**<br><br>**Subgraph inference attack** | **Infer basic graph properties** of a graph via its graph embedding<br><br>**Infer if a certain subgraph exists in a graph** via its graph embedding |

*All attacks discussed in this talk are simulated in the lab environment.

# Graph Neural Netwok (GNN)

**Node Embeddings**

A  B  C  D  E  F

- Mean pooling
- Max pooling

**Graph Embedding**

B A C D E F

Graph Classification

Graph Matching

Graph Visualization



| Original network | Pooled network at level 1 | Pooled network at level 2 | Pooled network at level 3 | Graph classification |

Toxicity Prediction

[1] Hierarchical Graph Representation Learning with Differentiable Pooling. Ying et. al.

# Property Inference Attack

**Scenario**
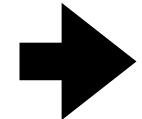
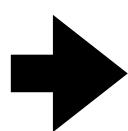Security Boundary

GNN → Graph Embeddings

GNN model:
**Graph classification**

Attacker's **capability**:
1. Embeddings of graphs (from training data) obtained from the target model
2. Can query the GNN model

Private Graph

GNN

<0.12, 0.19, 0.3, ..., 0.06>

<0.01, 0.08, 0.12, ..., 0.72>

...

<0.11, 0.09, 0.1, ..., 0.07>

# Property Inference Attack

Security Boundary
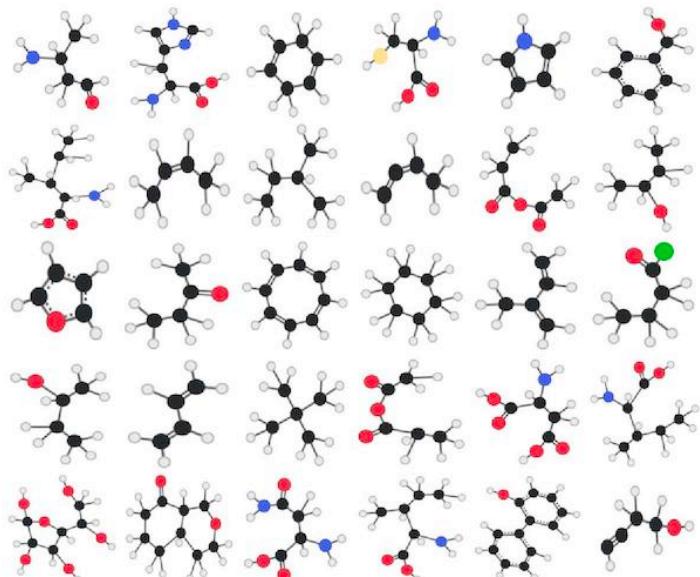
GNN

Graph Embeddings

GNN model:
**Graph classification**

Attacker's **capability**:
1. Embeddings of graphs (from training data) obtained from the target model
2. Can query the GNN model

<0.12, 0.19, 0.3, ..., 0.06>

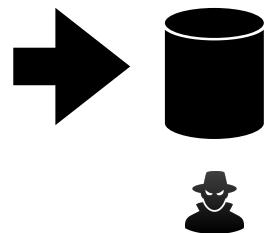<0.01, 0.08, 0.12, ..., 0.72>

...

<0.11, 0.09, 0.1, ..., 0.07>

This is a graph with ~4 nodes

infer

# Property Inference Attack

Security Boundary

## Scenario

GNN → Graph Embeddings

GNN model:
**Graph classification**

Attacker's **capability**:
1. Embeddings of graphs (from training data) obtained from the target model
2. Can query the GNN model

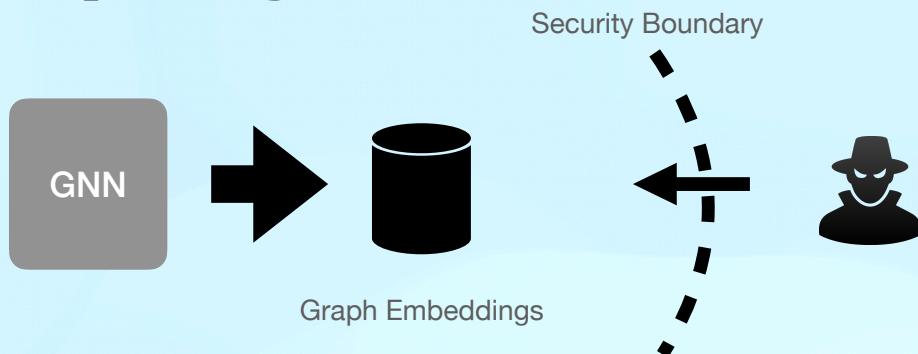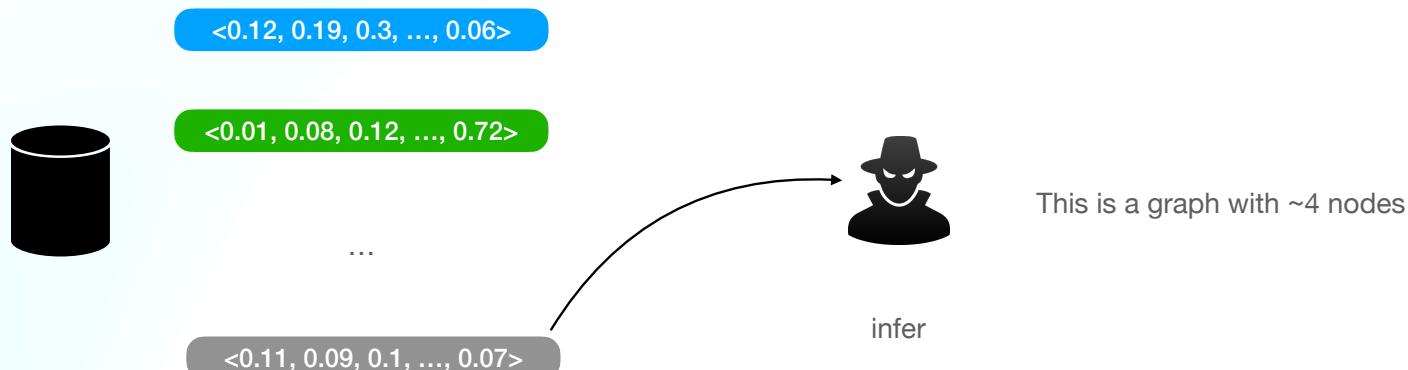**Remote access**

Auxiliary graphs

GNN

Graph Embeddings

<0.12, 0.19, 0.3, …, 0.06>

<0.01, 0.08, 0.12, …, 0.72>

Attack Model

Estimated

Ground Truth

**Cross-entropy loss**

40
0
[1-2][3-4][5-6] [7+]

100
0
[1-2] [3-4] [5-6] [7+]

60
0
[1-2]      [5-6]

100
0
[1-2] [3-4] [5-6] [7+]

**Local environment**

# Property Inference Attack

**Scenario**



Security Boundary
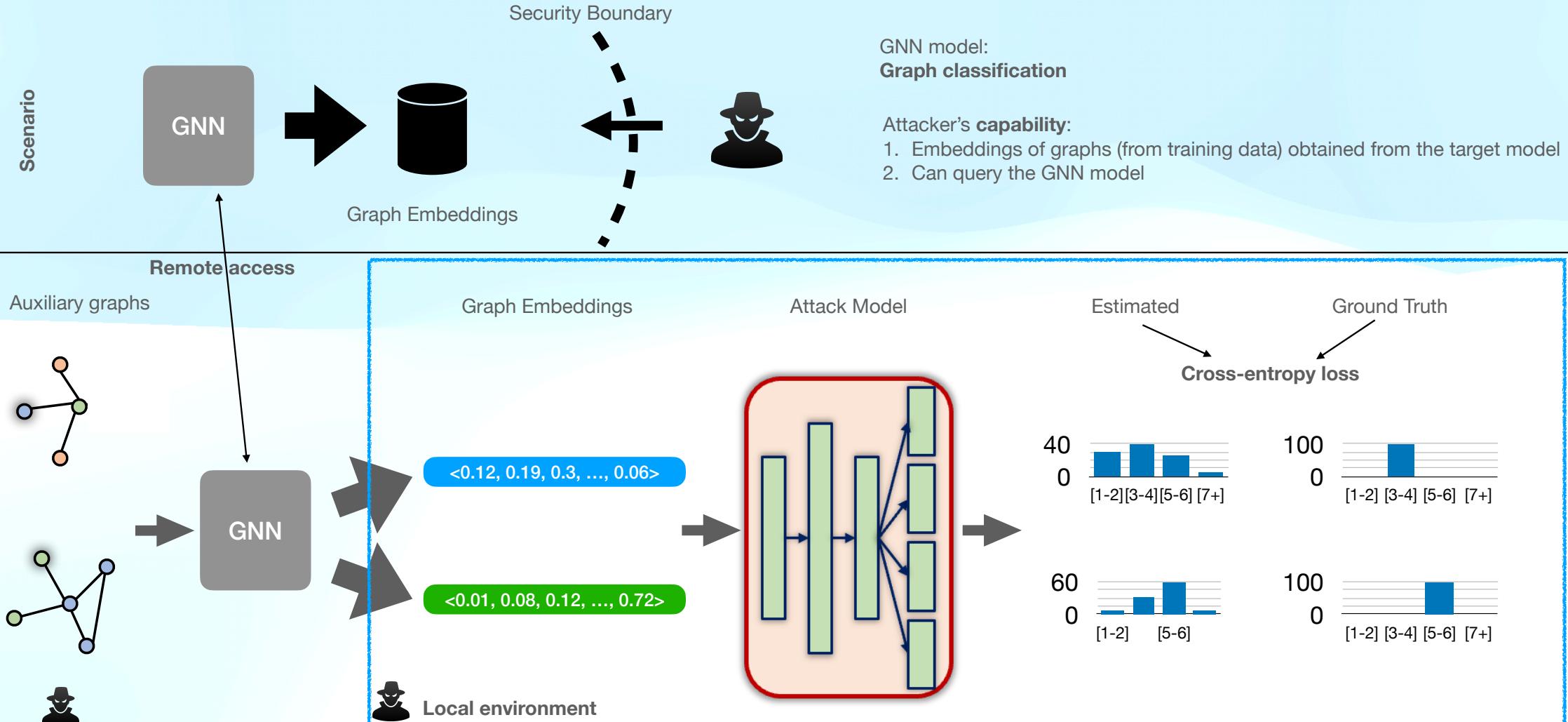
GNN → Graph Embeddings

GNN model:
**Graph classification**

Attacker's **capability**:
1. Embeddings of graphs (from training data) obtained from the target model
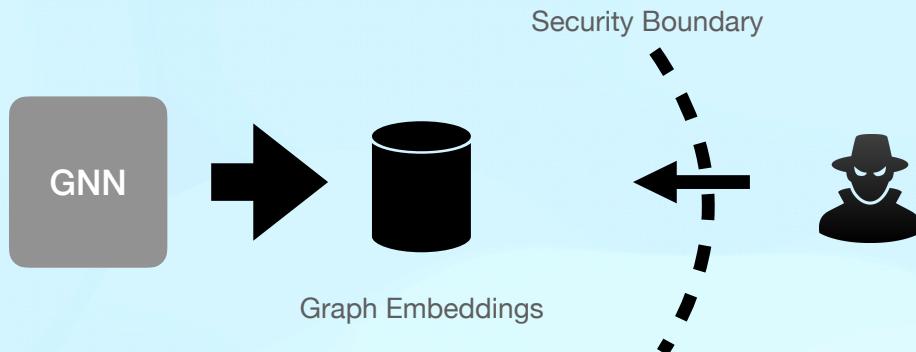2. Can query the GNN model

Graph Embeddings      Attack Model      Estimated

<0.11, 0.09, 0.1, ..., 0.07>

60
0
[1-2] [3-4] [5-6] [7+]

This is a graph with ~4 nodes

# Subgraph Inference Attack

Security Boundary

GNN

Graph Embeddings

GNN model:
**Graph classification**

Attacker's **capability**:
1. Embeddings of graphs (from training data) obtained from the target model
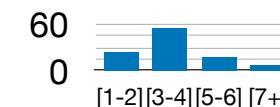2. Can query the GNN model

<0.12, 0.19, 0.3, ..., 0.06>

infer

This graph contains at least one

# Subgraph Inference Attack



Security Boundary

**Scenario**

GNN → Graph embeddings

GNN model:
**Graph classification**

Attacker's **capability**:
1. Embeddings of graphs (from training data) obtained from the target model
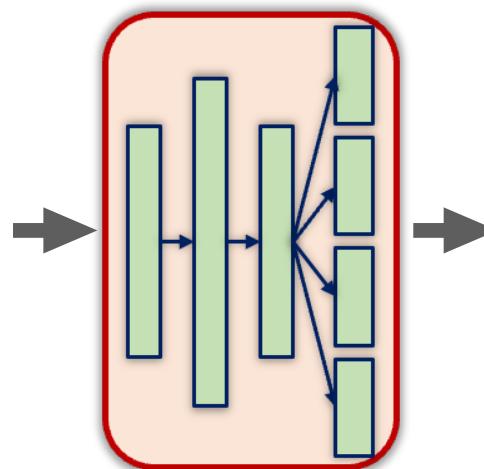2. Can query the GNN model

Remote access

Auxiliary Graphs

Graph Embeddings

GNN

<0.12, 0.19, 0.3, …, 0.06>
<0.01, 0.08, 0.12, …, 0.72>

+

**negative pair**
<0.12, 0.19, 0.3, …, 0.06>
<0.01, 0.08, 0.12, …, 0.72>

**positive pair**
<0.01, 0.08, 0.12, …, 0.72>
<0.01, 0.08, 0.12, …, 0.72>

Attack model

Subgraph embeddings

GNN → <0.01, 0.08, 0.12, …, 0.72>

# Subgraph Inference Attack

Security Boundary

GNN

Graph embeddings

GNN model:
**Graph classification**

Attacker's **capability**:
1. Embeddings of graphs (from training data) obtained from the target model
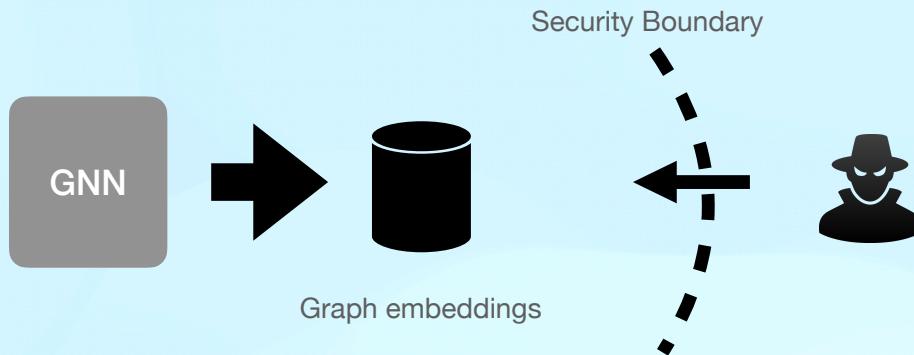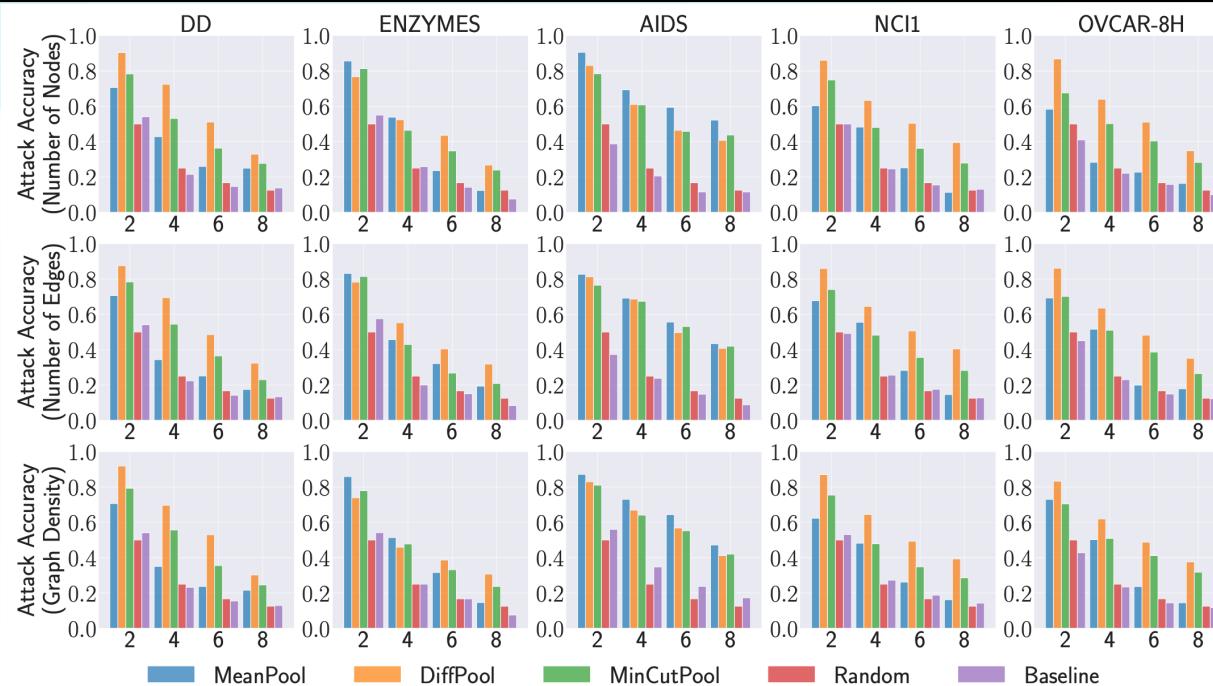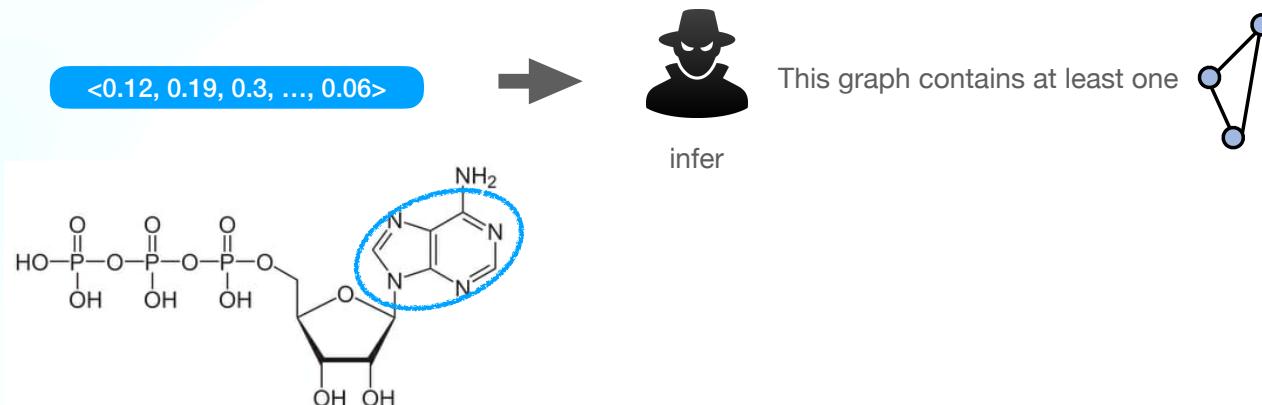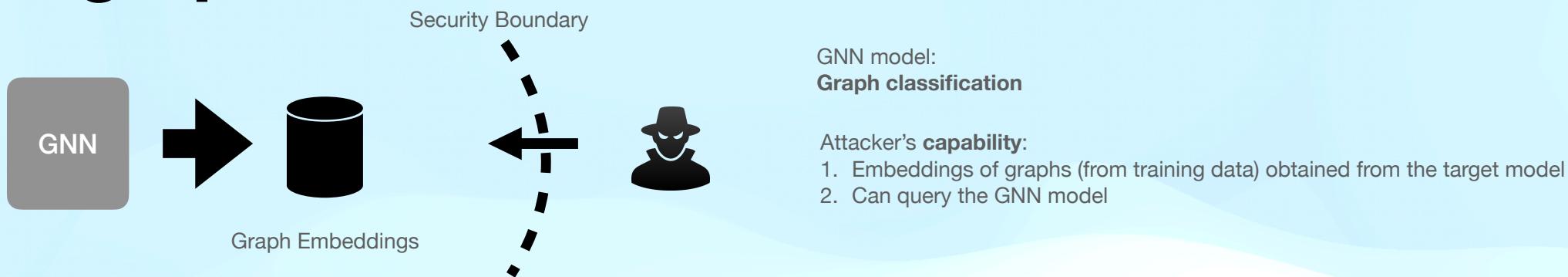2. Can query the GNN model

---

Graph Embeddings

Attack Model

<0.11, 0.09, 0.1, ..., 0.07>

<0.11, 0.09, 0.1, ..., 0.07>
<0.01, 0.08, 0.12, ..., 0.72>

80
0
Positive    Negative

GNN

<0.01, 0.08, 0.12, ..., 0.72>

# Subgraph Inference Attack

Scenario



Security Boundary

GNN → Graph embeddings

GNN model:
**Graph classification**

Attacker's **capability**:
1. Embeddings of graphs (from training data) obtained from the target model
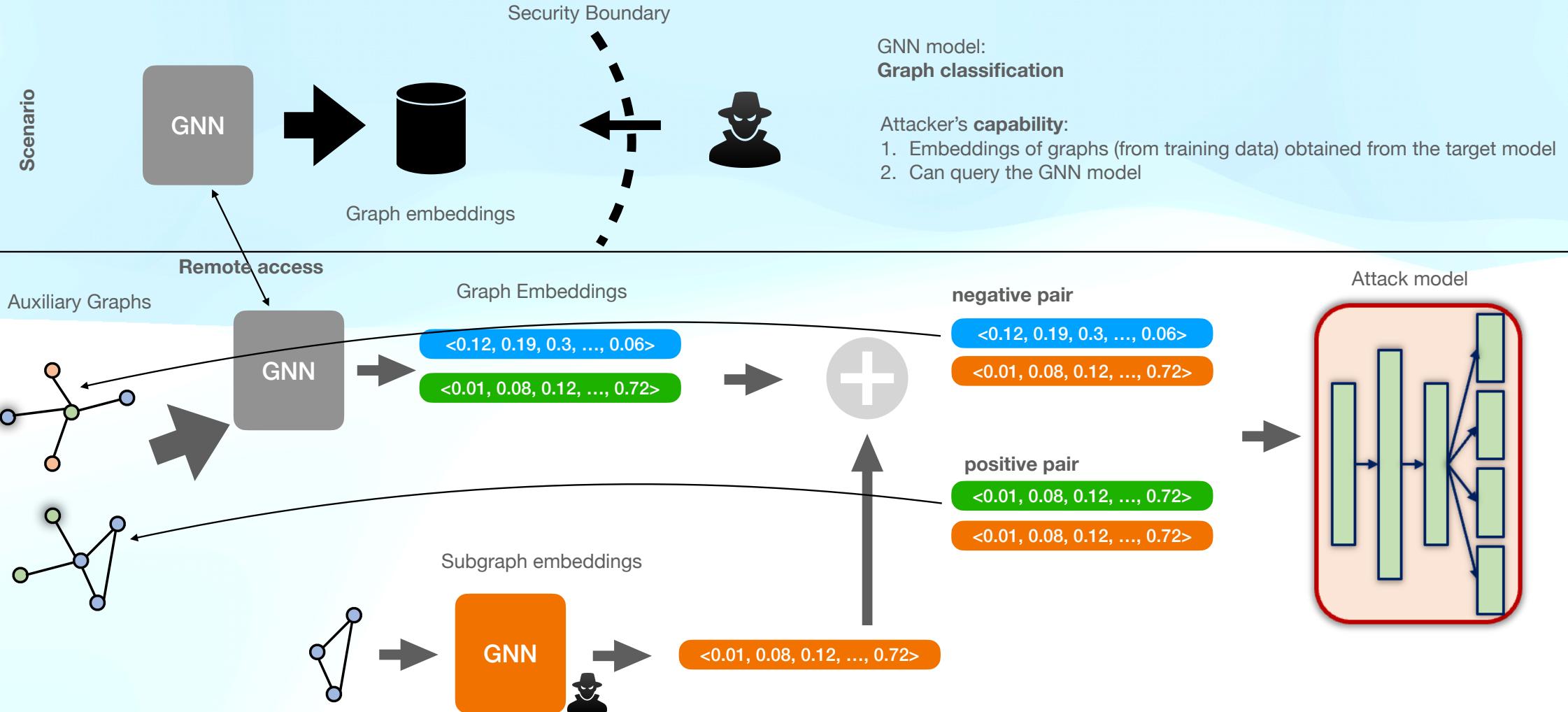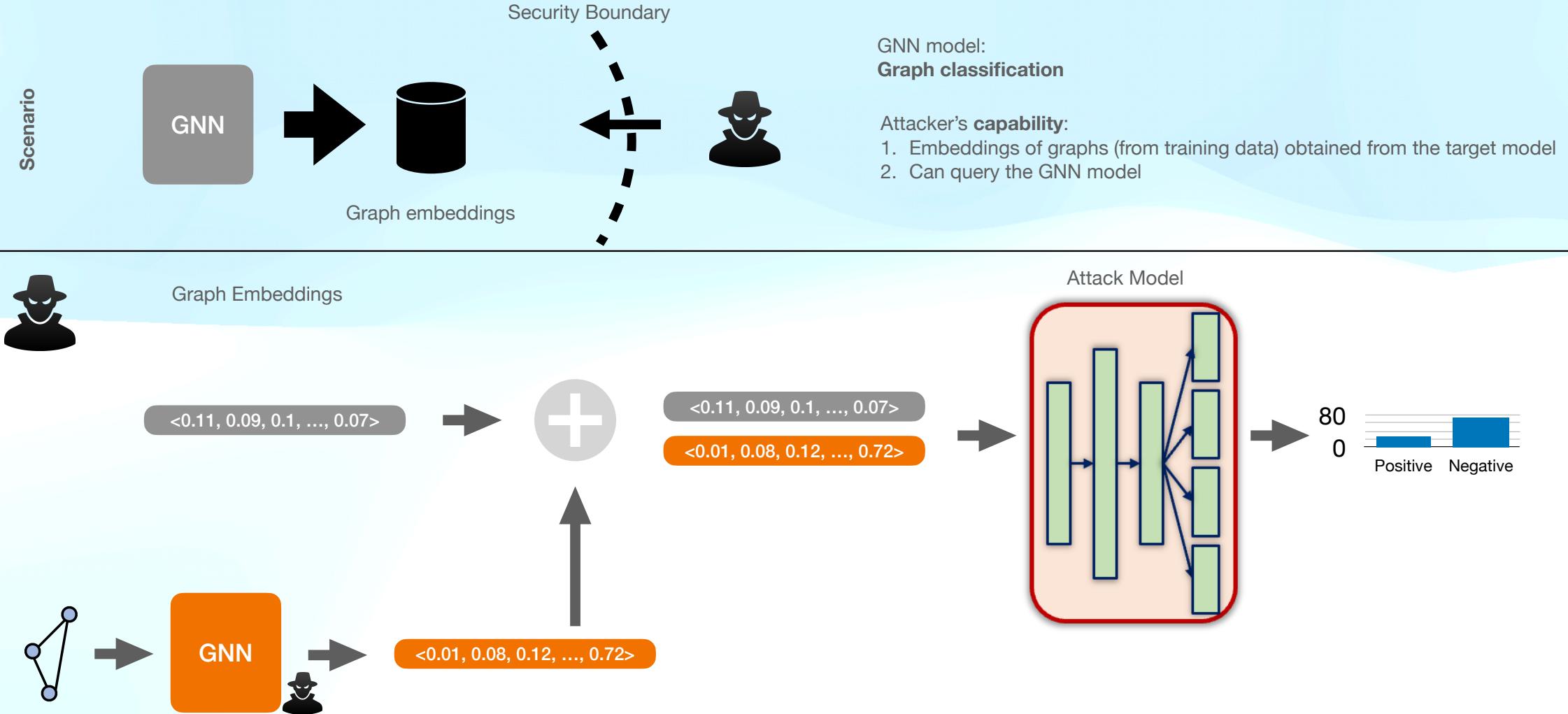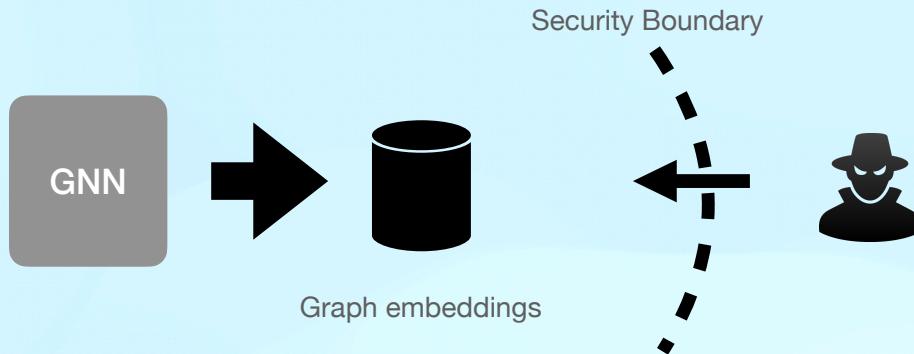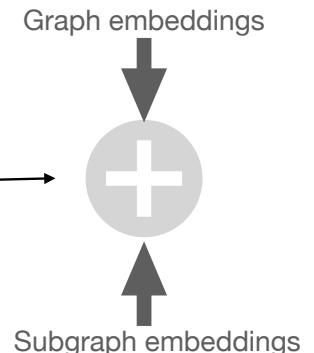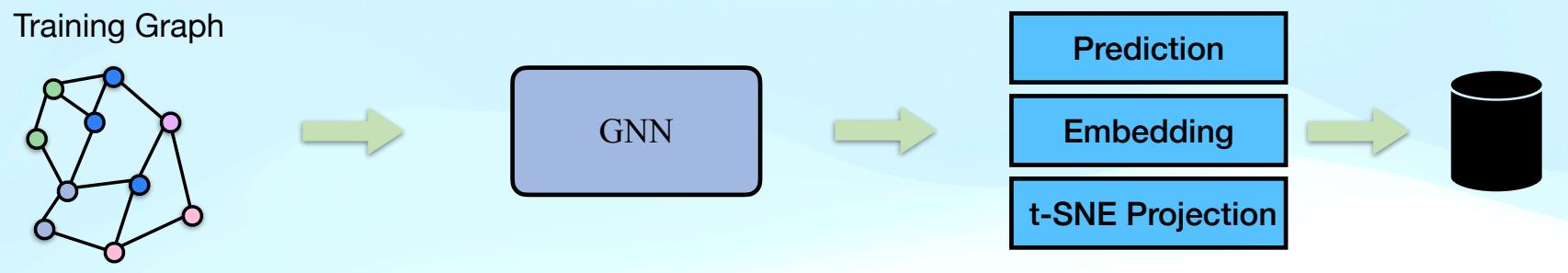2. Can query the GNN model

**AUC**

| Dataset | 0.8 Concat | 0.8 EDist | 0.8 EDiff | 0.6 Concat | 0.6 EDist | 0.6 EDiff | 0.4 Concat | 0.4 EDist | 0.4 EDiff | 0.2 Concat | 0.2 EDist | 0.2 EDiff |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DD | $0.53 \pm 0.01$ | $0.81 \pm 0.06$ | $\mathbf{0.88 \pm 0.01}$ | $0.51 \pm 0.01$ | $0.79 \pm 0.04$ | $\mathbf{0.87 \pm 0.01}$ | $0.52 \pm 0.01$ | $0.79 \pm 0.02$ | $\mathbf{0.85 \pm 0.01}$ | $0.50 \pm 0.02$ | $0.71 \pm 0.08$ | $\mathbf{0.80 \pm 0.00}$ |
| ENZYMES | $0.49 \pm 0.02$ | $0.63 \pm 0.10$ | $\mathbf{0.88 \pm 0.03}$ | $0.52 \pm 0.03$ | $0.71 \pm 0.10$ | $\mathbf{0.88 \pm 0.03}$ | $0.54 \pm 0.02$ | $0.56 \pm 0.07$ | $\mathbf{0.86 \pm 0.01}$ | $0.48 \pm 0.02$ | $0.53 \pm 0.03$ | $\mathbf{0.78 \pm 0.01}$ |
| AIDS | $0.51 \pm 0.01$ | $0.53 \pm 0.04$ | $\mathbf{0.78 \pm 0.04}$ | $0.55 \pm 0.01$ | $0.51 \pm 0.02$ | $\mathbf{0.76 \pm 0.05}$ | $0.54 \pm 0.01$ | $0.51 \pm 0.03$ | $\mathbf{0.73 \pm 0.06}$ | $0.56 \pm 0.02$ | $0.50 \pm 0.00$ | $\mathbf{0.76 \pm 0.05}$ |
| NCI1 | $0.51 \pm 0.00$ | $0.51 \pm 0.02$ | $\mathbf{0.70 \pm 0.06}$ | $0.49 \pm 0.02$ | $0.52 \pm 0.01$ | $\mathbf{0.67 \pm 0.06}$ | $0.50 \pm 0.01$ | $0.51 \pm 0.01$ | $\mathbf{0.64 \pm 0.03}$ | $0.49 \pm 0.01$ | $0.51 \pm 0.01$ | $\mathbf{0.64 \pm 0.00}$ |
| OVCAR-8H | $0.54 \pm 0.01$ | $0.63 \pm 0.12$ | $\mathbf{0.89 \pm 0.02}$ | $0.50 \pm 0.04$ | $0.69 \pm 0.09$ | $\mathbf{0.88 \pm 0.02}$ | $0.51 \pm 0.03$ | $0.74 \pm 0.02$ | $\mathbf{0.84 \pm 0.01}$ | $0.54 \pm 0.01$ | $0.60 \pm 0.13$ | $\mathbf{0.82 \pm 0.02}$ |

Graph embeddings

Subgraph embeddings

# Analysis

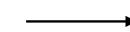Training Graph

GNN

Prediction

Embedding

t-SNE Projection

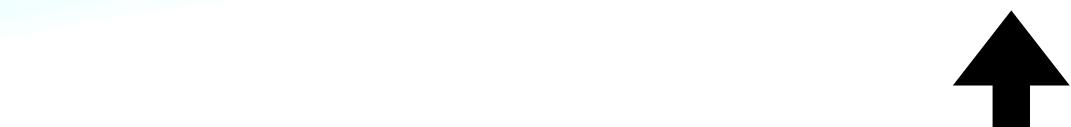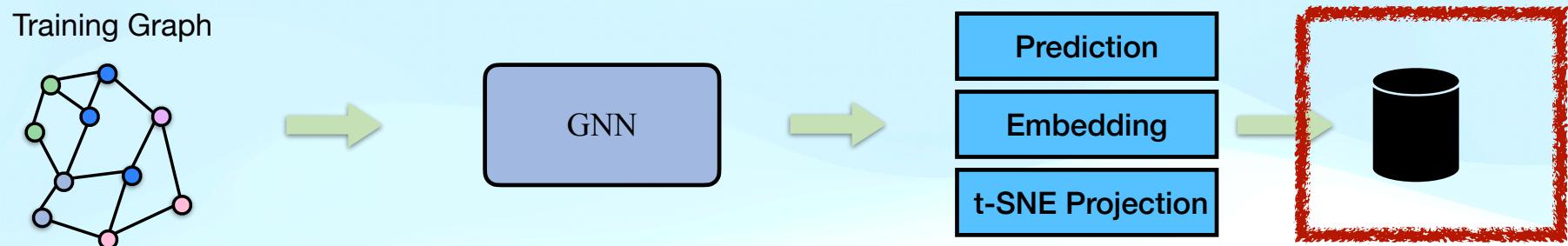Link re-identification attack → Training graph's node posterior scores

Property inference attack → Training graph's graph embeddings

# Analysis

Training Graph

GNN

Prediction

Embedding

t-SNE Projection

Link re-identification attack ✗→ Training graph's node posterior scores

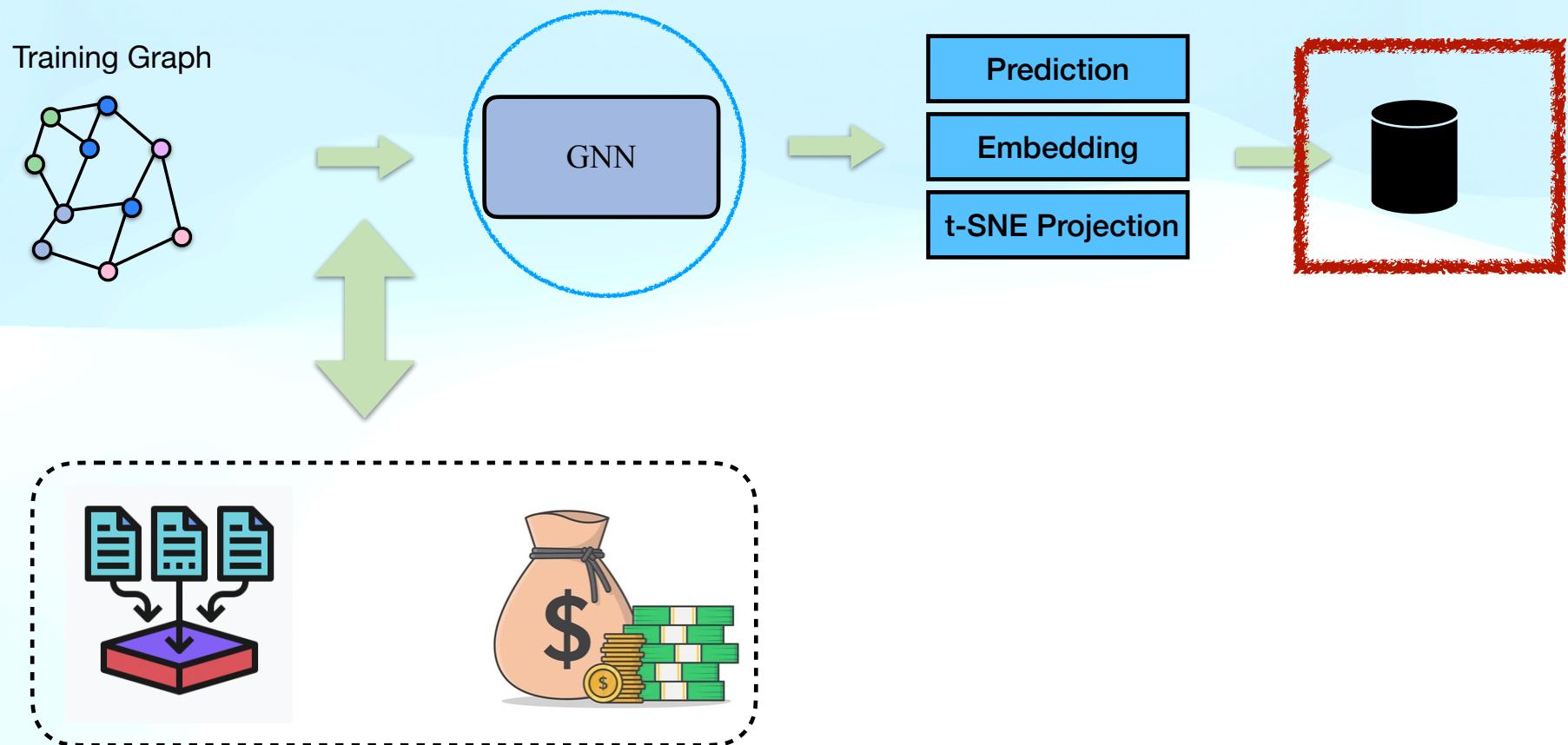Property inference attack ✗→ Training graph's graph embeddings

# Takeaways (1)

- **Secure** your infrastructure

- **Audit** your GNN-based machine learning pipeline

# What Is Next?

Training Graph

GNN

Prediction

Embedding

t-SNE Projection

# Overview *

|  | Graph | GNN |
|---|---|---|
| **Security** |  | **Model extraction attack** |
| **Privacy** |  | **Faithfully replicate the GNN functionality** |

# Model Stealing Attack

Training Graph

GNN

Prediction

Embedding

t-SNE Projection

Security Boundary

Query Graph

Downstream Applications

Customer

# Model Stealing Attack

**Scenario**

Security Boundary

GNN

GNN model:
**Node classification**

Attacker's **capability**:
1. Can query the GNN model via publicly accessible API

# Model Stealing Attack

# Model Stealing Attack

# Model Stealing Attack



**$G_O$**

Build

Target Model

**$M_T$**

Response

| **R** | | |
|---|---|---|
| **H** | **Θ** | **Υ** |
| Embedding | Prediction | t-SNE Projection |

Security Boundary

Node-level Query

...

**IDGL framework [1] / $k$NN**

**$X_Q$**

**$G_Q$**

Learn Discrete Graph Structure

[1] Chen, Yu, Lingfei Wu, and Mohammed Zaki. "Iterative deep graph learning for graph neural networks: Better and robust node embeddings." *Advances in neural information processing systems* 33 (2020)

# Model Stealing Attack



Security Boundary

**R**

| **H** | **Θ** | **Υ** |
|---|---|---|
| Embedding | Prediction | t-SNE Projection |

$\mathbf{G}_O$

Build

Target Model

$\mathbf{M}_T$

Response

Node-level Query

**IDGL framework [1] / $k$NN**

$\mathbf{X}_Q$

$\mathbf{G}_Q$

❶ Learn Discrete Graph Structure

$\hat{\mathbf{H}}_Q = \mathscr{F}(\mathbf{X}_Q, \mathbf{A}_Q)$

$\mathbf{M}_S$

$+$

$\mathscr{O}(\hat{\mathbf{H}}_Q)$

❷ Learn Surrogate Model

[1] Chen, Yu, Lingfei Wu, and Mohammed Zaki. "Iterative deep graph learning for graph neural networks: Better and robust node embeddings." *Advances in neural information processing systems* 33 (2020)

# Model Stealing Attack



Euclidean Space

$\mathbf{G}_Q$      $\mathbf{M}_T$

Embedding

Prediction

t-SNE Projection

$\mathbf{M}_S$      $\mathbf{G}_Q$

$$\hat{\mathbf{H}}_Q = \mathcal{F}(\mathbf{X}_Q, \mathbf{A}_Q)$$

$$\mathcal{L}_R = \frac{1}{n_Q}\|\hat{\mathbf{H}}_Q - \mathbf{R}\|_{2}$$

$$\hat{\mathbf{H}}_Q = \mathscr{F}(\mathbf{X}_Q, \mathbf{A}_Q)$$

$\mathbf{M}_S$

$+$

$\mathscr{O}(\hat{\mathbf{H}}_Q)$

❷ Learn Surrogate Model

[1] Chen, Yu, Lingfei Wu, and Mohammed Zaki. "Iterative deep graph learning for graph neural networks: Better and robust node embeddings." *Advances in neural information processing systems* 33 (2020)

# Model Stealing Attack



Euclidean Space

$\mathbf{G}_Q$

$\mathbf{M}_T$

Embedding

Prediction

t-SNE Projection

1  1

decision boundary

$\mathbf{M}_S$

$\mathbf{G}_Q$

2  2

$$\mathcal{L}_P = -\frac{1}{n_Q} \sum_{v \in \mathbf{G}_Q} \sum_{i \in |\mathbf{C}_Q|} c_i log[O(\mathbf{h}_v)_i]$$

$\hat{\mathbf{H}}_Q = \mathscr{F}(\mathbf{X}_Q, \mathbf{A}_Q)$

$\mathbf{M}_S$

$\mathscr{O}(\hat{\mathbf{H}}_Q)$

❷ Learn Surrogate Model

[1] Chen, Yu, Lingfei Wu, and Mohammed Zaki. "Iterative deep graph learning for graph neural networks: Better and robust node embeddings." *Advances in neural information processing systems* 33 (2020)

# Model Stealing Attack



Euclidean Space

$\mathbf{G}_Q$

$\mathbf{M}_T$

Embedding

Prediction

t-SNE Projection

$\mathbf{M}_S$

$\mathbf{G}_Q$

conduct the attack without knowing the target model's architecture

# Model Stealing Attack

# Model Stealing Attack



Euclidean Space

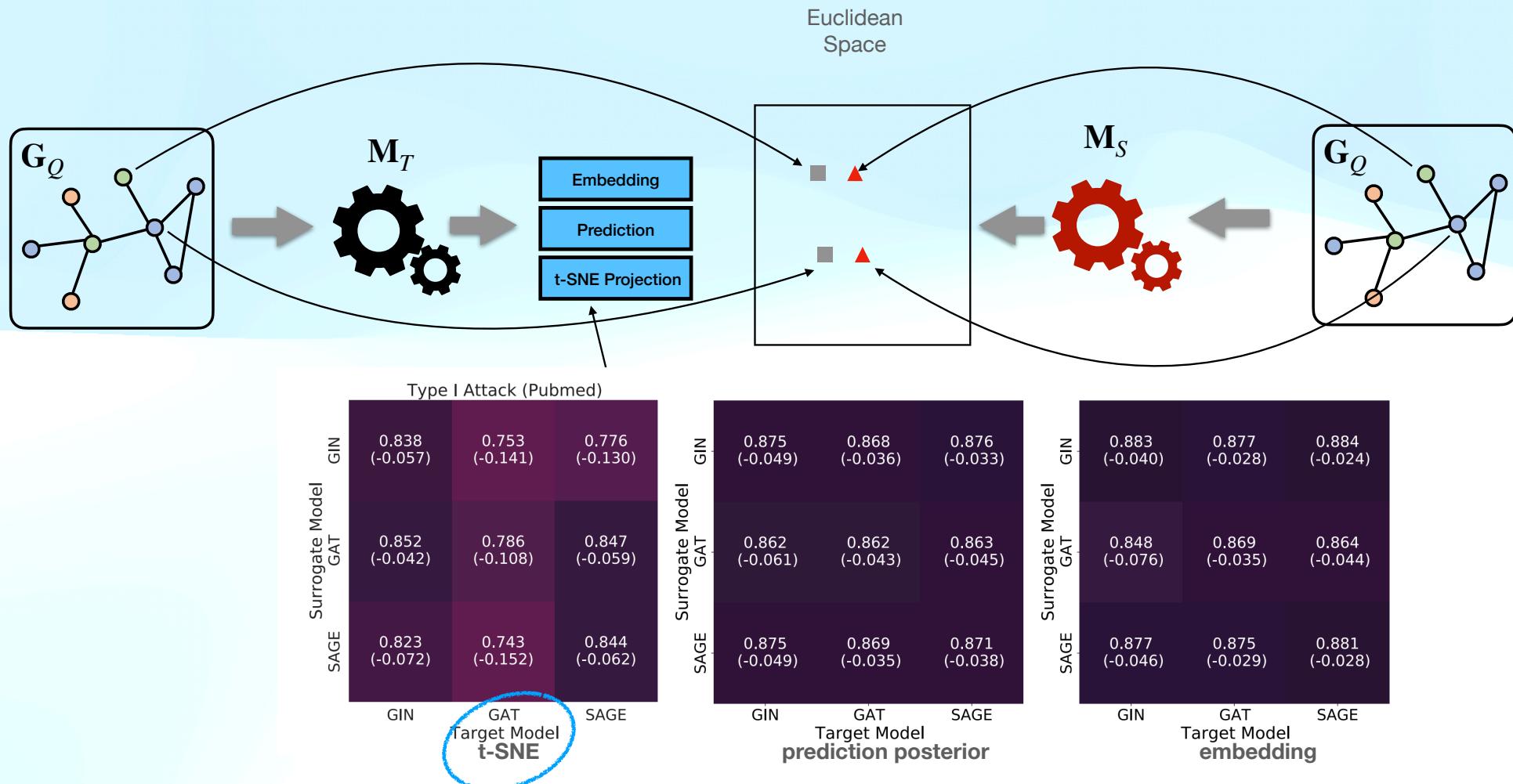$\mathbf{G}_Q$     $\mathbf{M}_T$     Embedding     Prediction     t-SNE Projection     $\mathbf{M}_S$     $\mathbf{G}_Q$

2 dimensional t-SNE projection can be the new attack surface

# Model Stealing Attack

# Takeaways (2)

- **<u>Secure</u>** your infrastructure

- **<u>Audit</u>** your GNN-based machine learning pipeline

- **<u>Monitor</u>** your model logs for anomalies

- **<u>Evaluate the security and privacy posture</u>** of your Graph Neural Network (GNN) models

# Code

- **Link re-identification attack**

  https://github.com/xinleihe/link_stealing_attack

- **Property/Subgraph inference attack**

  https://github.com/Zhangzhk0819/GNN-Embedding-Leaks

- **Model stealing attack**

  https://github.com/xinleihe/GNNStealing

# Thank You

Yang Zhang and his research group
CISPA Helmholtz Center for Information Security
**zhang@cispa.de**

Azzedine Benameur and Yun Shen
Spot by NetApp
**{Azzedine.Benameur, Yun.Shen}@netapp.com**