

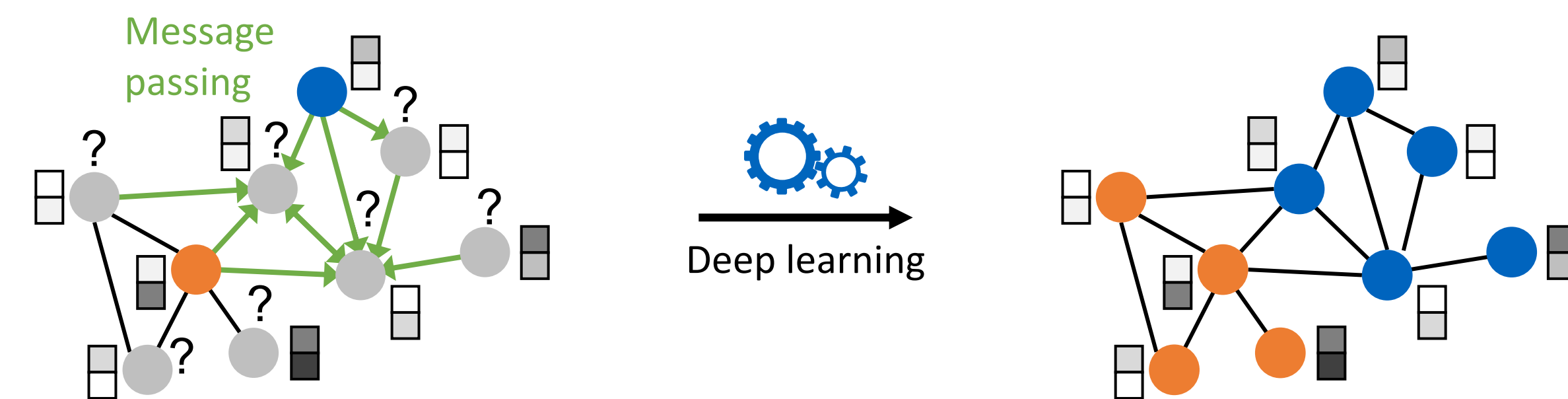
## Motivation

- Deep learning for graphs has brought great improvements on many tasks, e.g. node representation learning or node classification.
- Adversaries are a **real threat** in the applications where these models are likely to be used, e.g. search engines, recommender systems.
- Before deploying these models, we have to answer the question:

Are deep learning models for graphs robust with respect to **adversarial attacks**?

## Focus: Semi-supervised Node Classification

Partially labeled, attributed graph



- Message passing could lead to **higher robustness** due to local averaging; or **lower robustness** due to propagation effects in the network.

## Possible Attacks on Graphs

- In addition to feature vectors, adversarial attacks on graphs can **modify the graph structure**.
- Even **more critical**: a node can be attacked **without direct access** due to network effects.

**Target node**  $t \in V$ : node whose classification label the attacker wants to change

**Attacker nodes**  $S \subset V$ : nodes the attacker can modify

**Direct attack** ( $S = \{t\}$ )

- Modify the **target's** features
- Add connections to the **target**
- Remove connections from the **target**

**Example**

- Change website content
- Buy likes/followers
- Unfollow untrusted users

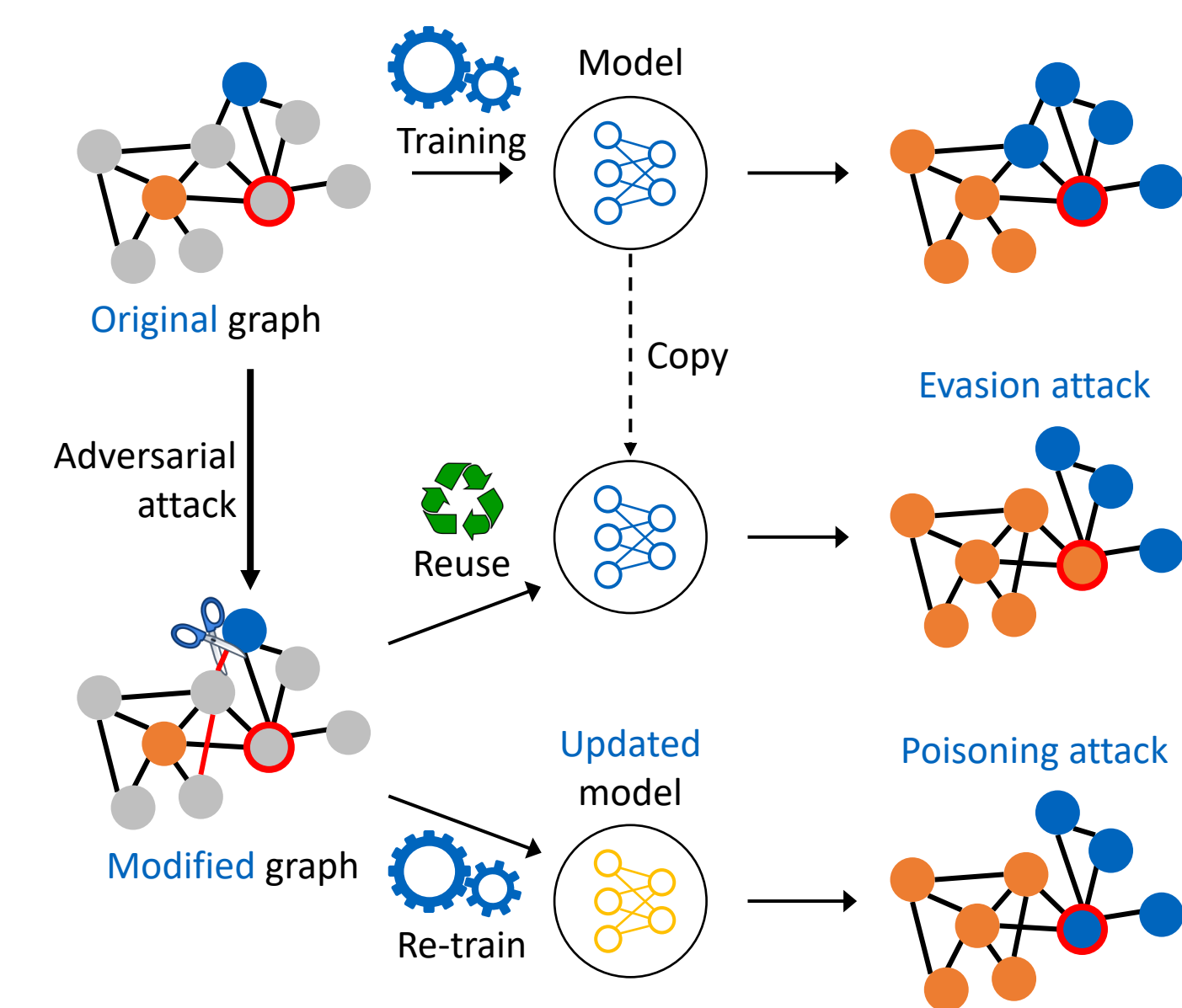
**Indirect attack** ( $t \notin S$ )

- Modify the **attackers'** features
- Add connections to the **attackers**
- Remove connections from the **attackers**

**Example**

- Hijack friends of target
- Create a link/spam farm

## Poisoning vs. Evasion Attacks



- Transductive learning**: data consists of labeled and unlabeled samples; all data used for training.
- Evasion attack**: Modify data to fool a **static classifier**.
- But**: modifications are on the training data (transductive setting).
- Re-training** can restore the predictions
- Important for a **realistic attack**: Impact **after model re-training** (poisoning).

## Formal Problem Definition

- Goal: **maximize classification loss** of a single target node.
- Measure **impact after training** the classifier on the modified data (**poisoning attack**).
- Enforce **unnoticeability constraints** for subtle perturbations.
- Challenge: **discrete data** (gradient less reliable because we need to extrapolate).

$$(A^c, X^c) = \arg \max_{A', X'} \max_{c \neq c_{old}} \log Z_{v,c}^* - \log Z_{v,c_{old}}^*$$

Classifier:  $where Z^* = f_{\theta^*}(A', X')$   $A \in \{0,1\}^{N \times N}$ : adjacency matrix

Training:  $with \theta^* = \arg \min_{\theta} \mathcal{L}(\theta; A', X')$   $X \in \{0,1\}^{N \times D}$ : node attributes

Unnoticeability:  $s.t. (A', X') \approx (A, X)$   $A'$ : modified structure

Budget constraint:  $and \|A' - A\|_0 + \|X' - X\|_0 < \Delta$   $X'$ : modified features

$v$ : target node

## Surrogate Model

- Linear **surrogate model** based on two-layer GCN.
- Enables computation of the **exact impact** of a perturbation **efficiently** and in **closed form**.
- Attacker chooses perturbation that **maximizes loss** on the surrogate model (one at a time).

Linearize classifier:  $Z = f_{\theta}(A, X) = \text{softmax}(\hat{A} \hat{X} W^{(1)}) W^{(2)}$

Simplified equation:  $\log Z' = \hat{A}^2 X W'$

**Structure** perturbations:  $\max_{\hat{A}} \mathcal{L}'(\log Z'_v) where \log Z'_v = [\hat{A}^2 C]_v$  **Constants**

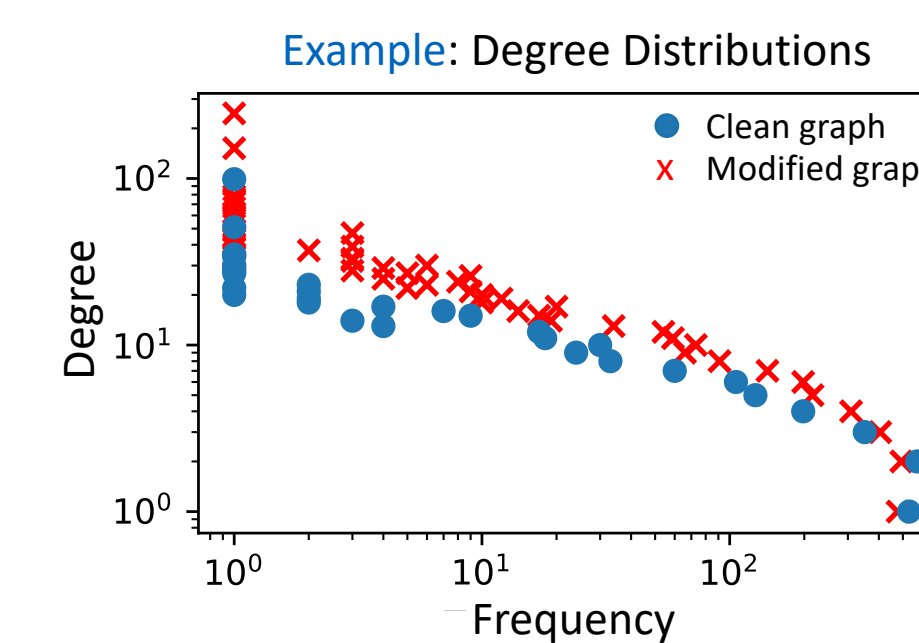
**Feature** perturbations:  $\max_{X'} \mathcal{L}'(\log Z'_v) where \log Z'_v = [C_1 X C_2]_v$

## Unnoticeability Constraints

$(A', X') \approx (A, X)$ : What are sensible measures of 'closeness' for graphs?

### Structure perturbations

- Fundamental property of graphs: **degree distribution**.
- Hypothesis test**: Were the original and modified degree distributions  $D$  and  $D'$  generated by the **same underlying powerlaw distribution**?



$H_0$ :  $D$  and  $D'$  come from the **same powerlaw** distribution.  
 $H_1$ : they come from **separate powerlaw** distributions.

Approx. power law exponent:  $\alpha(D) \approx 1 + |D| \cdot \left[ \sum_{d_i \in D} \log \frac{d_i}{d_{min} - 0.5} \right]^{-1}$

Log likelihood:  $l(D) = |D| \cdot \log \alpha [1 + \log d_{min}] + [\alpha + 1] \sum_{d_i \in D} \log d_i$

Hypothesis likelihoods:  $l(H_0) = l(D \cap D'); l(H_1) = l(D) + l(D')$

Test statistic:  $\Lambda(D, D') = -2 \cdot l(H_0) + 2 \cdot l(H_1)$

- For large  $N$ ,  $\Lambda$  follows a  $\chi^2$  distribution with one degree of freedom.
- We choose a **very conservative** p-value constraint ( $p=0.95$ ), i.e. we can detect **small changes**.
- Bookkeeping** enables **incremental** updates to all terms in constant time.

### Feature perturbations

- Constraint to prevent addition of **unrealistic features** (words).
  - Co-occurrence graph** of the node attributes in the data:
- $$C = (F, E), F \in \{0,1\}^P, E \subseteq F \times F$$
- $(f_1, f_2) \in E$ :  $f_1$  and  $f_2$  co-occur in at least one node in the dataset.  
Node  $u$ 's features:  $S_u = \{j \mid X_{uj} \neq 0\}$

Probability of reaching a feature by a **random walker** on  $C$  starting at node  $u$  in one step:

$$p(i|S_u) = \frac{1}{|S_u|} \sum_{j \in S_u} 1/d_j \cdot E_{ij} > \sigma; \text{ in our experiments we set } \sigma = 0.5 \cdot \frac{1}{|S_u|} \sum_{j \in S_u} 1/d_j.$$

- Encourages** addition of features that **co-occur** with many of node  $u$ 's features.
- Discourages** addition of features that only co-occur with unspecific features, i.e. **stopwords**.

Top inserted words to papers of the class neural networks:

| w/ constraint | w/o constraint |
|---------------|----------------|
| probabilistic | efforts        |
| bayesian      | david          |
| inference     | family         |

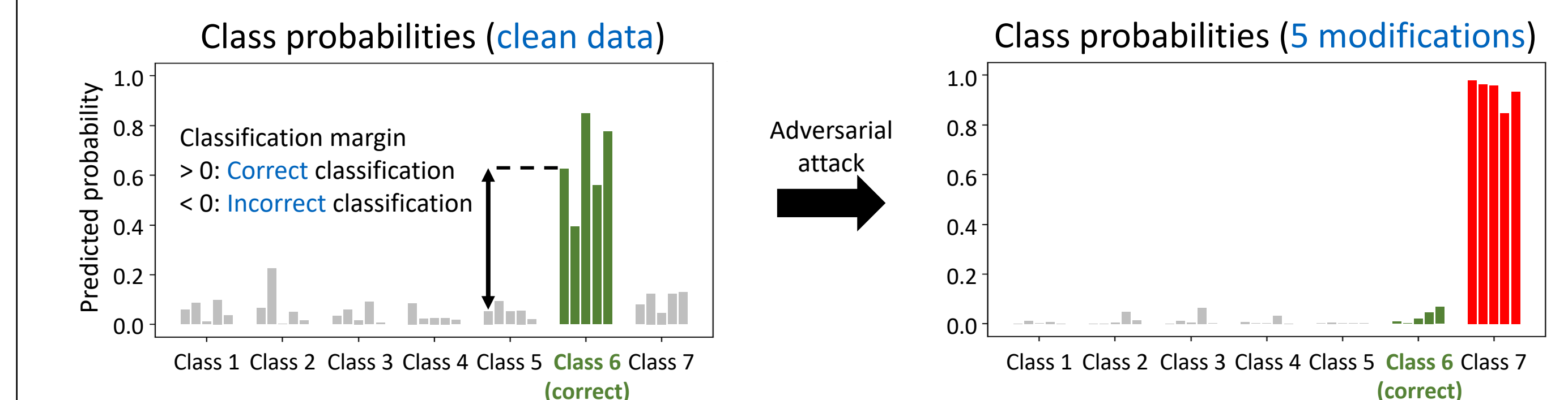
## Experimental Results

- Transfer experiments**: Graph Convolutional Network (GCN), Column Network, DeepWalk.
- Perturbation budget is  $d+2$ , where  $d$  is the target's degree.
- Evaluation on 5 different splits; 10x re-training per attack.

### Approach

- Attacks via **linear surrogate model** while enforcing **unnoticeability constraints**.
- Train deep learning models on the **poisoned data** and evaluate the drop in performance.
- No access** to the classifier is needed for the attack.

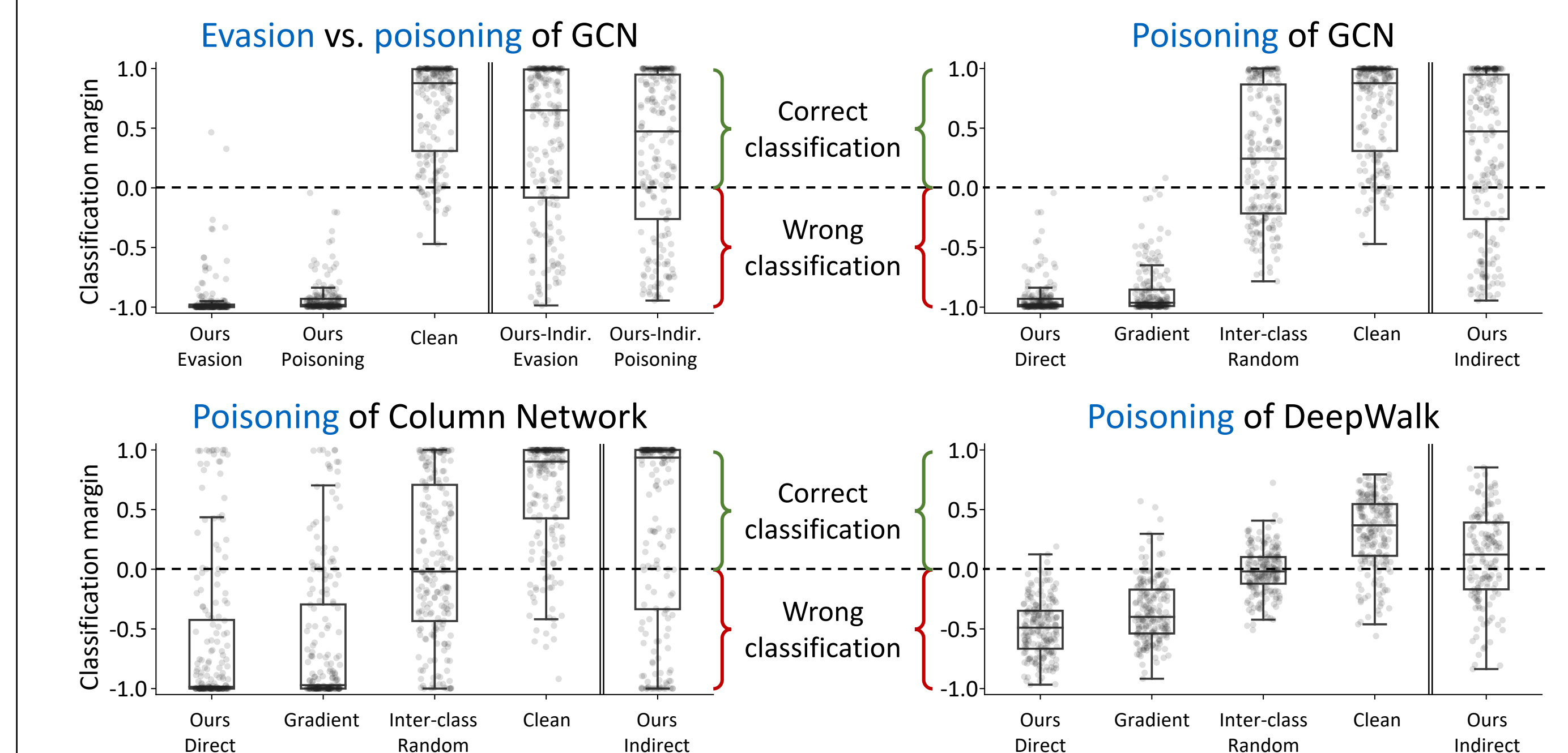
### Example Attack



## Results

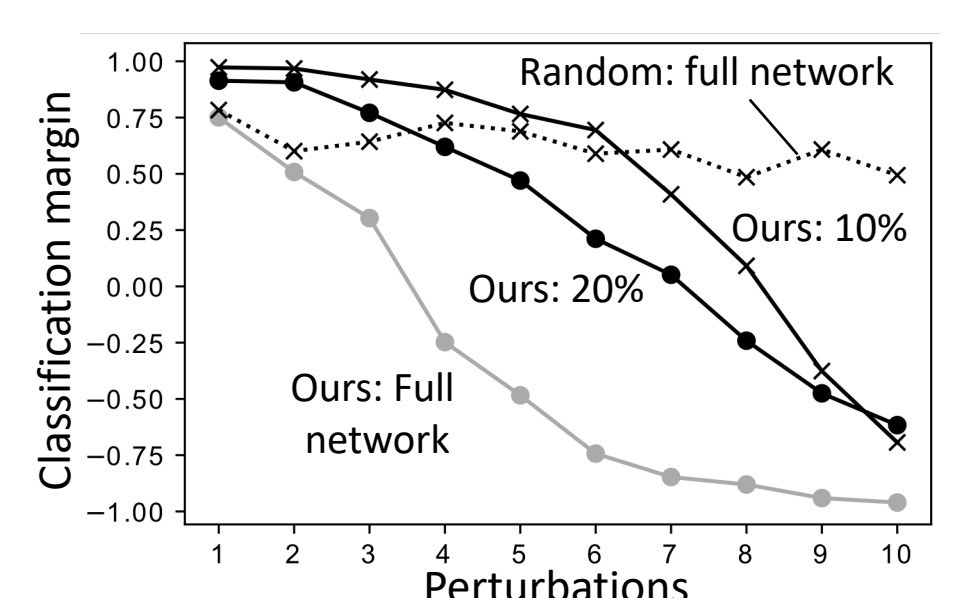
|                    | Cora        |             |             | Citeseer    |             |             | PolBlogs    |             |             |
|--------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Attack method      | GCN         | CLN         | DW          | GCN         | CLN         | DW          | GCN         | CLN         | DW          |
| Clean              | 0.90        | 0.82        | 0.84        | 0.88        | 0.71        | 0.76        | 0.93        | 0.63        | 0.92        |
| Ours               | <b>0.01</b> | <b>0.17</b> | <b>0.02</b> | <b>0.02</b> | <b>0.20</b> | <b>0.01</b> | <b>0.06</b> | <b>0.47</b> | <b>0.06</b> |
| Gradient           | 0.03        | 0.18        | 0.10        | 0.07        | 0.23        | 0.05        | 0.41        | 0.55        | 0.37        |
| Inter-class Random | 0.61        | 0.52        | 0.46        | 0.60        | 0.52        | 0.38        | 0.36        | 0.56        | 0.30        |
| Ours-Indirect      | 0.67        | 0.68        | 0.59        | 0.62        | 0.54        | 0.48        | 0.86        | 0.62        | 0.91        |

Table: Share of correct classifications of target nodes **after attack** and **re-training**



### Limited knowledge

- Only provide a **small part of the network** around the target node to the surrogate model to attack.
- Evaluation with GCN on the **complete graph**.
- Attacks are successful even in this restricted setting.



## Conclusion

- We propose an **efficient algorithm** for adversarial attacks on deep learning for graphs.
- No access** to the classifier is needed for the attack due to the surrogate model approach.
- Our attacks are successful even under **restrictive attack scenarios**, e.g. no access to target node or limited knowledge about the graph.

