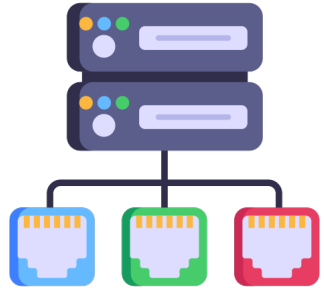


CiDer: A Black-box Approach to Classify Node with Certified Robustness Guarantees

Xiaoyu Liang, Haohua Du, Wen Ma, Ye Tian, and Xiaoya Xu



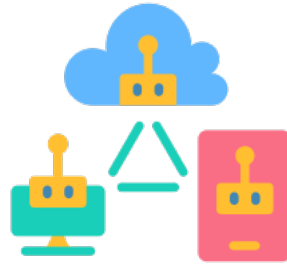
Graph-structured data



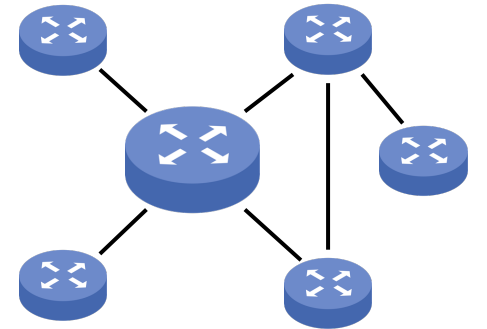
Router network



IoT



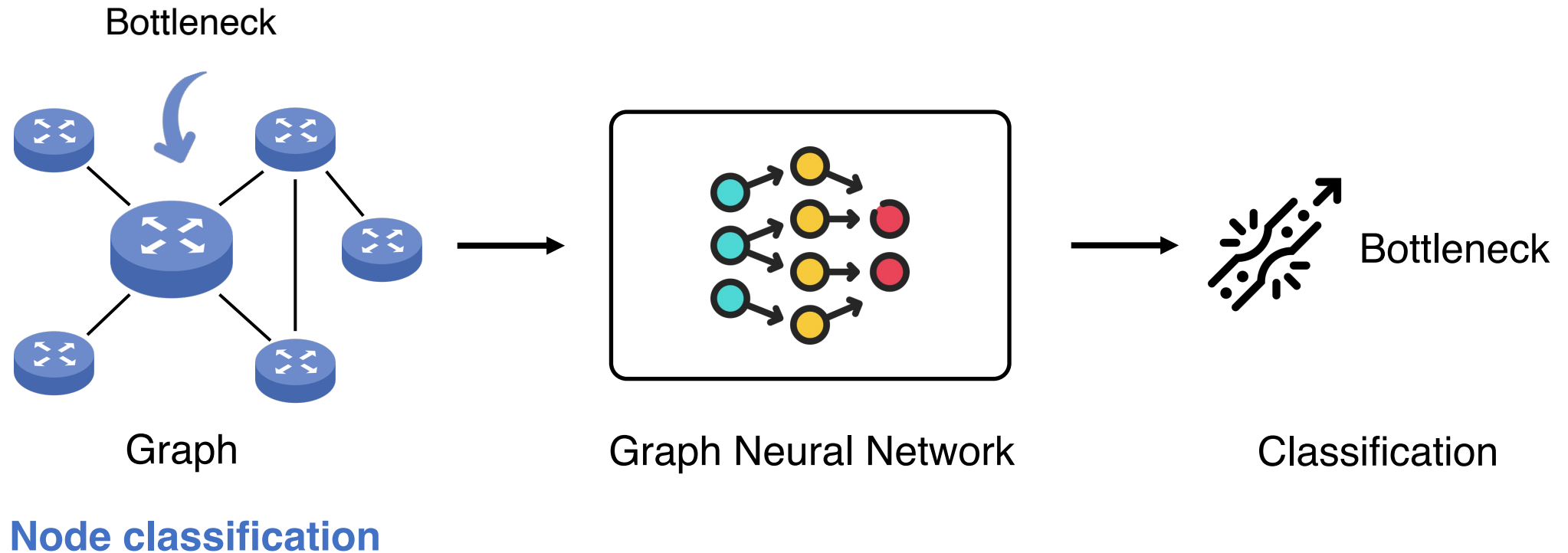
Edge computing



Graph

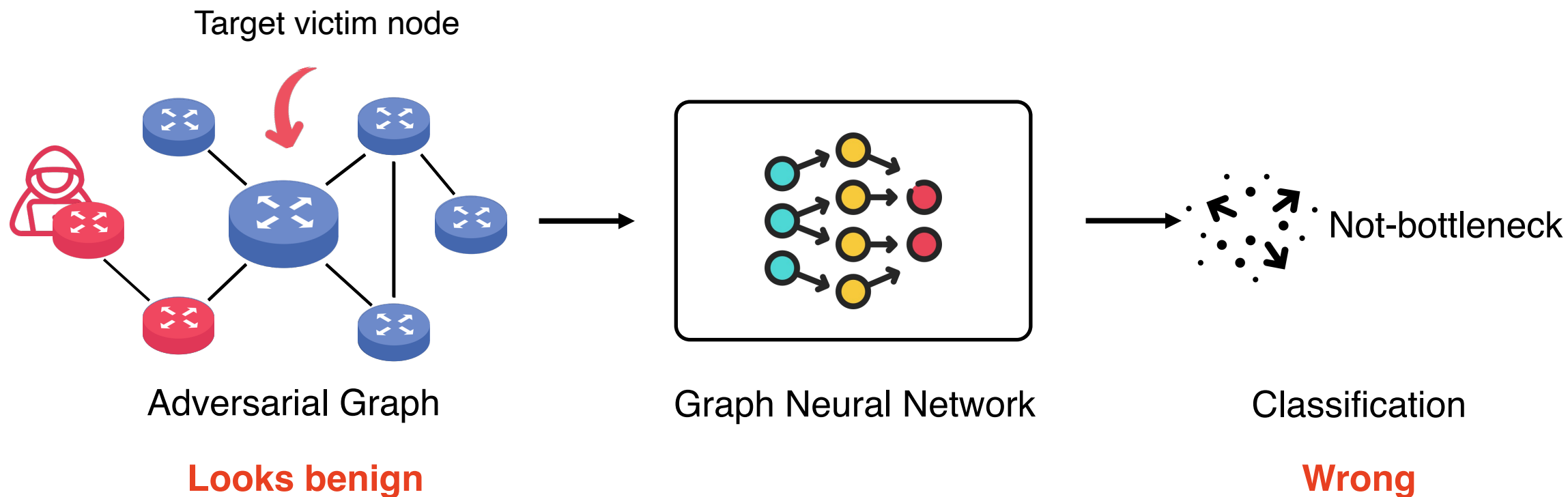
Real-world relationships are often abstracted into **graphs**.

Graph Node Classification



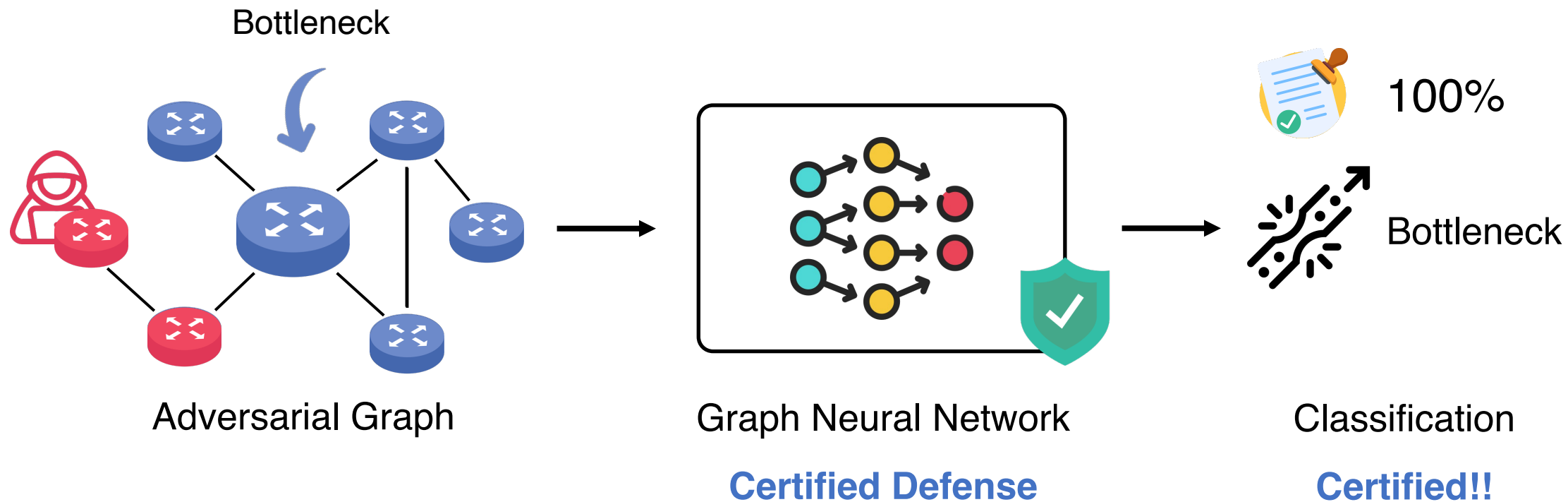
Graph neural network (**GNN**) is used for **node classification**.

Graph Evasion Attack



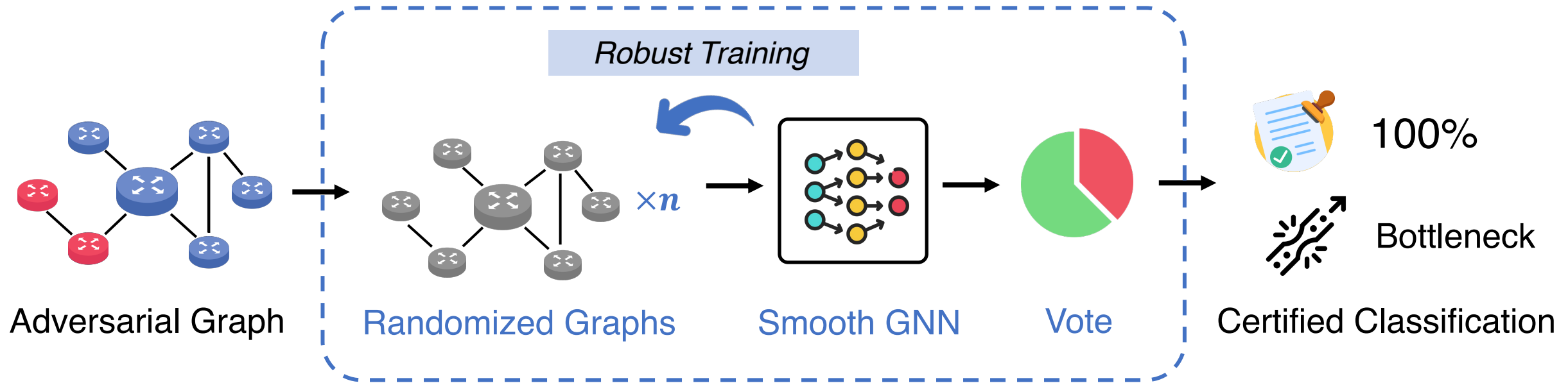
Attackers on the graph network could lead to misclassification by GNN.

Robustness Certificate



Robustness certificate gives provable robustness guarantee of the classification.

Randomized Smoothing is successful, but

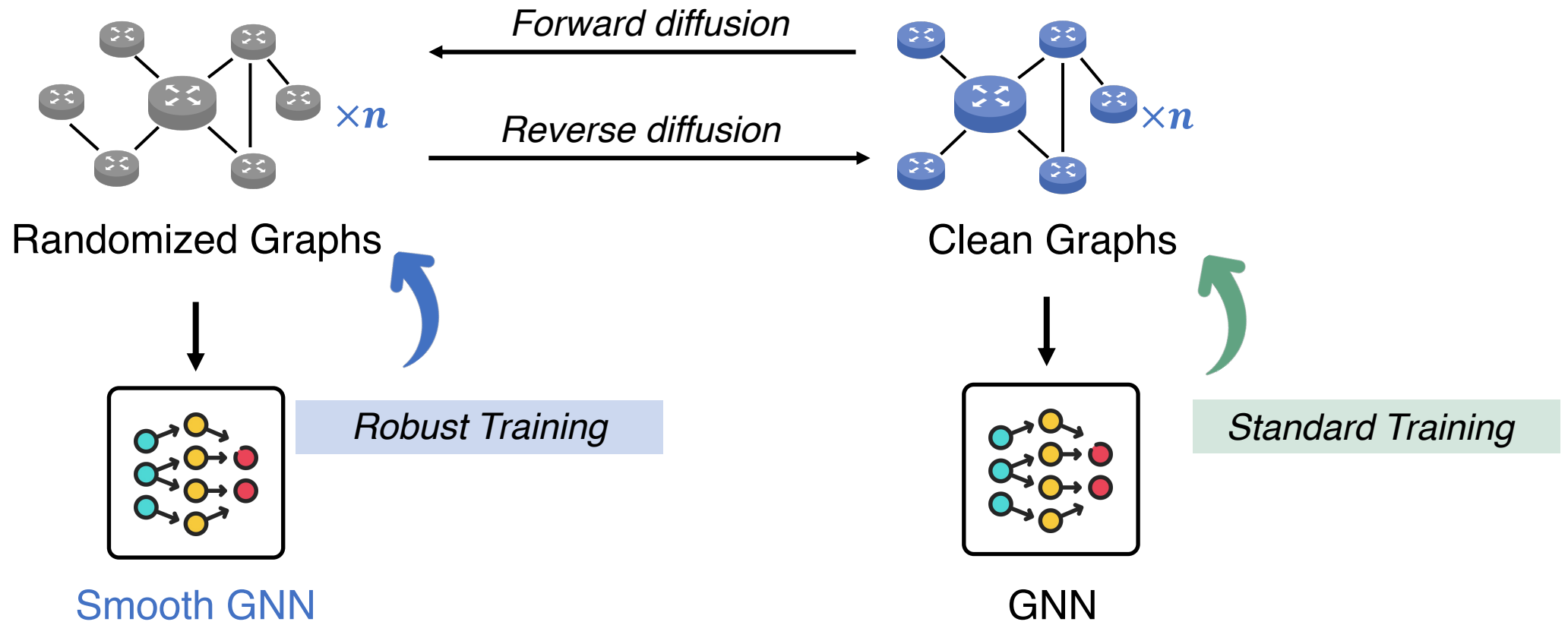


Randomized Smoothing

(Cohen *et al.* 2019, Bojchevski *et al.* 2020)

1. Need **fine-tuning** on the noisy data, and
2. **Accuracy drop** on clean samples.

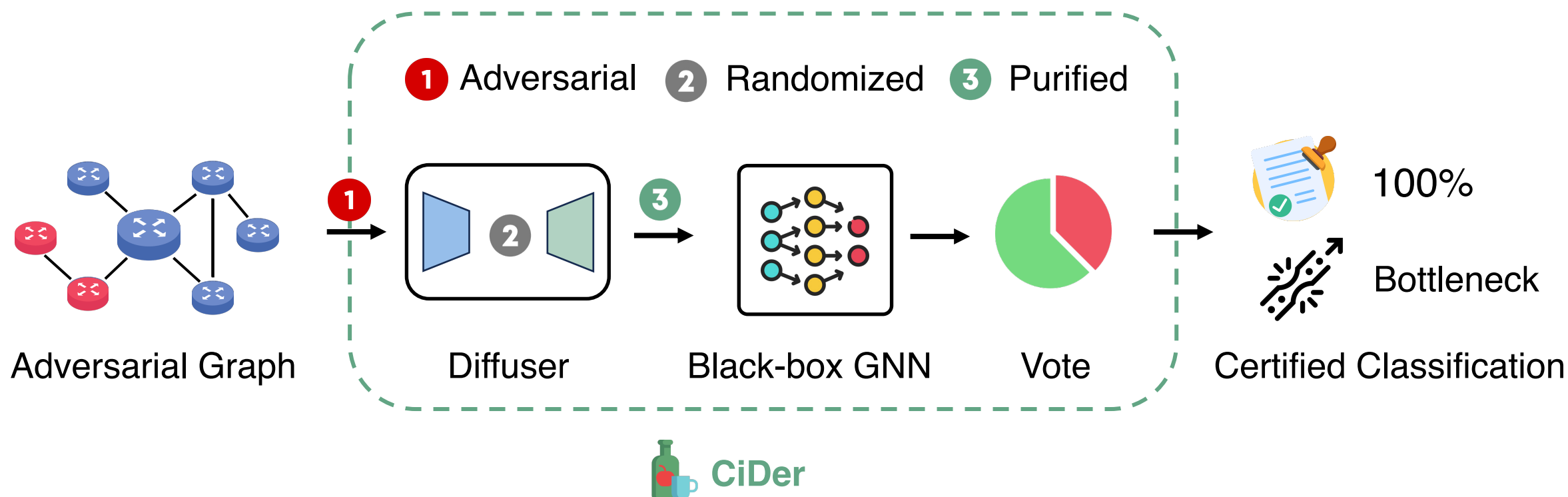
How to certify node classification with black-box GNNs?



Use **diffusion model** as a denoiser to purify the graph while maintain randomness.

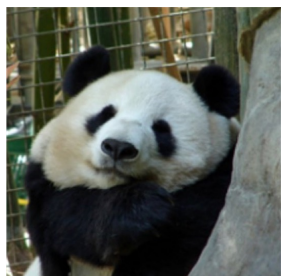
Our Idea: Graph diffusion denoised smoothing

Graph Diffusion Denoiser: Leverage the *denoising* capability of discrete diffusion models (Austin *et al.* 2021) to select the *core features* of different node types and generate the certificate with RS.



Diffusion denoiser works on images, but

1 Adversarial

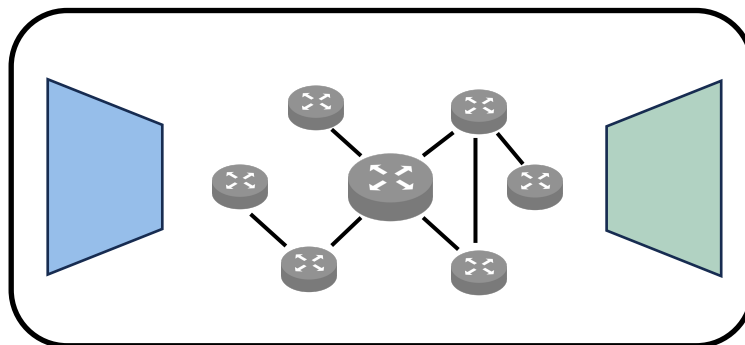
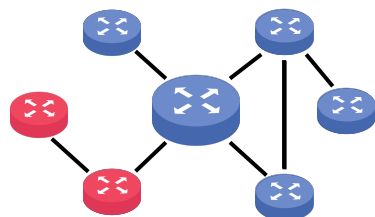
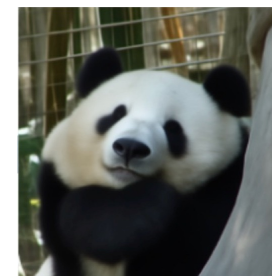


2 Randomized

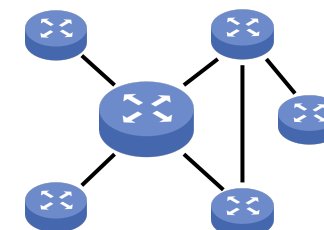


DDIM

3 Purified



Graph Joint Diffuser



Euclidean

Pixel-wise

Fixed-size

Non-Euclidean

Node-edge

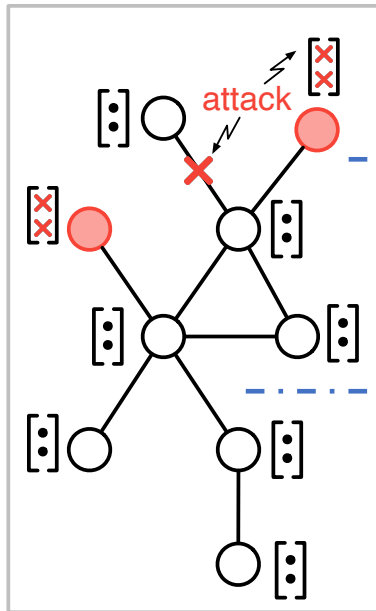
Large graph

Subgraph Extractor



Graph modification attack

○ benign ● attacked

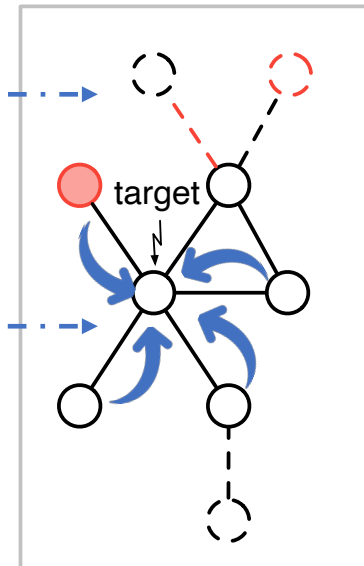


adversarial graph

$$\tilde{G} = (\tilde{X}, \tilde{A})$$

1 $X' = \{u \mid d(u, v) \leq L\}$

2 $A' = \{(u, w) \mid u, w \in X'\}$



subgraph $G' = (X', A')$ of target v

Large attributed graph

Efficiency gain

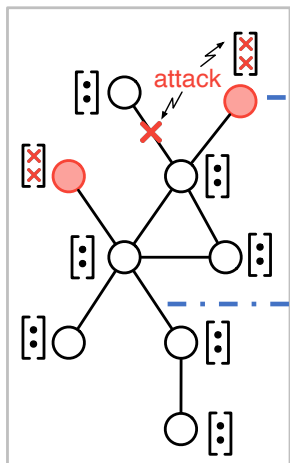
Subgraph for MPNN

Graph Joint Diffuser



Graph modification attack

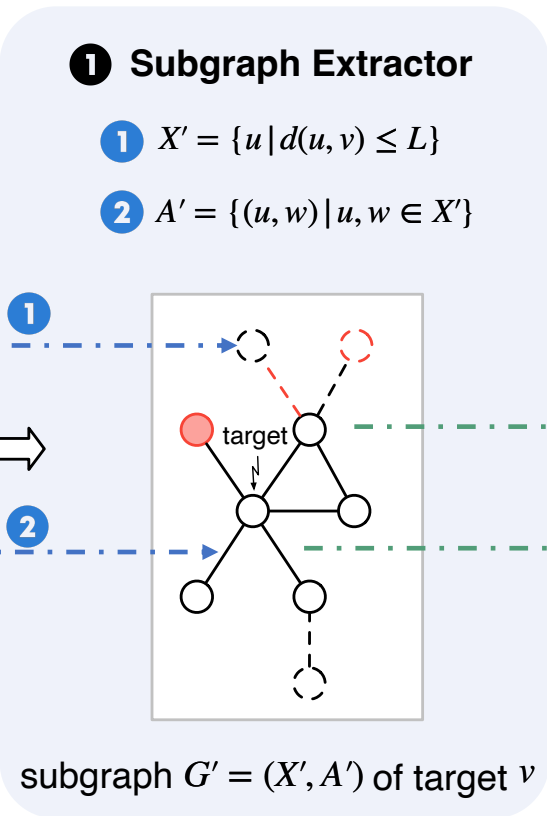
○ benign ● attacked



adversarial graph
 $\tilde{G} = (\tilde{X}, \tilde{A})$

1 Subgraph Extractor

- 1 $X' = \{u \mid d(u, v) \leq L\}$
- 2 $A' = \{(u, w) \mid u, w \in X'\}$



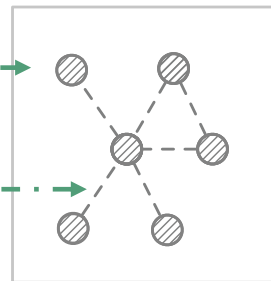
subgraph $G' = (X', A')$ of target v

diffusion
randomization

Non-Euclidean graph structure

$$2 \quad X_t \sim X' \oplus \epsilon_{t_x}^X$$

$$1 \quad A_t \sim A' \oplus \epsilon_{t_a}^A$$



$$\epsilon_{t_z}^Z \sim \text{Ber}(p = (\overline{\beta_{t_z} m_1^Z})^{1-Z} (\overline{\beta_{t_z} m_0^Z})^Z) \text{ with } Z \in \{X, A\}$$

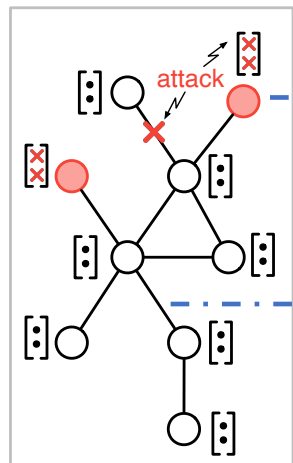
Data-dependent

Graph Joint Diffuser



Graph modification attack

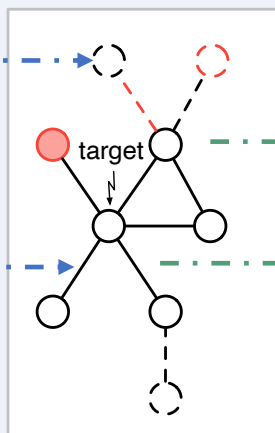
○ benign ● attacked



adversarial graph
 $\tilde{G} = (\tilde{X}, \tilde{A})$

1 Subgraph Extractor

- 1 $X' = \{u \mid d(u, v) \leq L\}$
- 2 $A' = \{(u, w) \mid u, w \in X'\}$



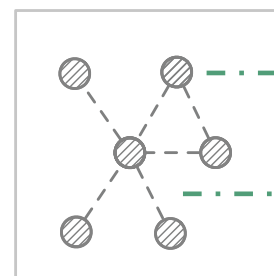
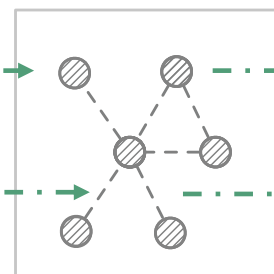
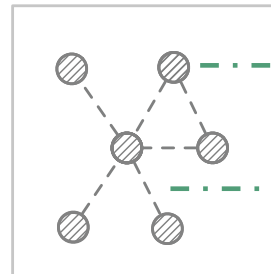
subgraph $G' = (X', A')$ of target v

*diffusion
randomization*

2 $X_t \sim X' \oplus \epsilon_{t_x}^X$

1 $A_t \sim A' \oplus \epsilon_{t_a}^A$

n samples



diffused graphs

$G_t = (X_t, A_t)$

*reverse
purification*

Node-edge relation

MLP

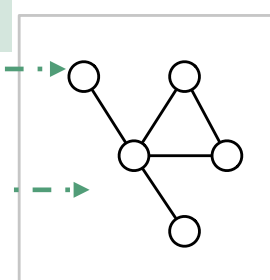
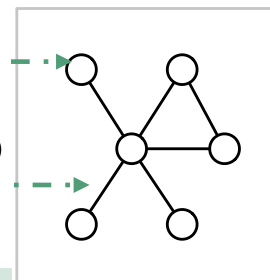
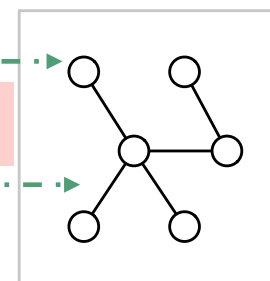
3 $\bar{X} \sim X_t \oplus \epsilon_{\theta}(X_t, t_x)$

4 $\bar{A} \sim A_t \oplus \epsilon_{\theta}(A_t, t_a \mid \bar{X})$

MPNN

Step-wise denoise

n samples



denoised graphs

$\bar{G} = (\bar{X}, \bar{A})$

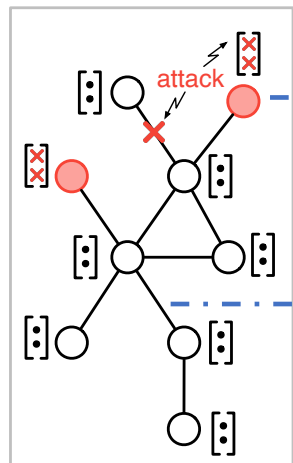
Certificate Generator

diffusion timestep = randomization parameter



Graph modification attack

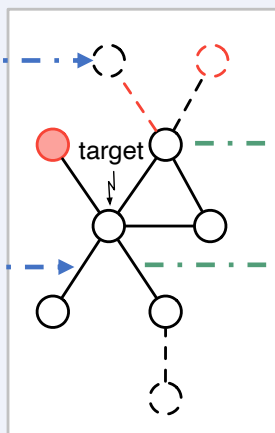
○ benign ○ attacked



adversarial graph
 $\tilde{G} = (\tilde{X}, \tilde{A})$

1 Subgraph Extractor

- 1 $X' = \{u \mid d(u, v) \leq L\}$
- 2 $A' = \{(u, w) \mid u, w \in X'\}$

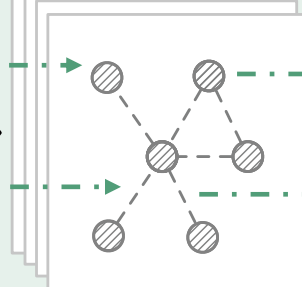


subgraph $G' = (X', A')$ of target v

2 Denoising Graph Joint Diffuser

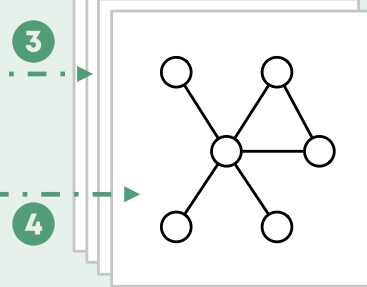
- 1 $A_t \sim A' \oplus \epsilon_{t_a}^A$
- 2 $X_t \sim X' \oplus \epsilon_{t_x}^X$
- 3 $\bar{X} \sim X_t \oplus \epsilon_\theta(X_t, t_x)$
- 4 $\bar{A} \sim A_t \oplus \epsilon_\theta(A_t, t_a \mid \bar{X})$

n samples



diffused graphs
 $G_t = (X_t, A_t)$

n samples



denoised graphs
 $\bar{G} = (\bar{X}, \bar{A})$



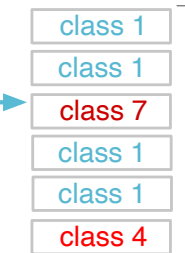
robustness certificate

$$f_\theta(v \mid \tilde{G}) = y^*, \quad \|\tilde{G} - G\| \in \Delta$$



GNN $f_\theta(v \mid G)$

n predictions

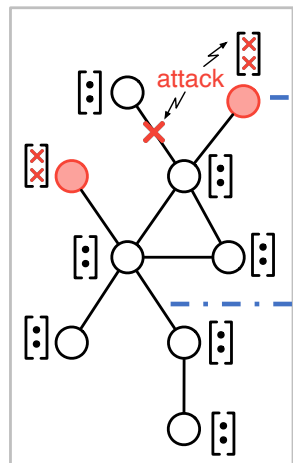


Certified robustness via Diffusion model on graph



Graph modification attack

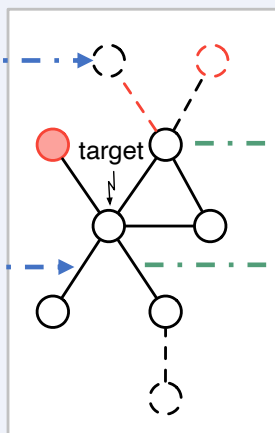
○ benign ● attacked



adversarial graph
 $\tilde{G} = (\tilde{X}, \tilde{A})$

1 Subgraph Extractor

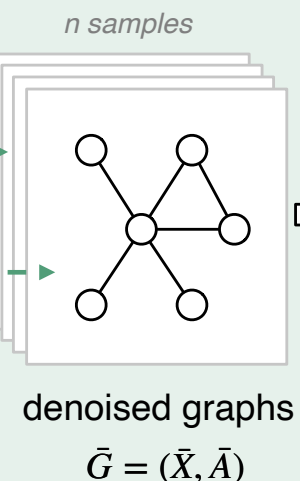
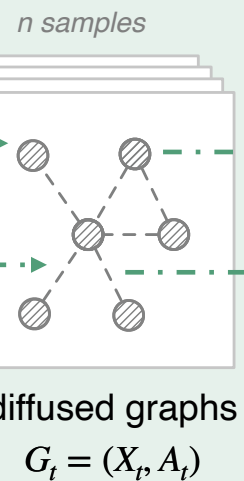
- 1 $X' = \{u \mid d(u, v) \leq L\}$
- 2 $A' = \{(u, w) \mid u, w \in X'\}$



subgraph $G' = (X', A')$ of target v

2 Denoising Graph Joint Diffuser

- 1 $A_t \sim A' \oplus \epsilon_{t_a}^A$
- 2 $X_t \sim X' \oplus \epsilon_{t_x}^X$
- 3 $\bar{X} \sim X_t \oplus \epsilon_{\theta}(X_t, t_x)$
- 4 $\bar{A} \sim A_t \oplus \epsilon_{\theta}(A_t, t_a \mid \bar{X})$



3 Certificate Generator



robustness certificate

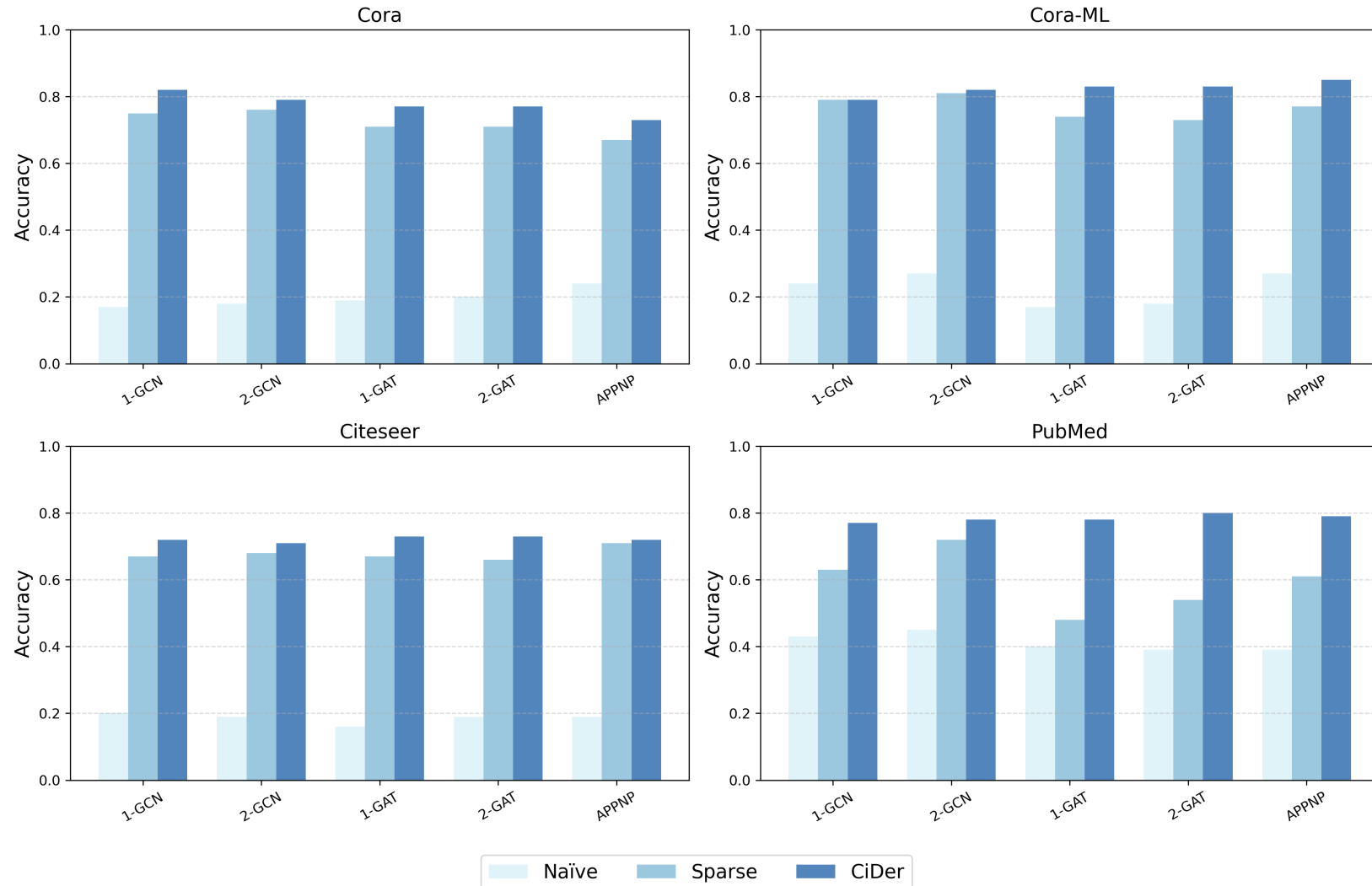
$$f_{\theta}(v \mid \tilde{G}) = y^*, \quad \|\tilde{G} - G\| \in \Delta$$



n predictions



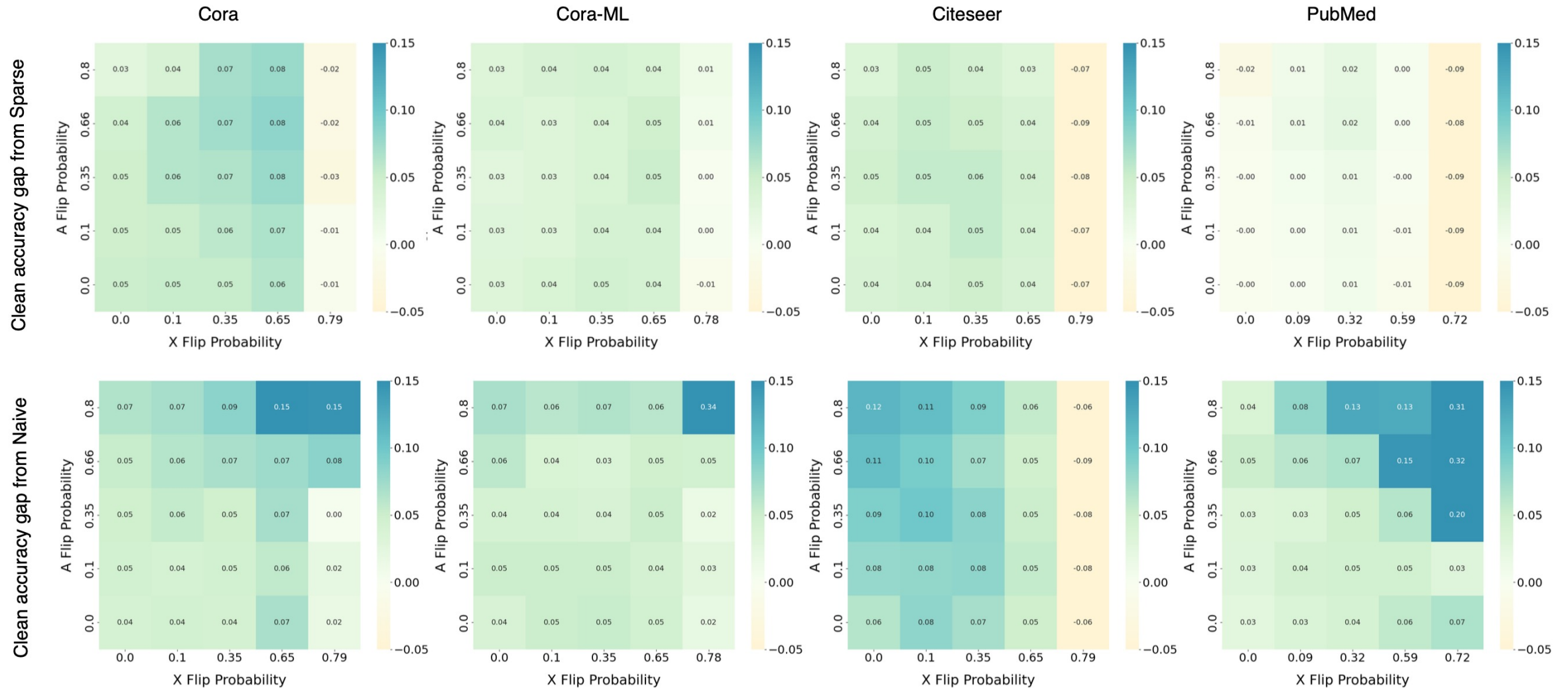
Evaluation – Clean Accuracy



Greater gain on large dataset

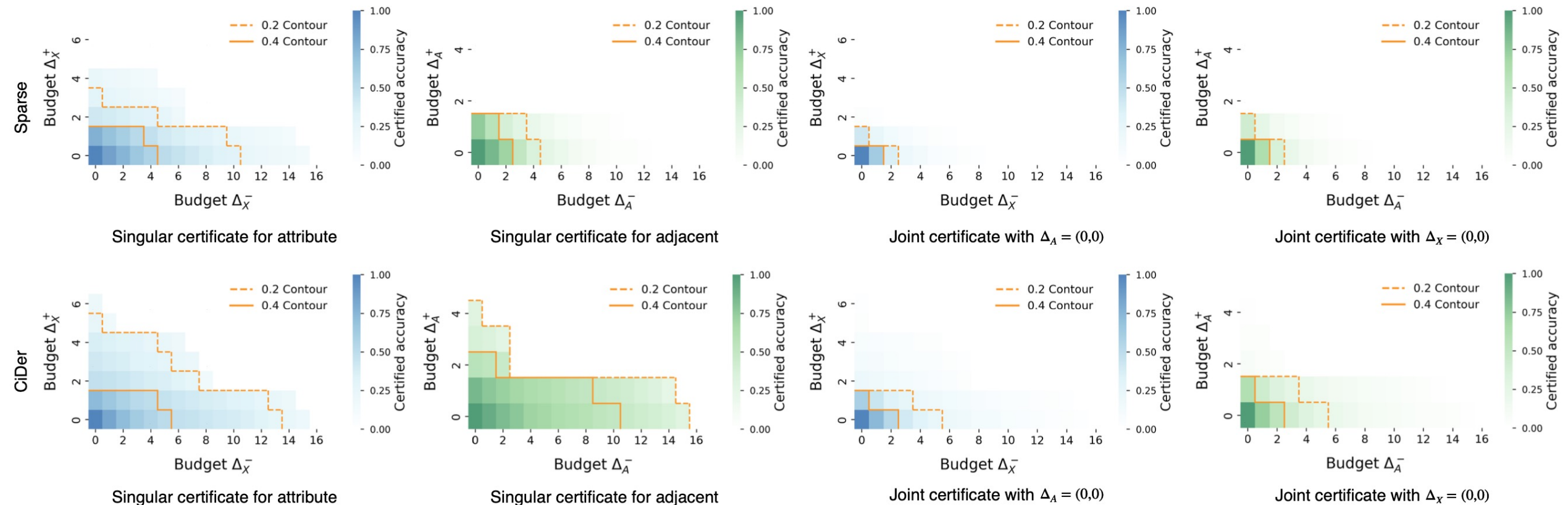
Higher accuracy across datasets, GNN architectures, and GNN layers.

Evaluation – Clean Accuracy



Higher accuracy across different perturbations and even than un-attacked GNN.

Evaluation – Certified Accuracy



Higher certified accuracy and **larger certifiable budget** under structure perturbation.

Summary



- **Black-box Certification:** Black-box adversary, black-box GNN.
- **Graph Joint Diffuser:** Adversarial graph purification with discrete diffusion model.
- **High Accuracy:** CiDer can provide robustness certificates with high accuracy.

- **Code:** <https://github.com/xiaoyuu-liang/CiDer>

Thank you!

Xiaoyu Liang

- xiaoyuliang@buaa.edu.cn
- Looking for a PhD position in Spring/Fall 2026.

Haohua Du, Wen Ma, Ye Tian, and Xiaoya Xu