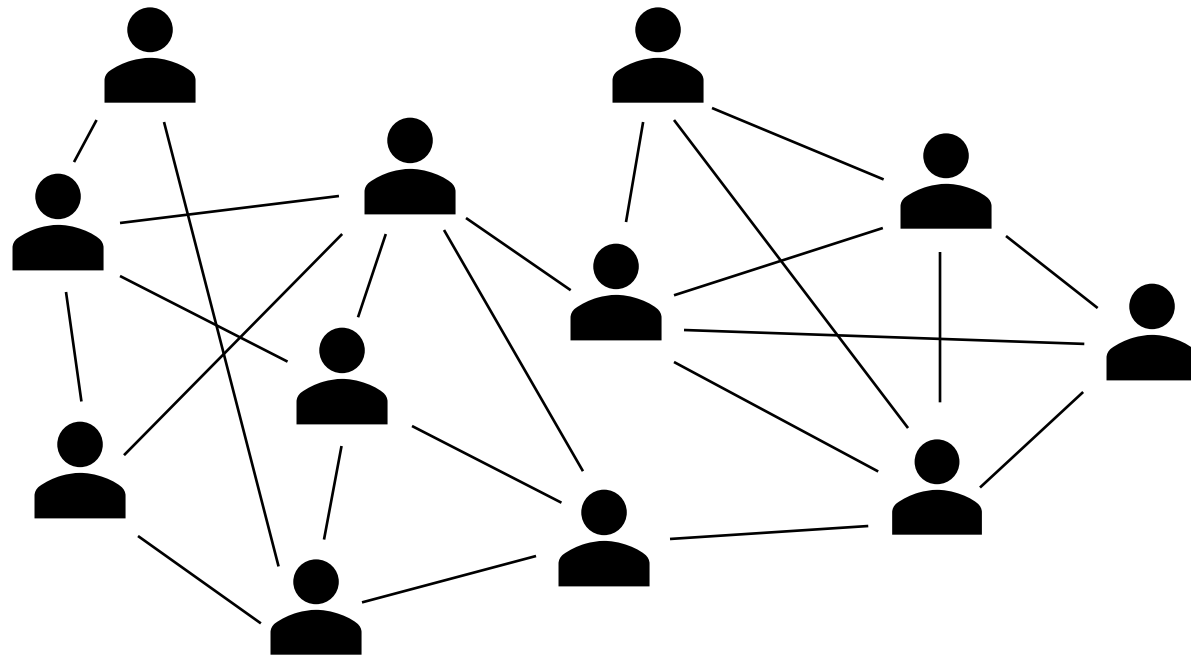# Sybil Detection using GNNs
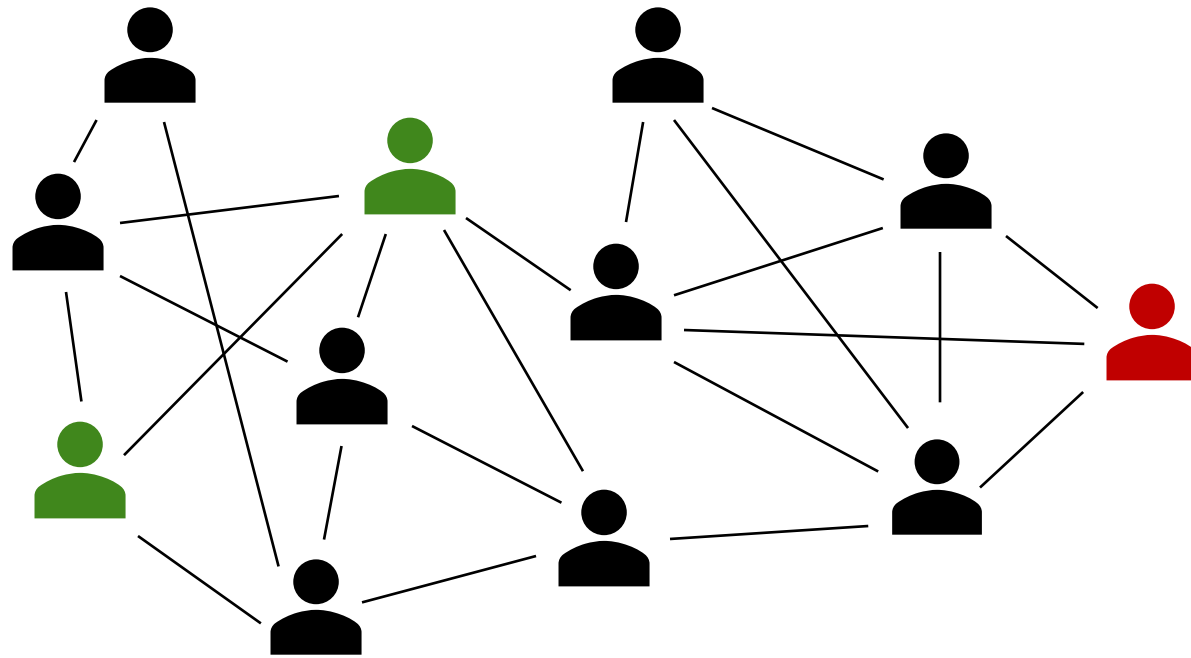
**Master Thesis Presentation**
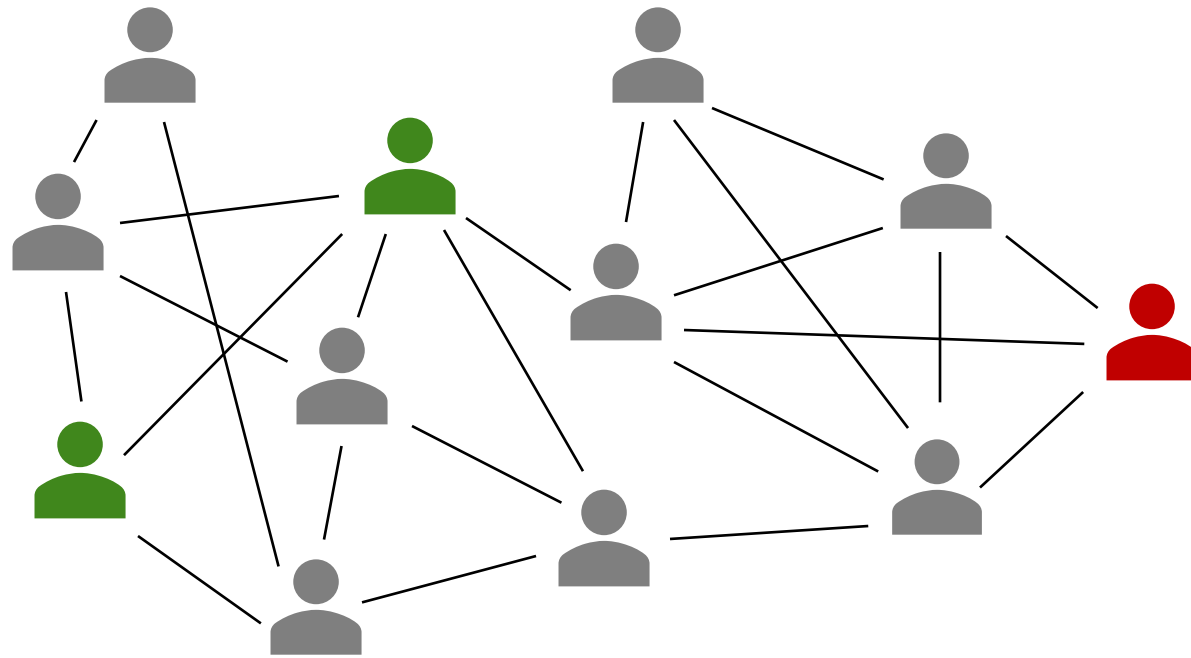
**Stuart Heeb**
September 12th, 2024
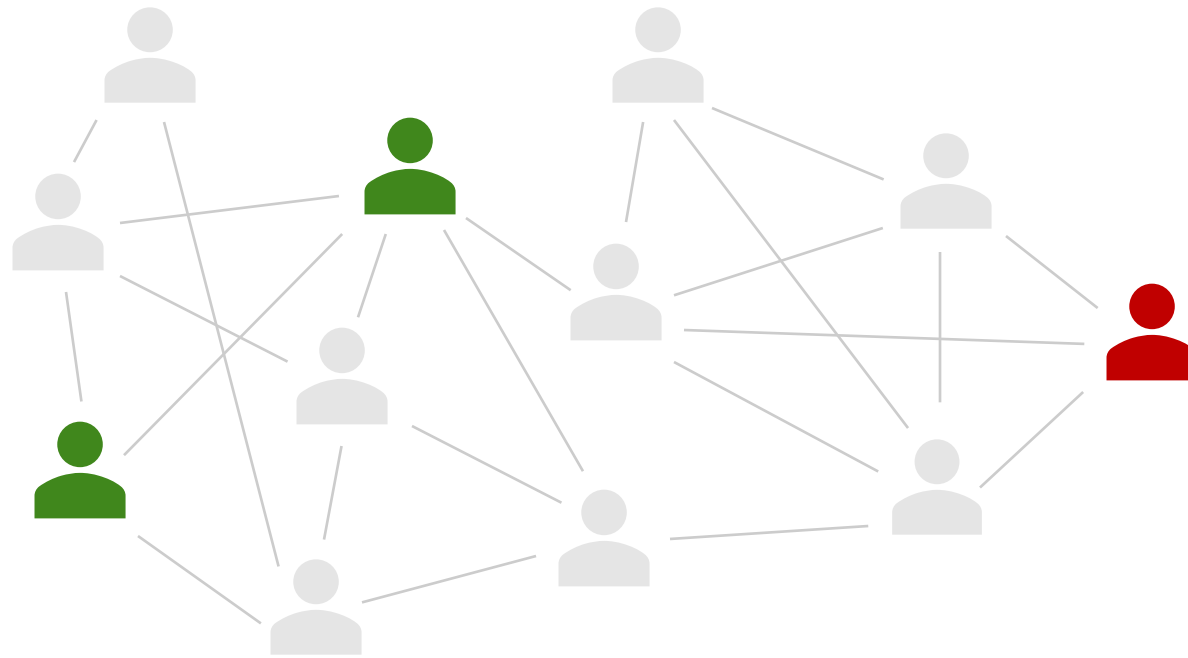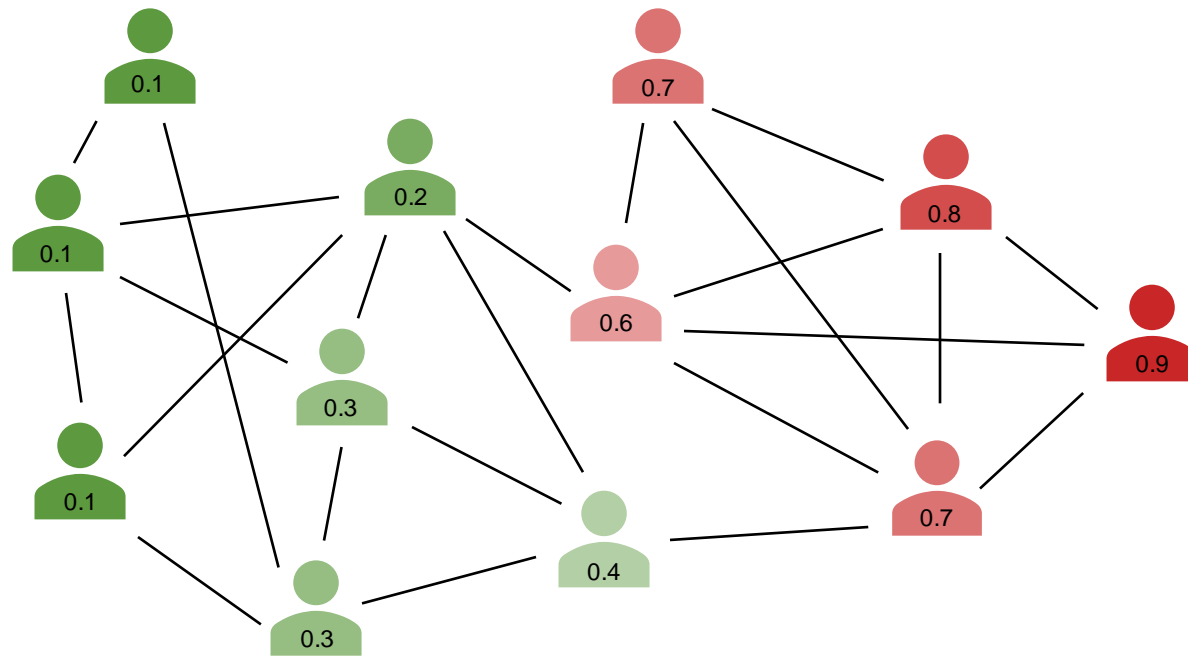
**… Sybil Detection Algorithm**

$$P(x_u = 1 \mid \mathbf{x}_V)$$

Attack edges

Sybil region

Honest region

Honest region

Attack edges

Sybil region

# Why are Sybils a Problem?

Spread misinformation

Degrade trust

Influence public opinion

# Modern Sybils



**AI-generated**

**Actual photograph**

# Which picture is real?
## As in, an image of the real world taken by a physical camera?

# **Structure-based** Sybil Detection

- Content becoming harder to (even manually) distinguish

- Other advantages

  – Privacy concerns

  – Enhanced generalizability

# Background & Related Work

## SYBILRANK

2012

- Random Walks (RW)
- Uses only honest labels
- Computationally efficient
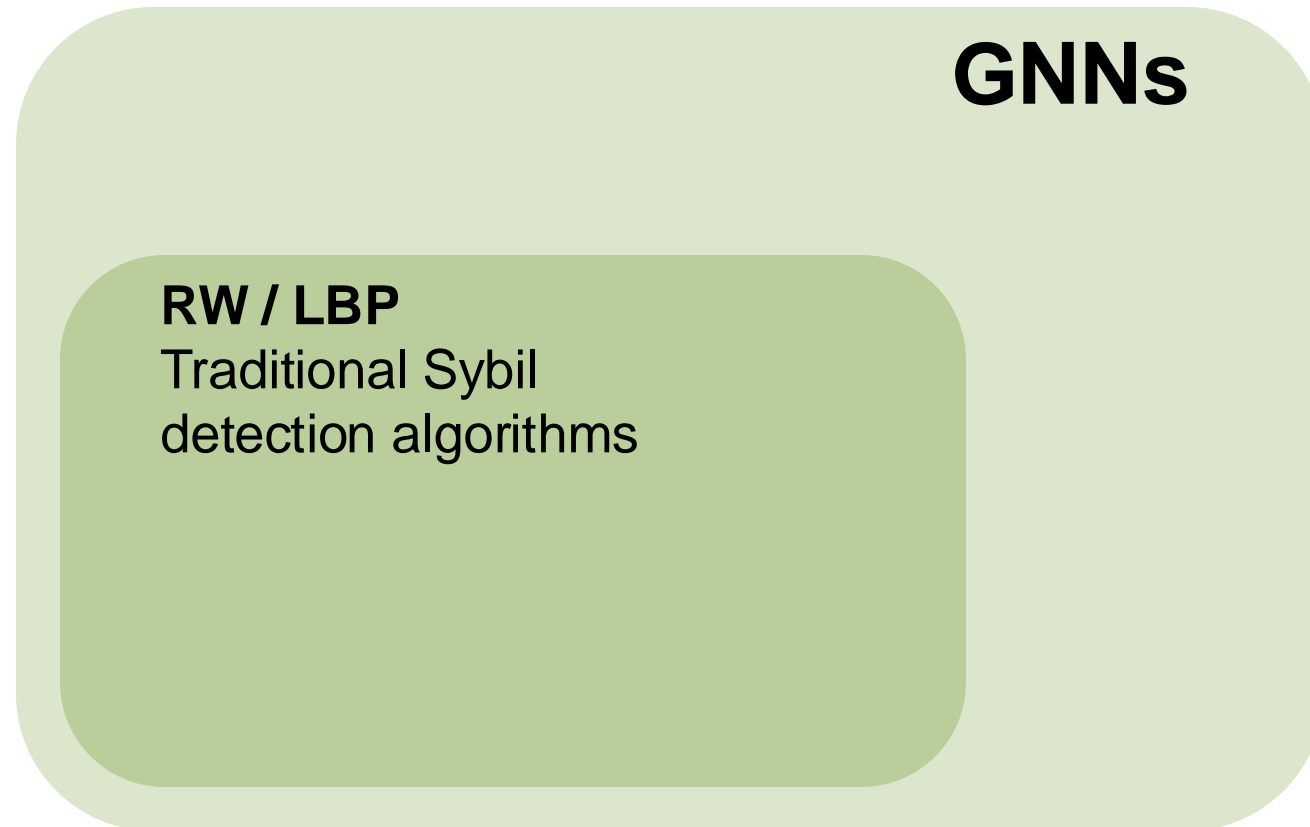
## SYBILBELIEF

2014

- Loopy Belief Propagation (LBP)
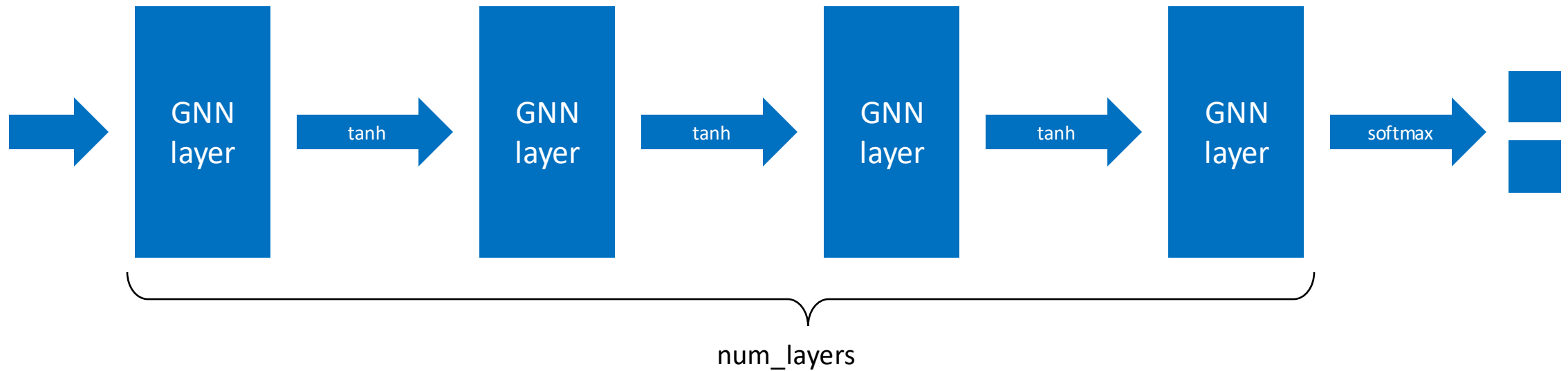- Uses both honest and sybil labels

## SYBILSCAR

2019

- Local rule-based propagation
- Uses both honest and sybil labels
- Combines benefits of RW and LBP

# Why use GNNs?



**GNNs**

**RW / LBP**
Traditional Sybil
detection algorithms

**GNNs as a «one size fits all»?**
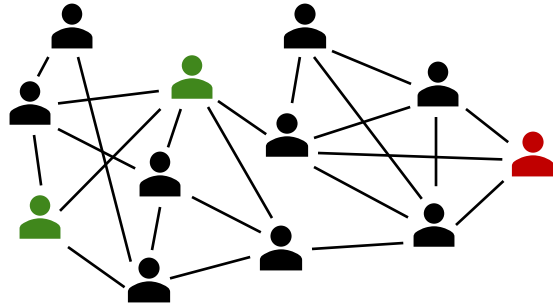
# GNNs for Sybil Detection



SYBILGCN    SYBILRGCN    SYBILGAT

# GNNs for Sybil Detection
**Approaches**

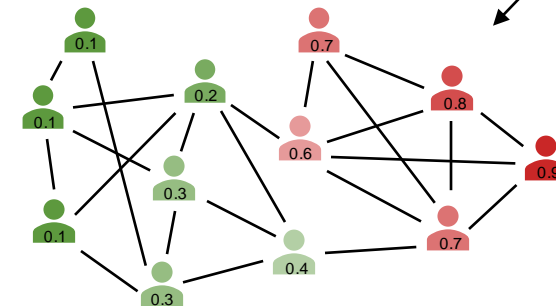

Train and **directly** predict

Train, **then** predict

*e.g., sampled subgraph*

*could be multiple graphs*

**Without pre-training**
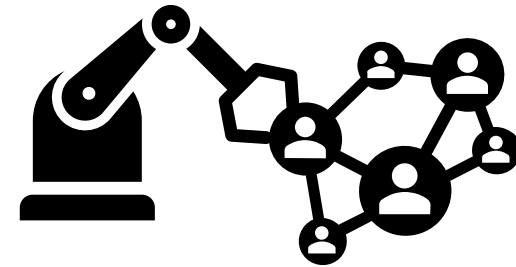
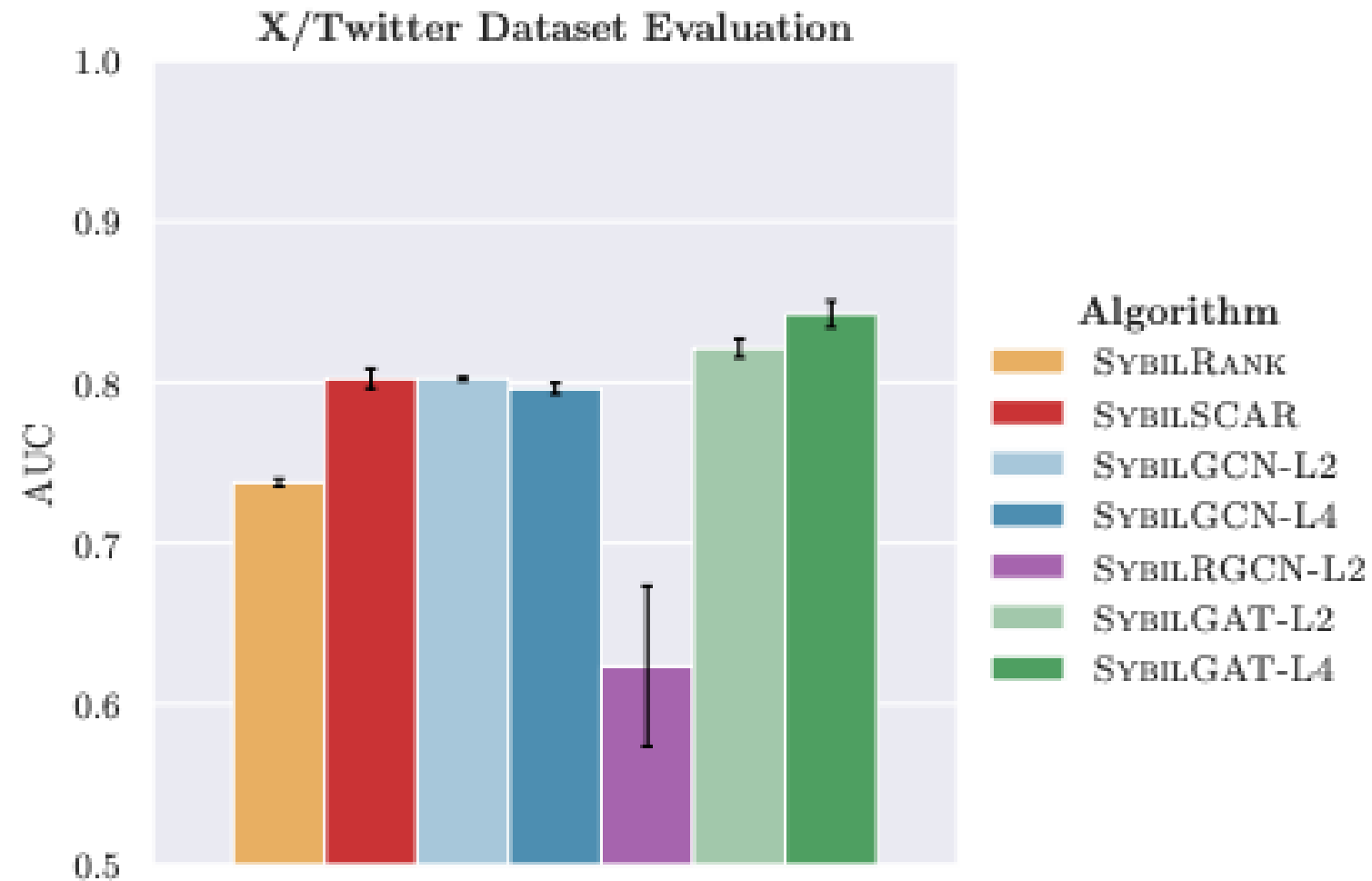**With pre-training**

# Data



270K nodes
6.8M edges
**labeled**



4K nodes
88K edges
unlabeled



Synthetic
networks

# Results: Performance Comparison

# Results: Social Networks Synthetization



**Honest region**

**Attack edges**

**Sybil region**

**How to synthesize attack edges?**

# Results: Robustness
## Label Noise Level

# Results: Robustness
## Label Noise Level

# Adversarial Attack



**Attack edges**

**Sybil region**

**Honest region**

**How to place attack edges in a targeted fashion?**

$$p_t \qquad \boldsymbol{p}$$

# Adversarial Attack

**Targeted attack edge**

$\text{Uniform}(N_k(\ \textbf{H}_{\textbf{train}}\ ))$

With prob. $p_k$

$\text{Uniform}(\ \textbf{S}\ )$

With prob. $p_t$

$$\boldsymbol{p} = [0.25, 0.25, 0.5] \in [0, 1]^K$$

direct hit     neighbor     2-hop neighbor

# Results: Adversarial Attacks

Increasing expected number of targeted attack edges

**Attack**

| | |
|---|---|
| (blue) | $p_t = 0.05, \; p = [0.25, 0.25, 0.5]$ |
| (orange) | $p_t = 0.10, \; p = [0.25, 0.25, 0.5]$ |
| (green) | $p_t = 0.15, \; p = [0.25, 0.25, 0.5]$ |
| (red) | $p_t = 0.20, \; p = [0.25, 0.25, 0.5]$ |

# Results: Adversarial Attacks

# Conclusion

- $\textsc{Sybil}\mathrm{GCN}$ and $\textsc{Sybil}\mathrm{GAT}$ algorithms outperform the baselines in almost all scenarios, including real-world X/Twitter dataset

- $\textsc{Sybil}\mathrm{RGCN}$ excels when attack complexity is high

- Synthesized social networks

**ETH** *zürich*

# Limitations & Future Work

- More advanced GNN architecture somewhat unexplored

- Robustness to lack of training data and label noise

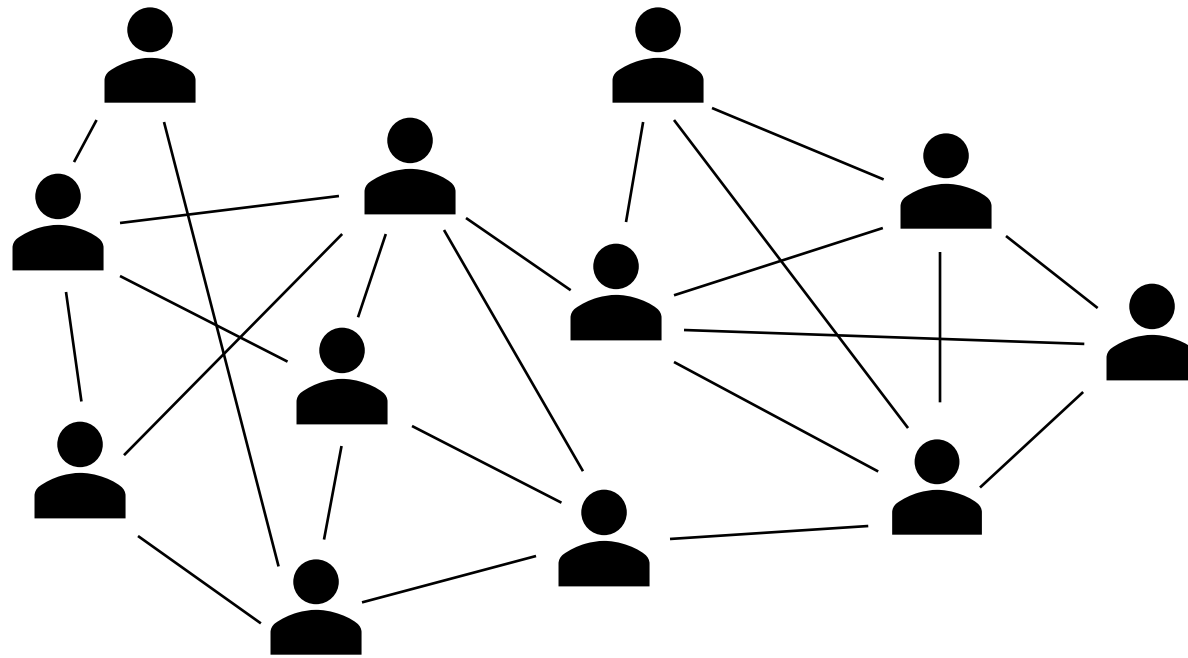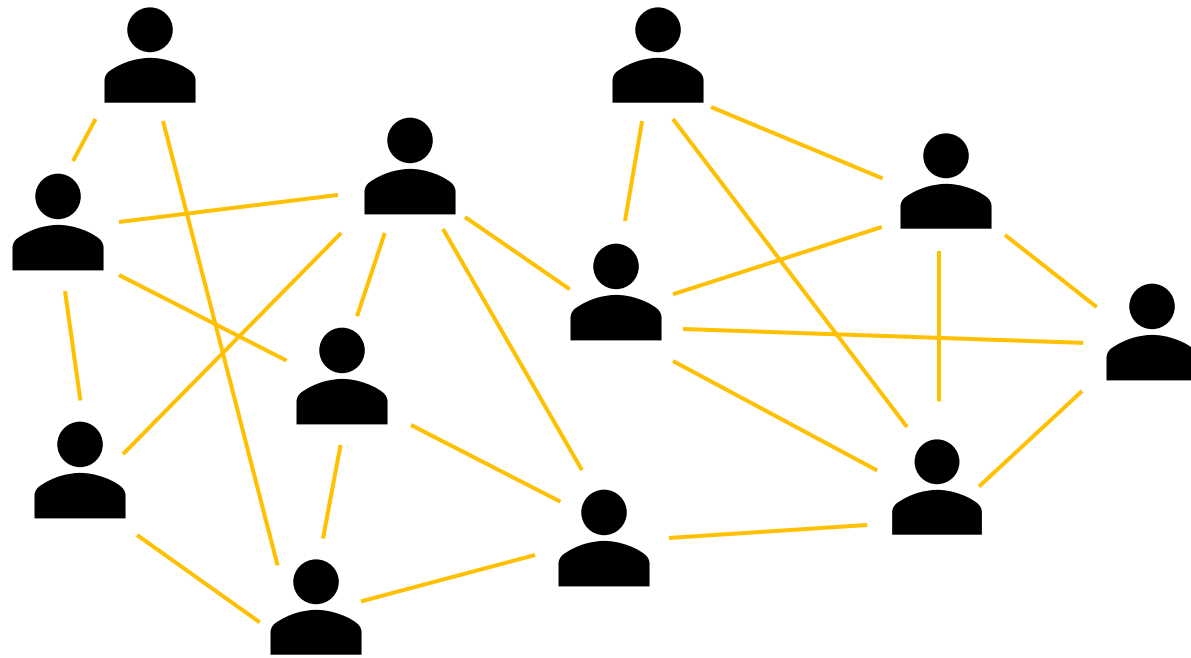- More data

- Theoretical analysis
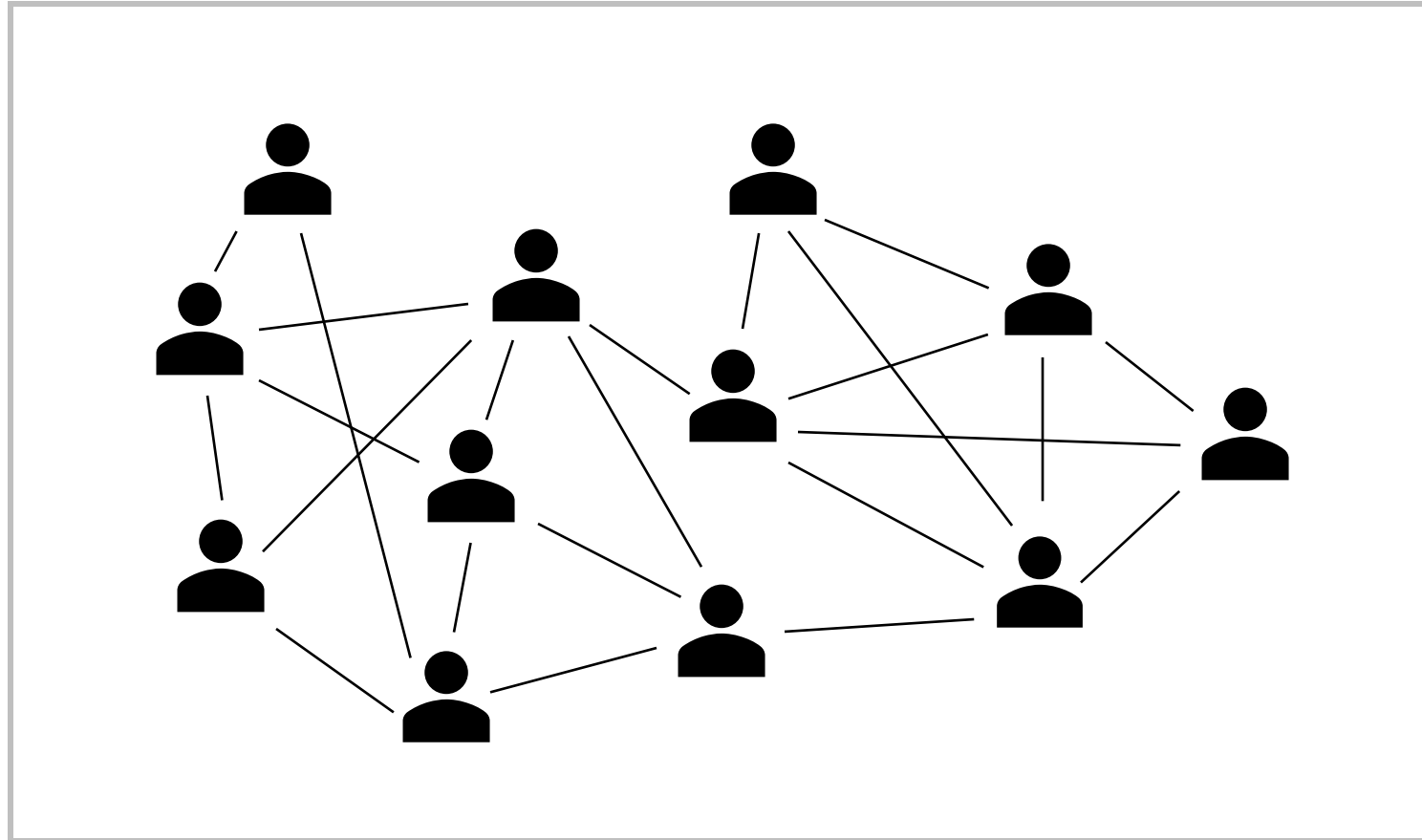
**ETH**_zürich_

**Thank you for listening!**

DITET    DINFK
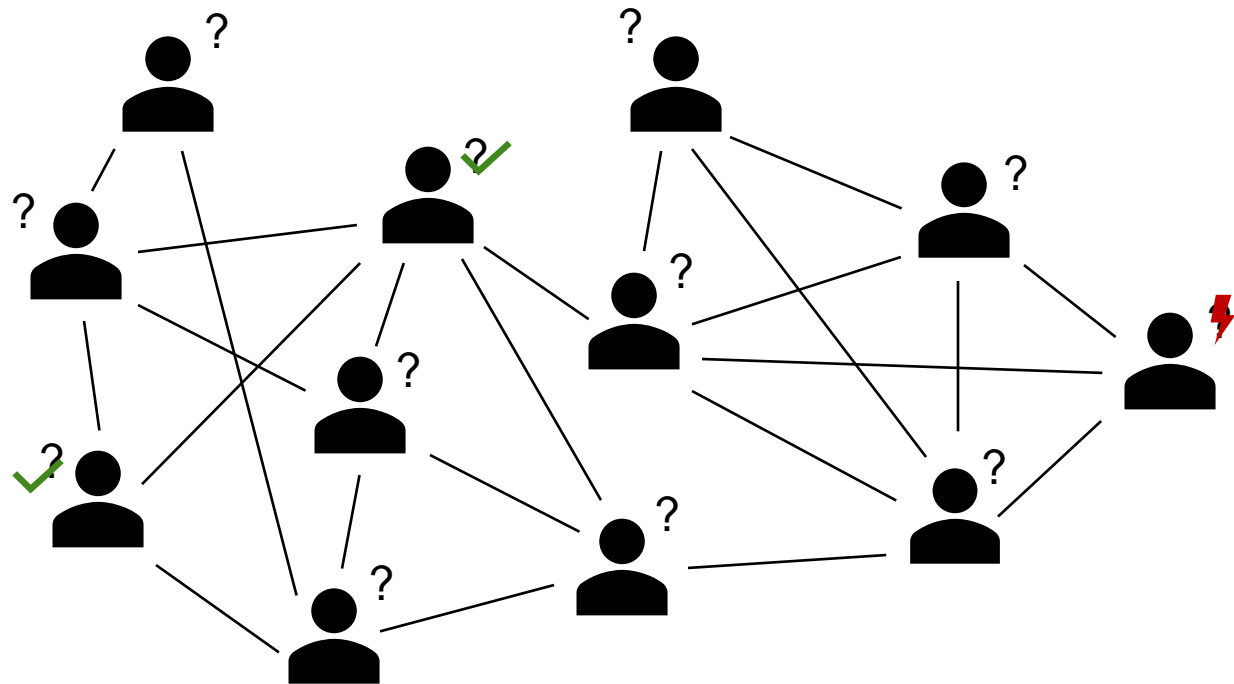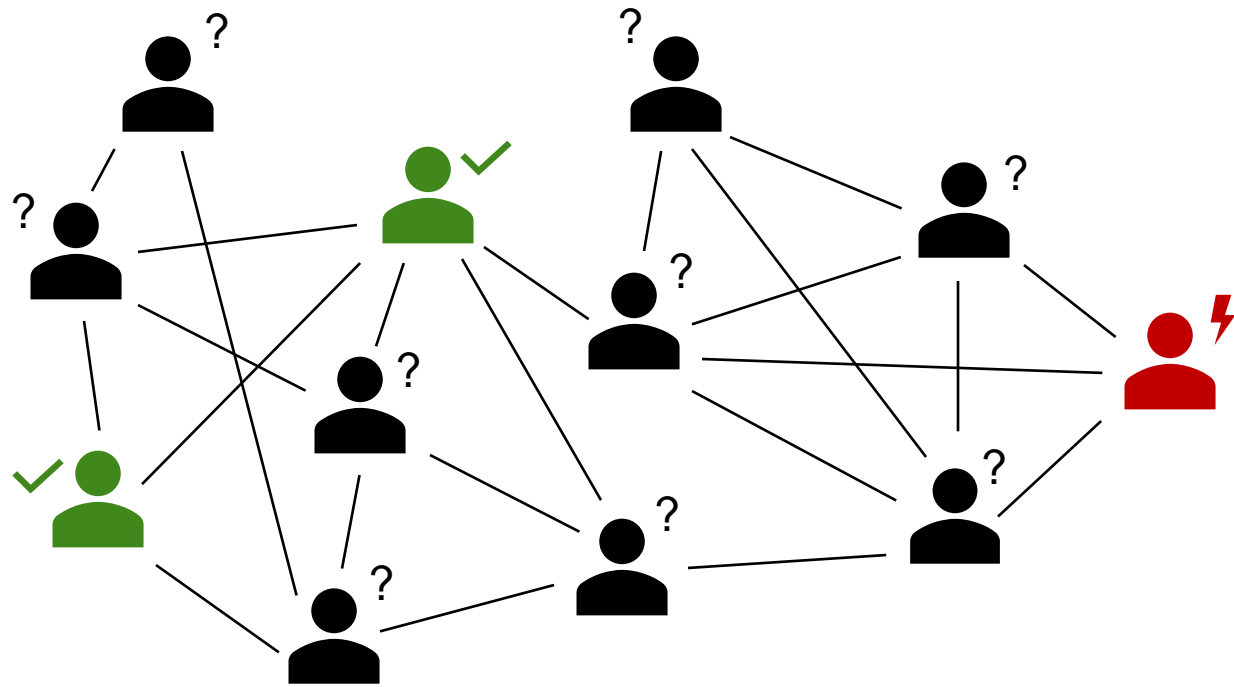
Stuart Heeb

stuart.heeb@inf.ethz.ch

# Backup Slides

trust

trust

**Users on a social media platform**

# Background & Related Work

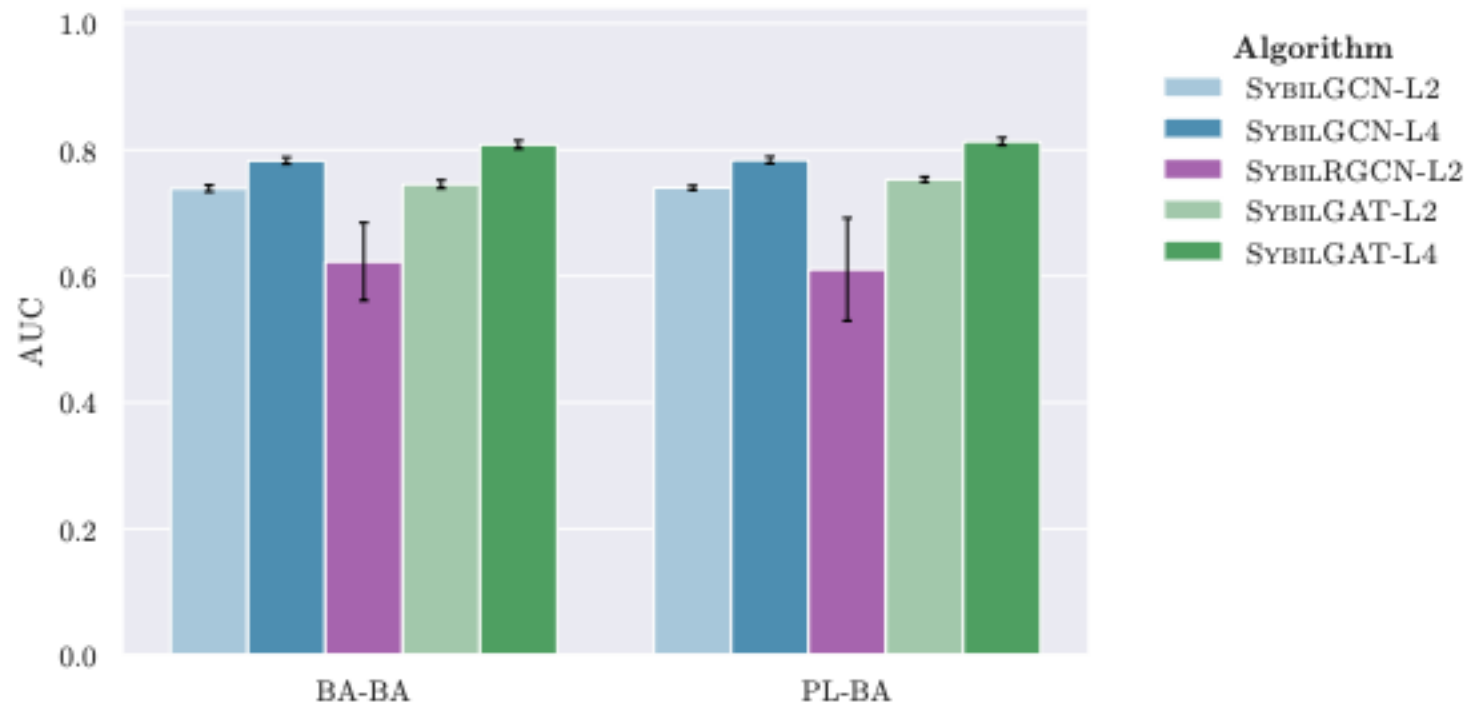|  | **SybilRank (2012)** | **SybilBelief (2014)** | **SybilSCAR (2019)** |
|---|---|---|---|
| **Approach** | Random Walks (RW) | Loopy Belief Propagation (LBP) | Local rule-based propagation |
| **Labels (honest / Sybil)** | Only honest | Both | Both |
| **Guaranteed convergence** | Yes | No | Yes |
| **Computational complexity** | O(n log n) | O(n) per iteration | O(n log n) |
| **Main advantage** | Computationally efficient | Can use both label types | Combines benefits of RW and LBP |

**ETH** *zürich*

# Modern Sybils

«The Eiffel Tower is a famous landmark in Paris, France. It stands tall at 330 meters (1,083 feet) with a square base. It is named after the engineer Gustav Eiffel.»

«The Eiffel Tower stands in Paris, France. It reaches a height of 330 meters (1,083 feet). The tower is named after Gustave Eiffel, the engineer whose company designed and built it.»
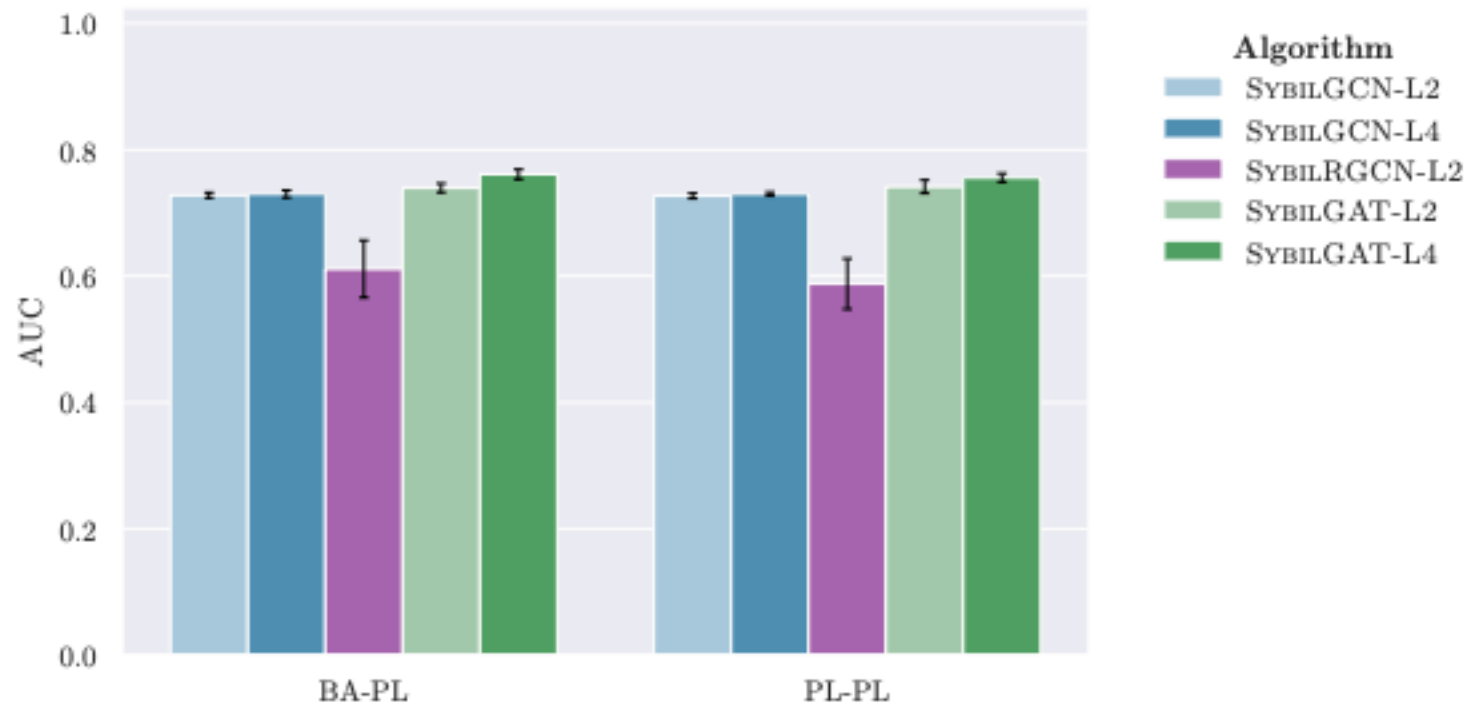
**Claude 3.5 Sonnet**
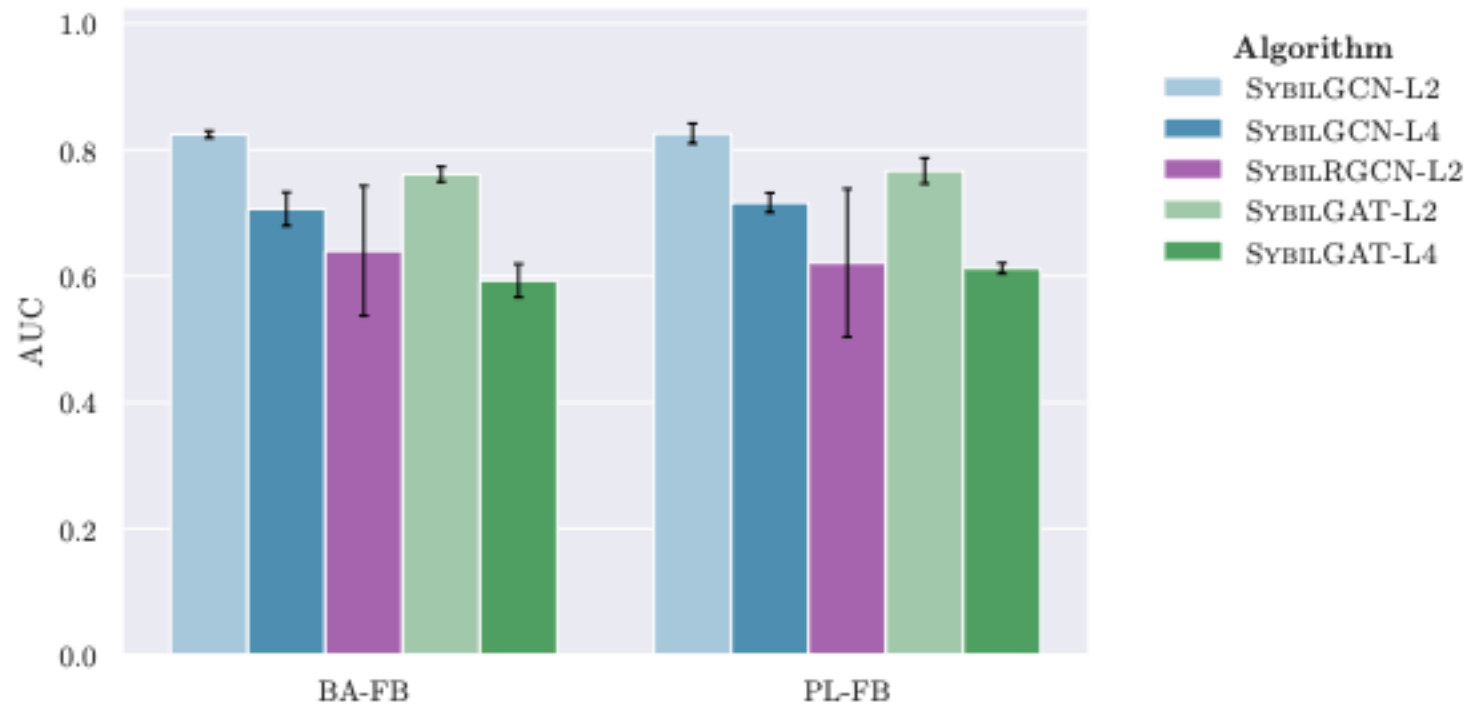*«write exactly 3 very short sentences about the eiffel tower, including its location, height and name origin»*

# Results: Pre-training on Small Network
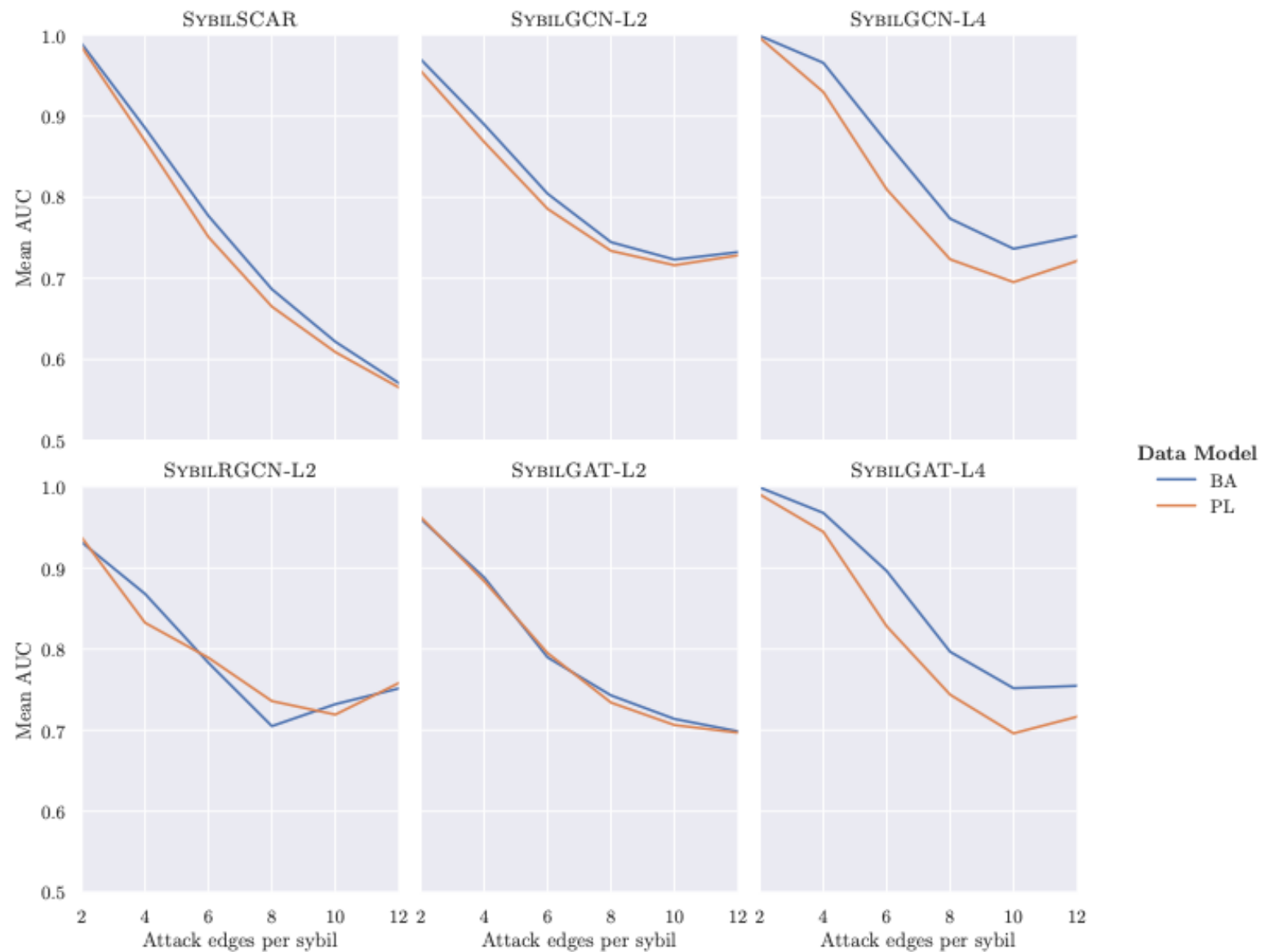
# Results: Pre-training on Small Network

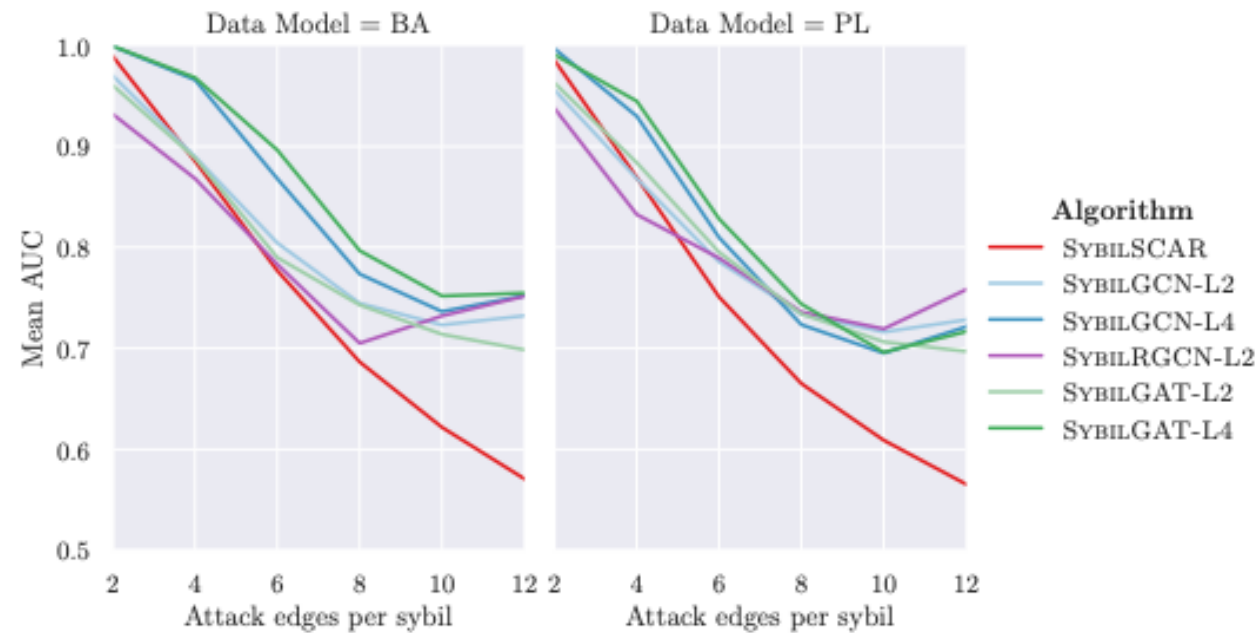# Results: Pre-training on Small Network
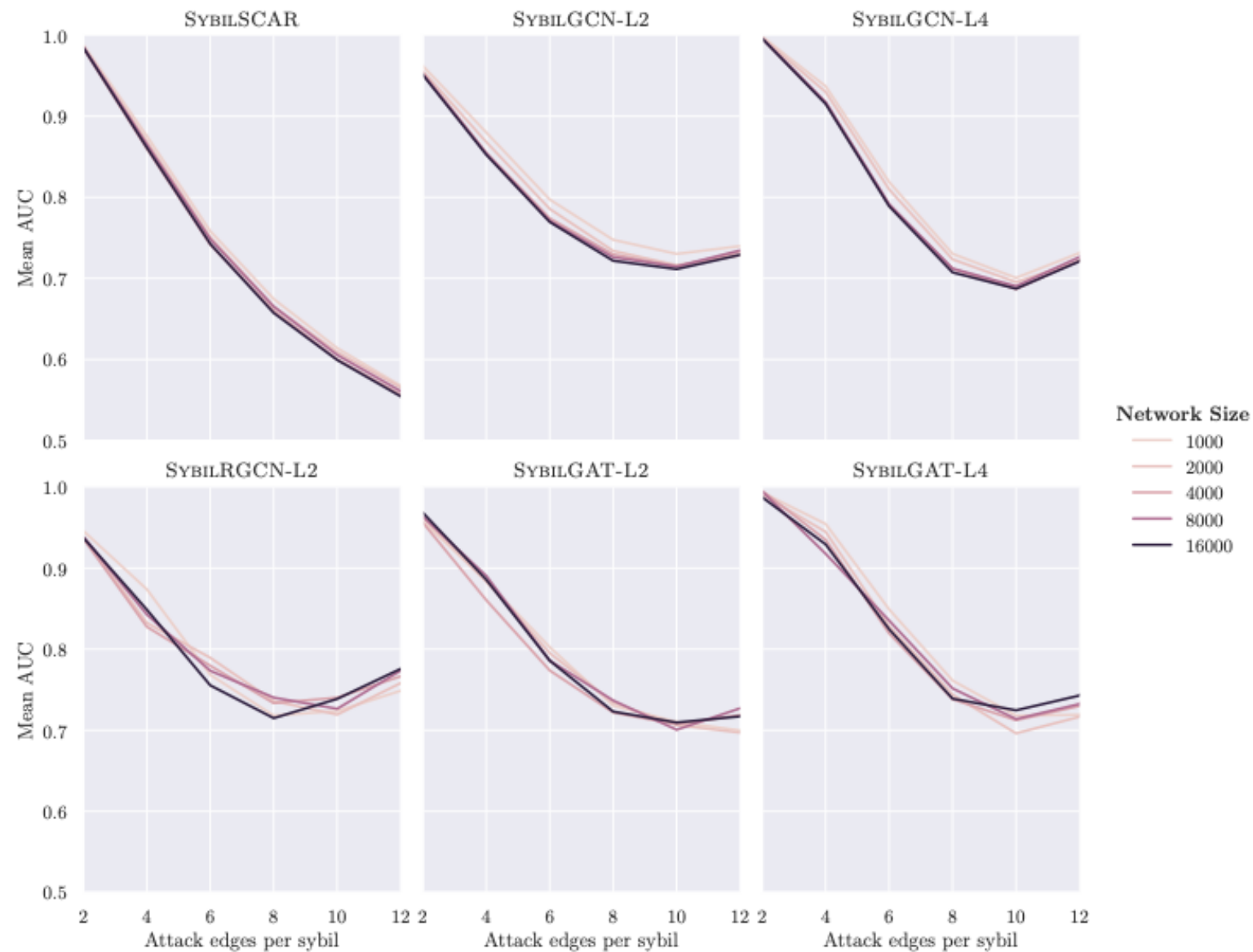
# Results: Robustness
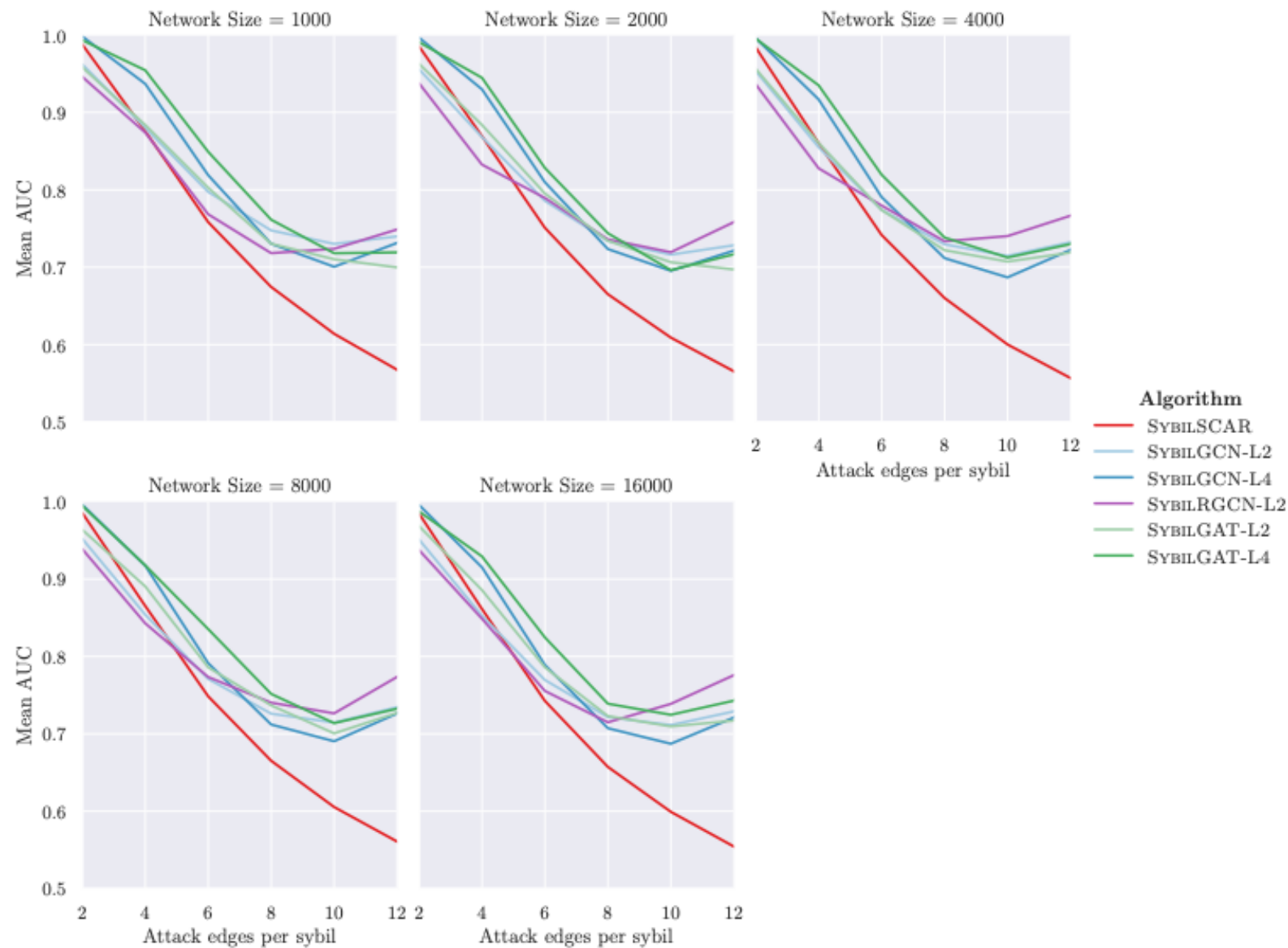## Data Model

# Results: Robustness
**Data Model**

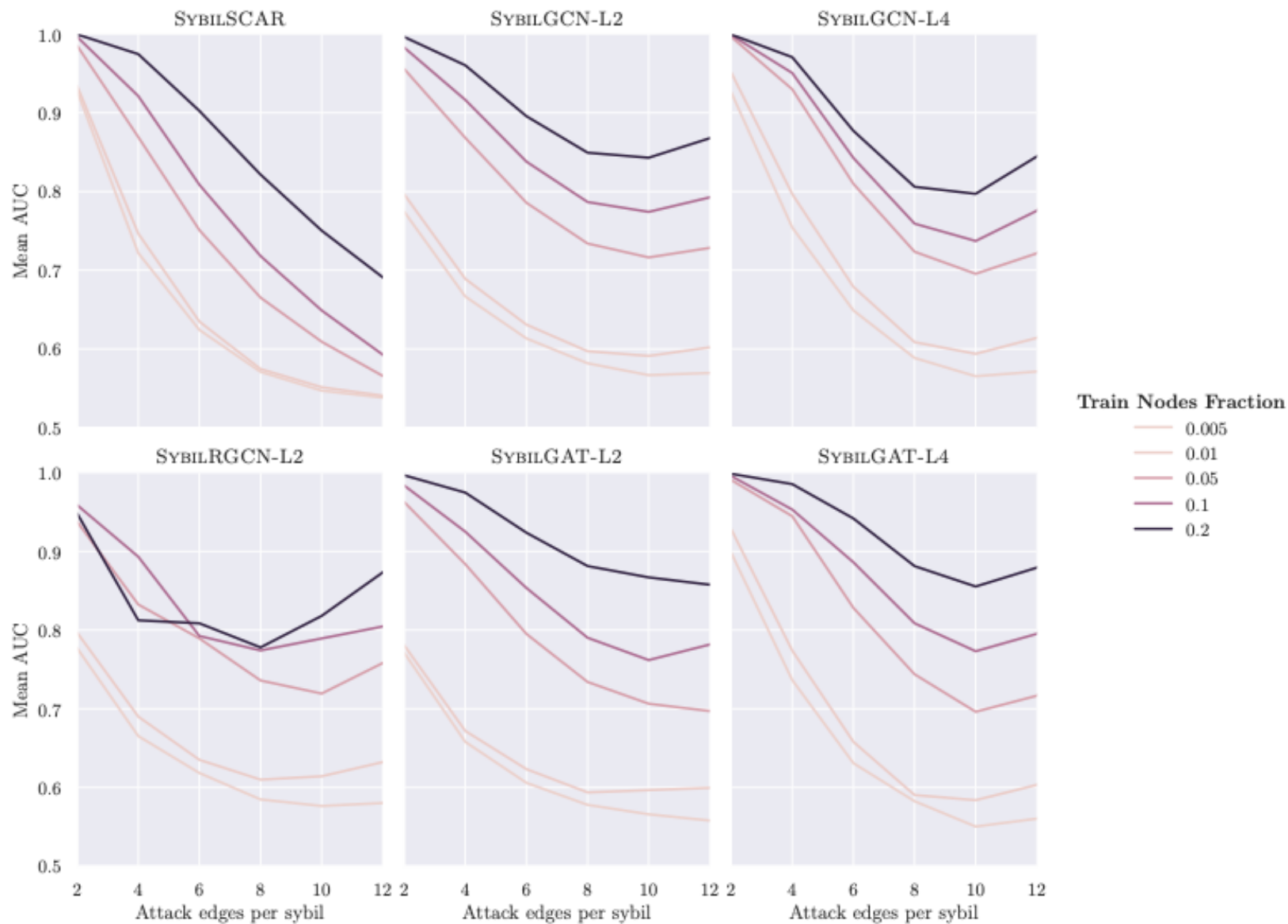# Results: Robustness
## Network Size
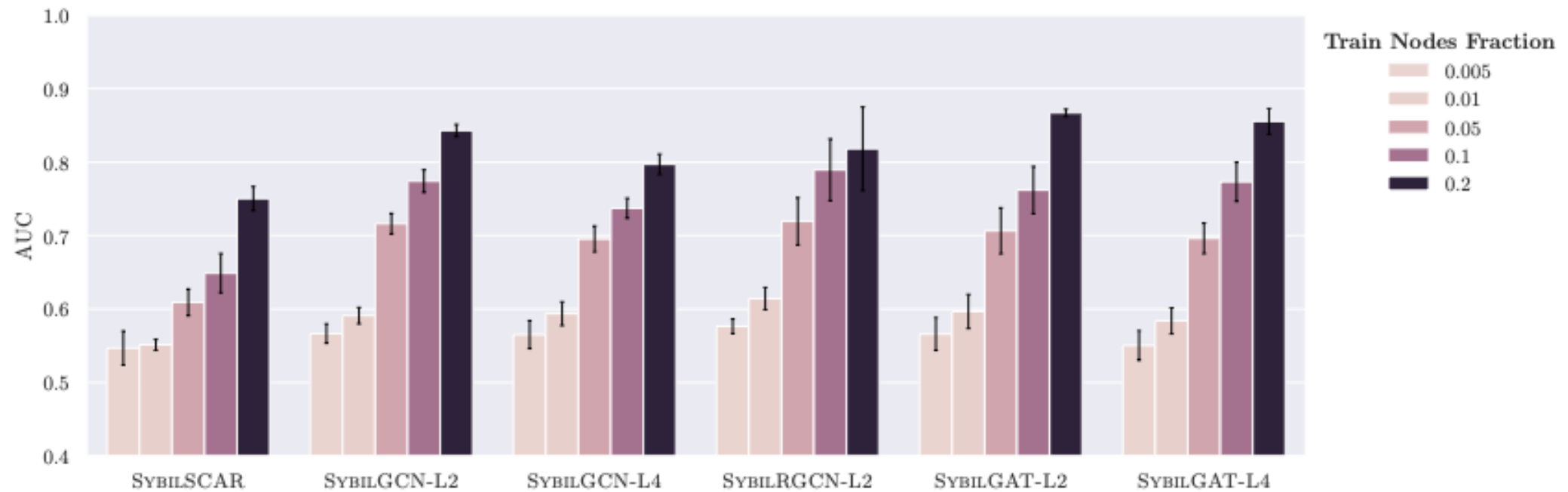
# Results: Robustness
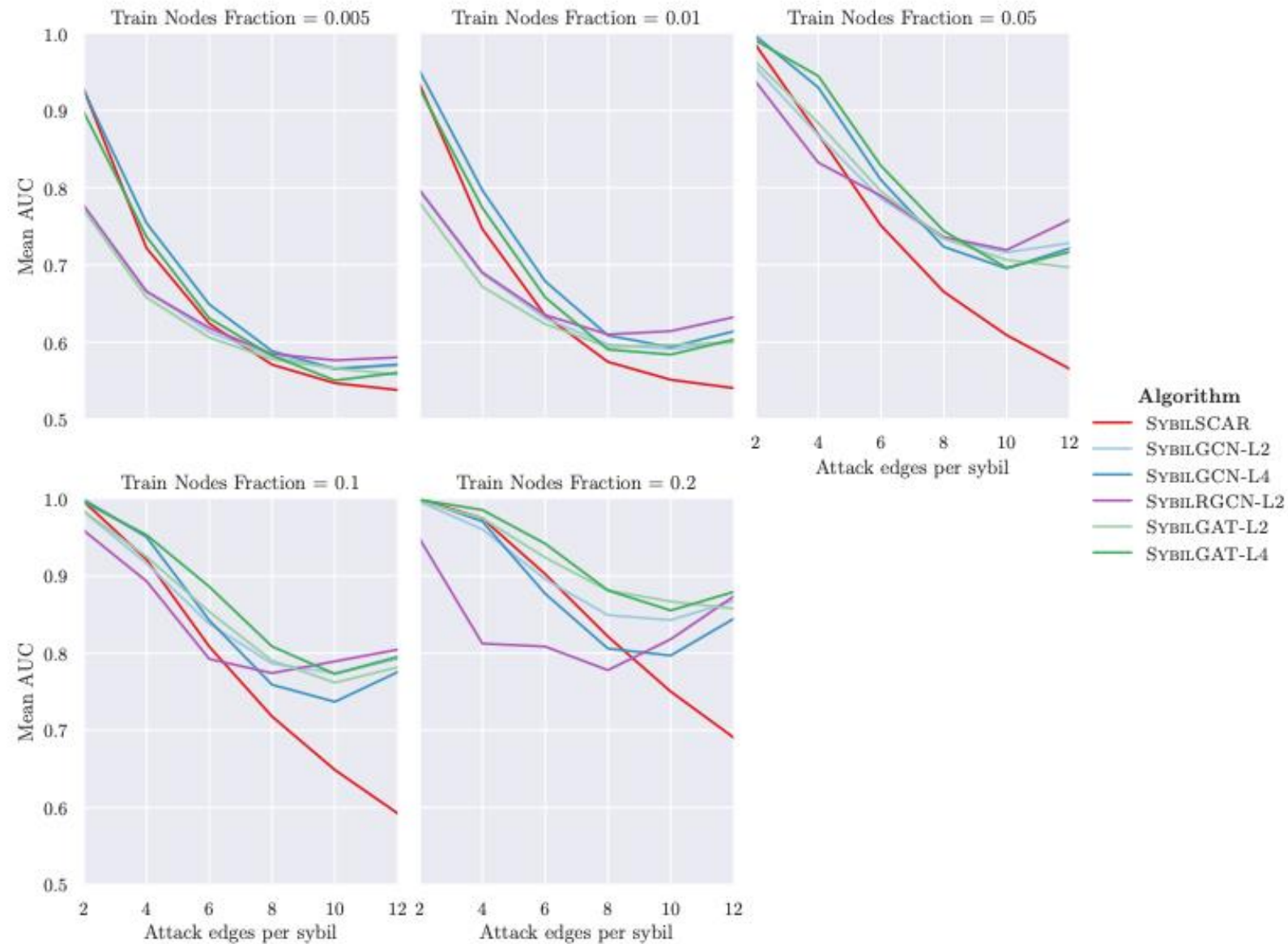## Network Size

# Results: Robustness
## Training Set Size
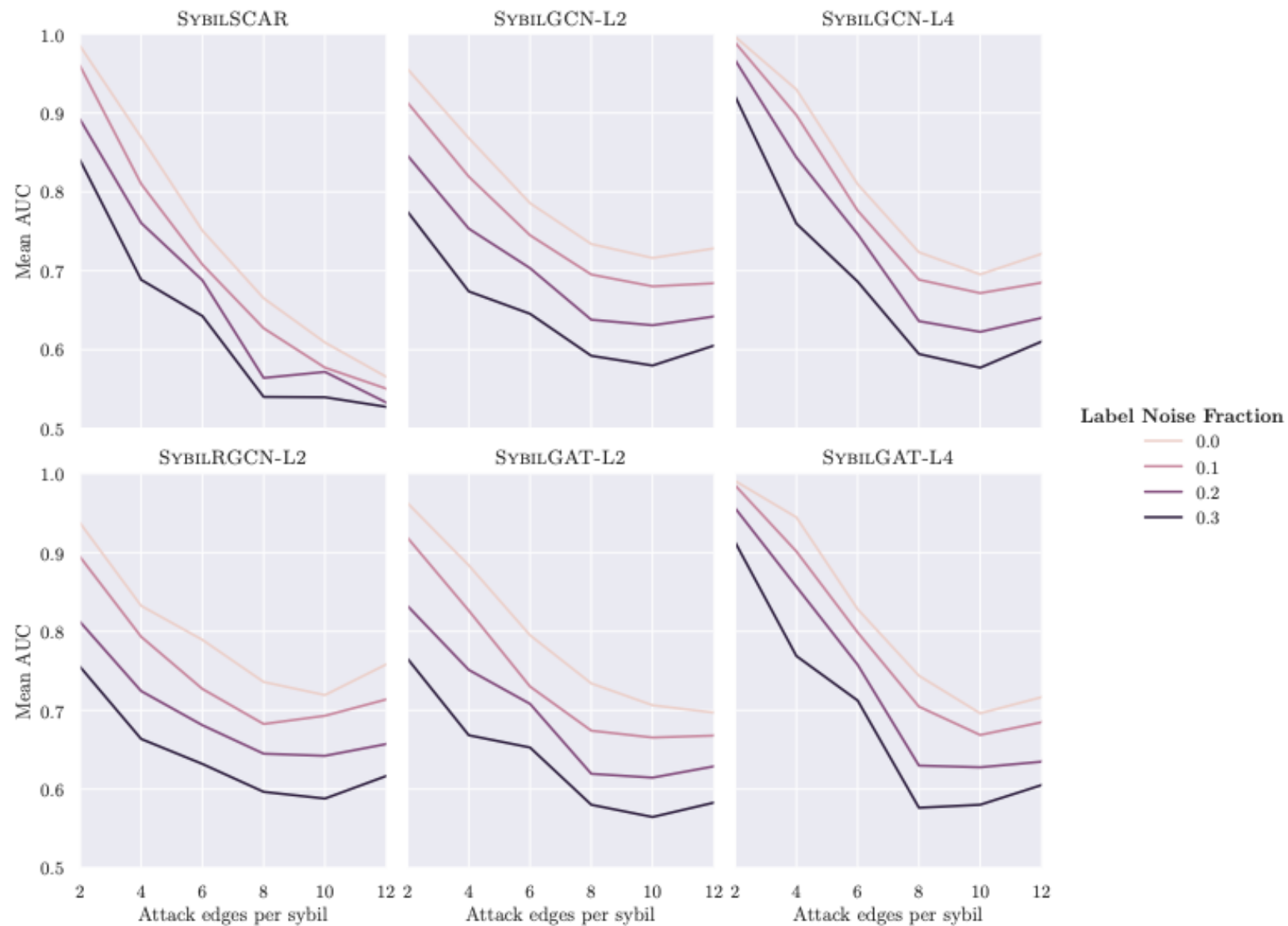
# Results: Robustness
## Training Set Size
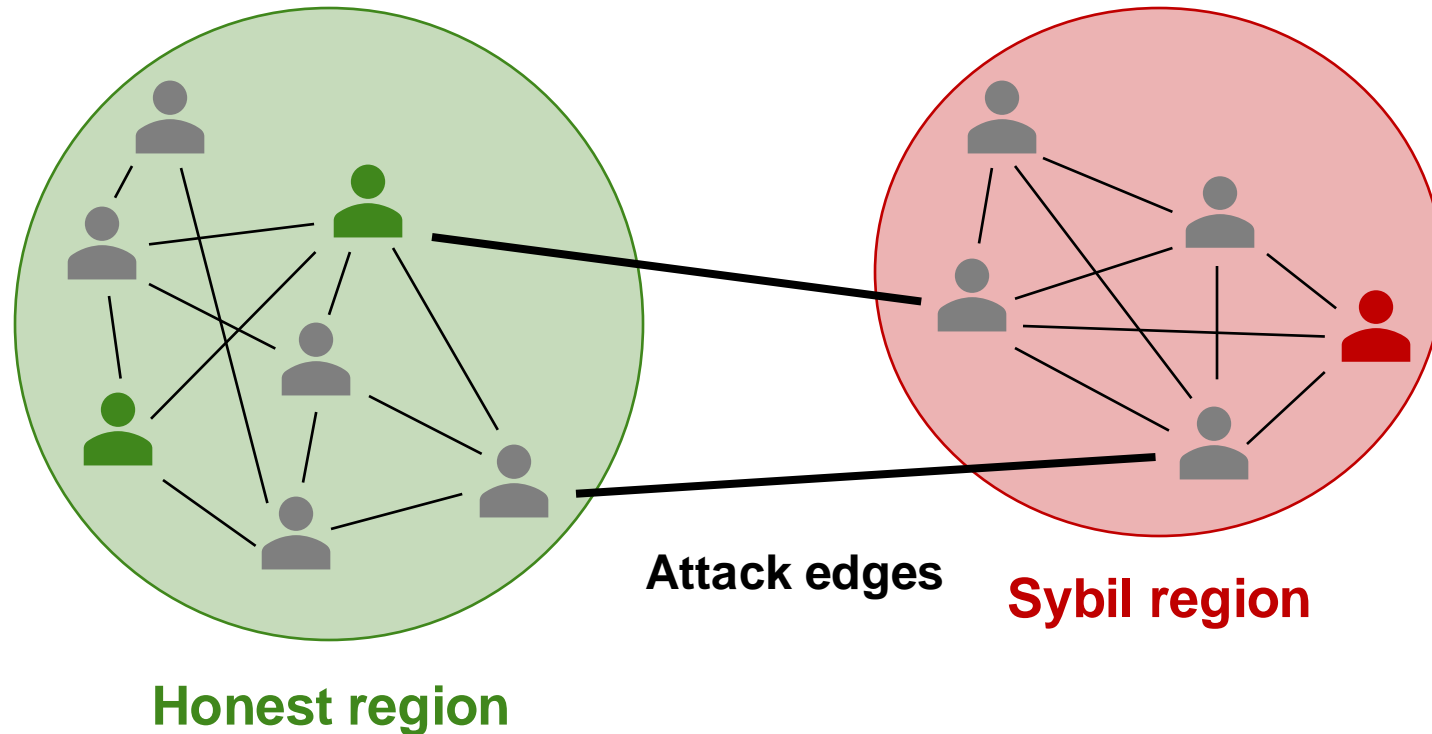
# Results: Robustness
## Training Set Size

# Results: Robustness
## Label Noise Level

# Adversarial Attack
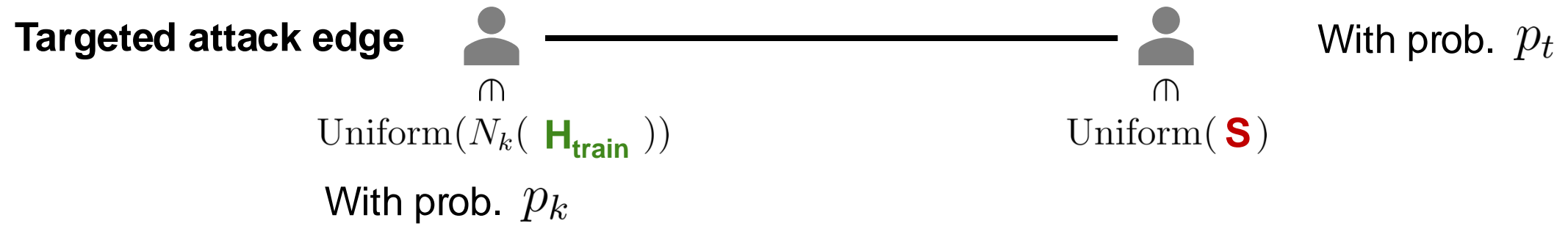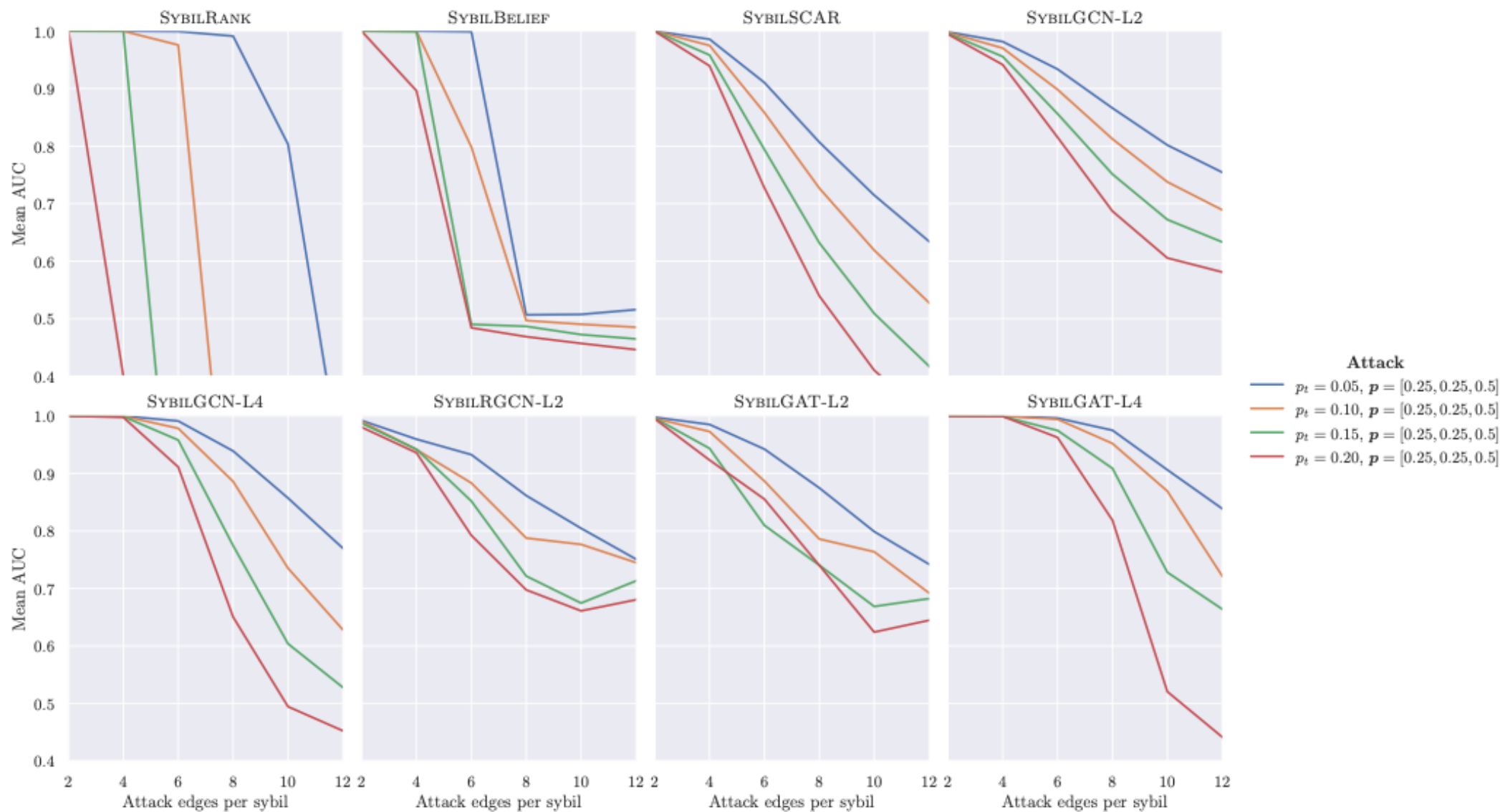
**Random attack edge**

$$\text{Uniform}(\mathbf{H}) \quad\quad\quad\quad \text{Uniform}(\mathbf{S})$$

With prob. $1 - p_t$



**Attack edges**

**Sybil region**

**Honest region**

# Adversarial Attack

**Targeted attack edge**



With prob. $p_t$

$\text{Uniform}(N_k(\ \mathbf{H_{train}}\ ))$

$\text{Uniform}(\ \mathbf{S}\ )$

With prob. $p_k$

**Attack edges**

**Sybil region**

**Honest region**

# Results: Adversarial Attacks

# Results: Pre-training Before Attack