

TL;DR: First certificate w.r.t. graph perturbations for a general class of models including Label Prop. and GNNs.

GNNs are vulnerable to Adversarial Attacks

Semi-Supervised Node Classification: Given a few labelled nodes predict the classes of the remaining nodes in the graph
Targeted Attack: Perturb the graph to **misclassify** a target node

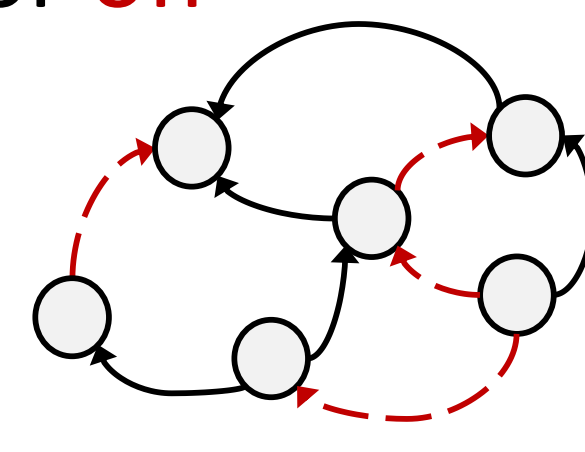
Research Questions

Certification: How to **verify** if a graph-based model is robust?
Robust Training: How can we **improve** certified robustness?

Flexible Threat Model

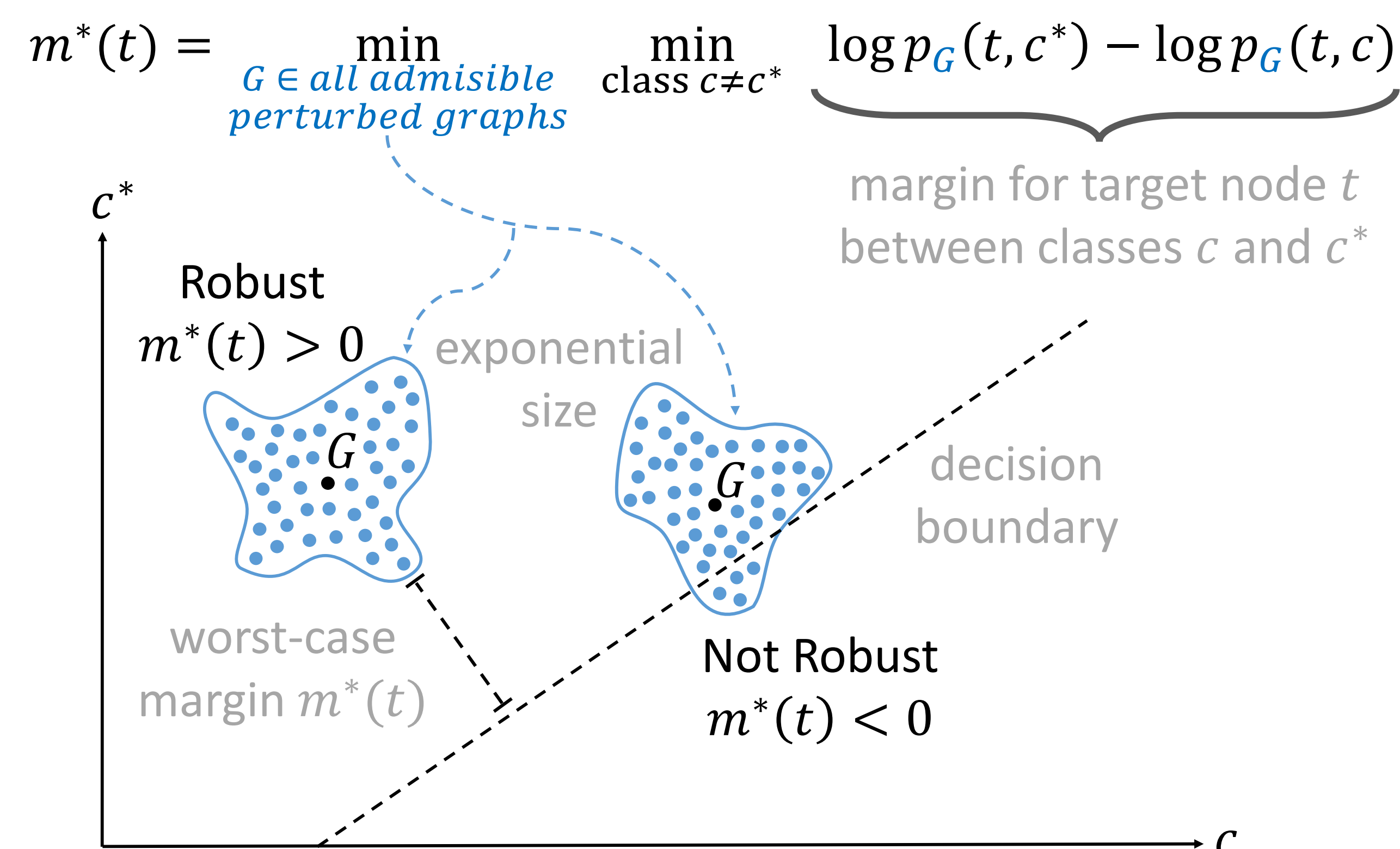
Attacker controls **fragile** edges they can turn **on** or **off**

Global Budget: perturb at most B edges in total
Local Budget: at most b_v edges for each node v



Robustness Certificate

Guarantee that the prediction does not change under **any** admissible perturbation of the input graph



Family of Models based on PageRank

Predictions are a linear function of (Personalized) PageRank

$$\log p_G(t, c) = \pi_G(t)^T h(c)$$

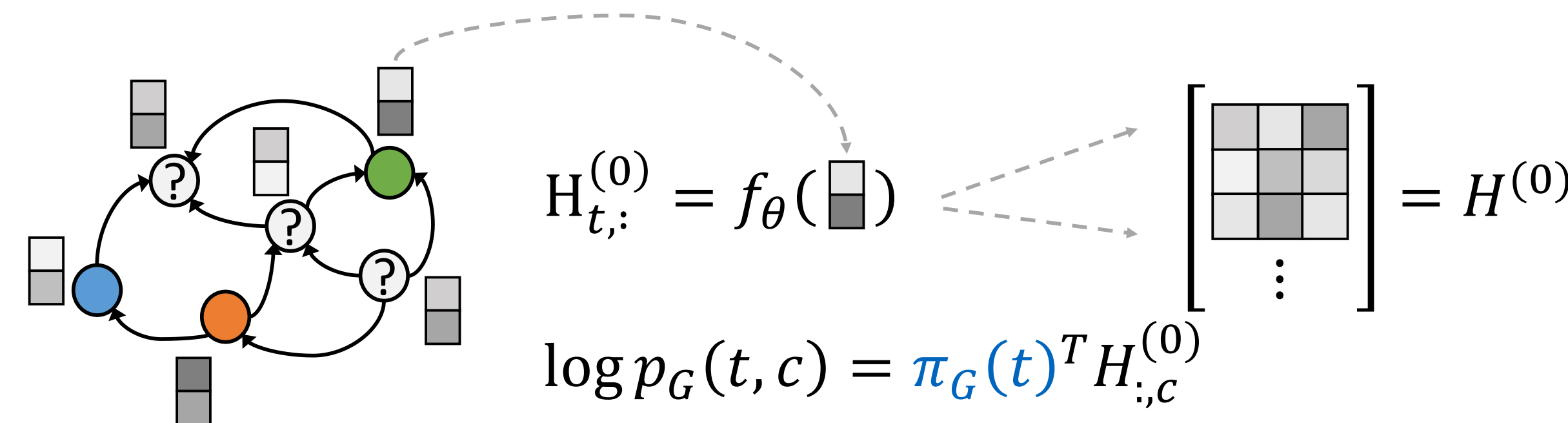
Personalized PageRank $\pi(t)$: Stationary distribution of a random walker teleporting back to node t with probability α

EX 1 - Label Propagation: repeatedly diffuse initial beliefs $H^{(0)}$

$$H^{(0)} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \\ \vdots & \vdots & \vdots \end{bmatrix} \quad H^{(t+1)} = (1 - \alpha)D^{-1}AH^{(t)} + \alpha H^{(0)}$$

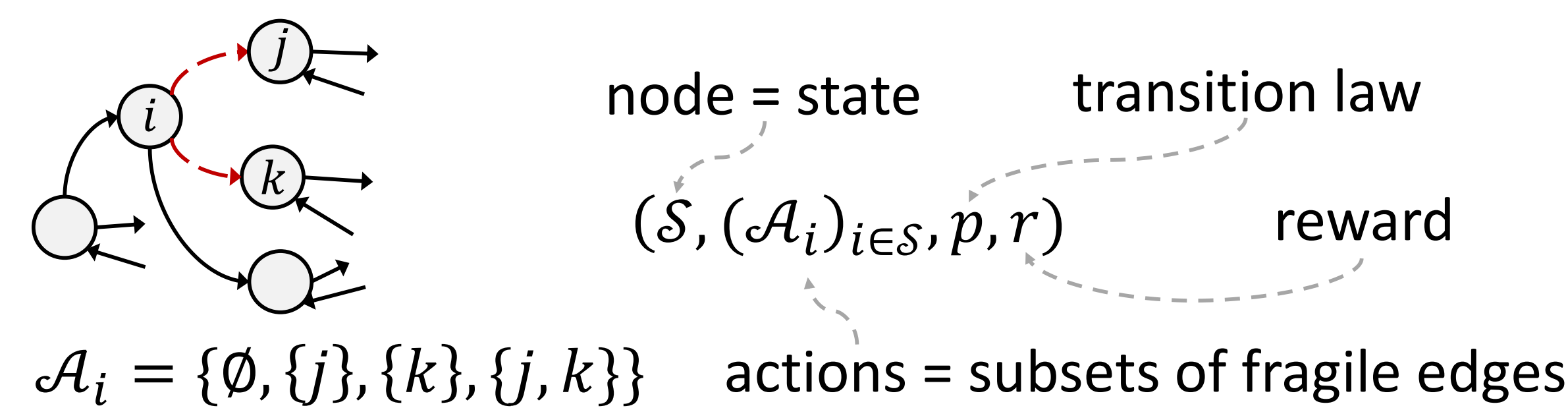
$$\log p_G(t, c) = H_{t,c}^{(\infty)} = \pi_G(t)^T H_{:,c}^{(0)}$$

EX 2 - Graph Neural Network (PPNP): first map node features to initial beliefs with a NN f_θ then diffuse with $\pi_G(t)$



Certificate \Leftrightarrow PageRank Optimization \Leftrightarrow MDP

Computing certificates amounts to finding optimal PageRank which can be done efficiently via a Markov Decision Process



Local budget: find **optimal** fragile edges with policy iteration

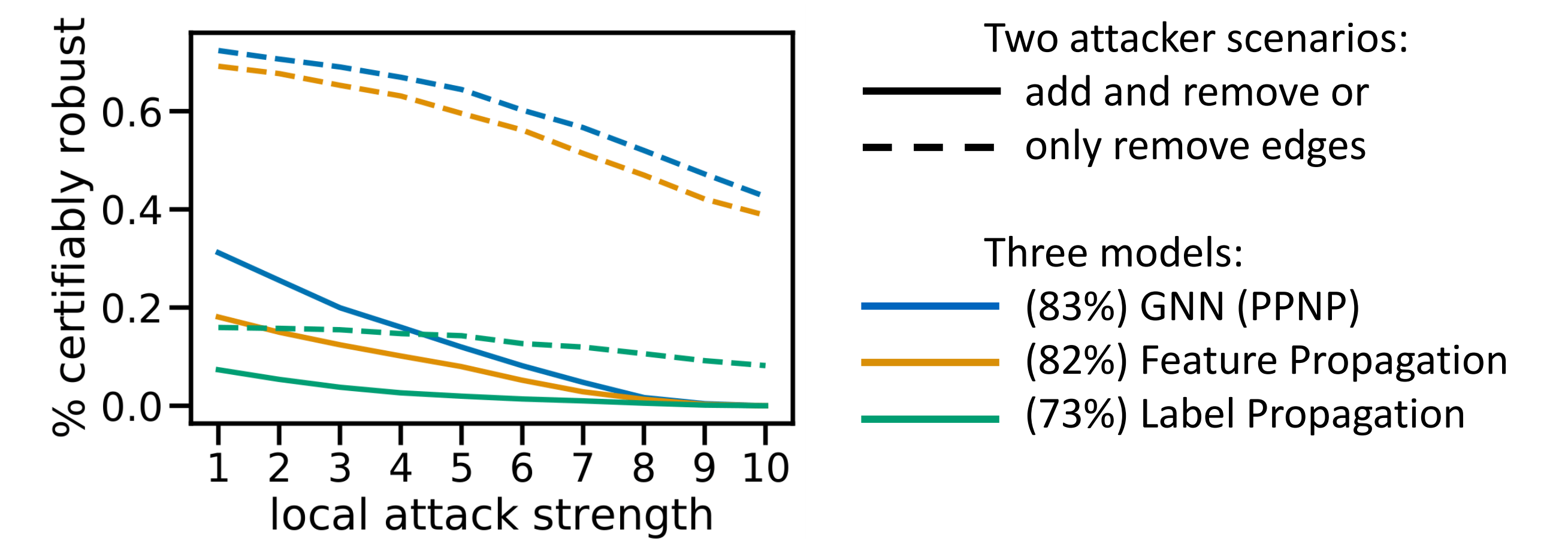
$$r = H_{c^*}^{(0)} - H_c^{(0)}$$

set reward to the logit difference

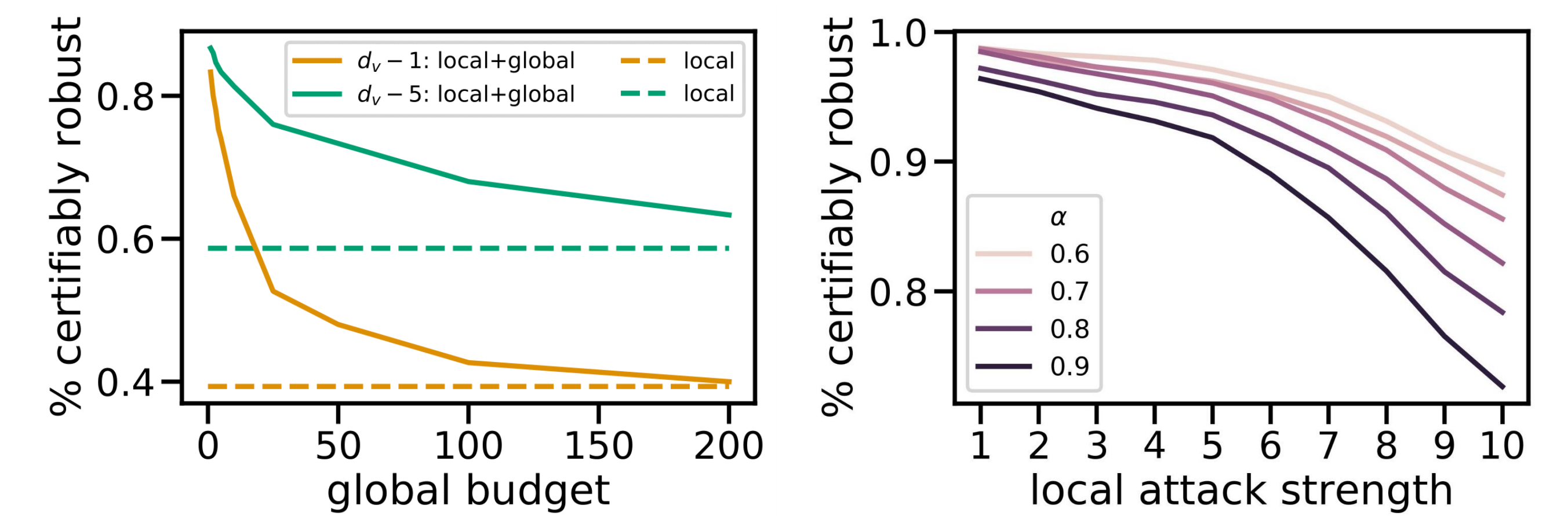
Local + Global budget: **NP-Hard**, augment graph & solve a QP

Certification Results

GNNs are more robust than Label/Feature Propagation



To increase ratio of certified nodes: decrease budget / lower α



Robust Training

Use worst-case margin during training to learn robust weights

Hinge-loss penalty: maximize the worst-case margin

Robust cross-entropy: worst-case instead of standard logits

$$\mathcal{L}_{CEM} = \mathcal{L}_{CE} + \sum_{c \neq c^*} \max(0, M - m^*) \quad \mathcal{L}_{RCE} = \mathcal{L}_{CE}(-m^*)$$

Robust training increases ratio of certified nodes and accuracy

