

TL;DR

- Label poisoning for GNNs is plagued by serious evaluation pitfalls.
- Existing attacks render ineffective post fixing these fallacies.
- We introduce two new simple yet effective family of attacks that are significantly stronger (up to 8%) than previous strongest attacks.

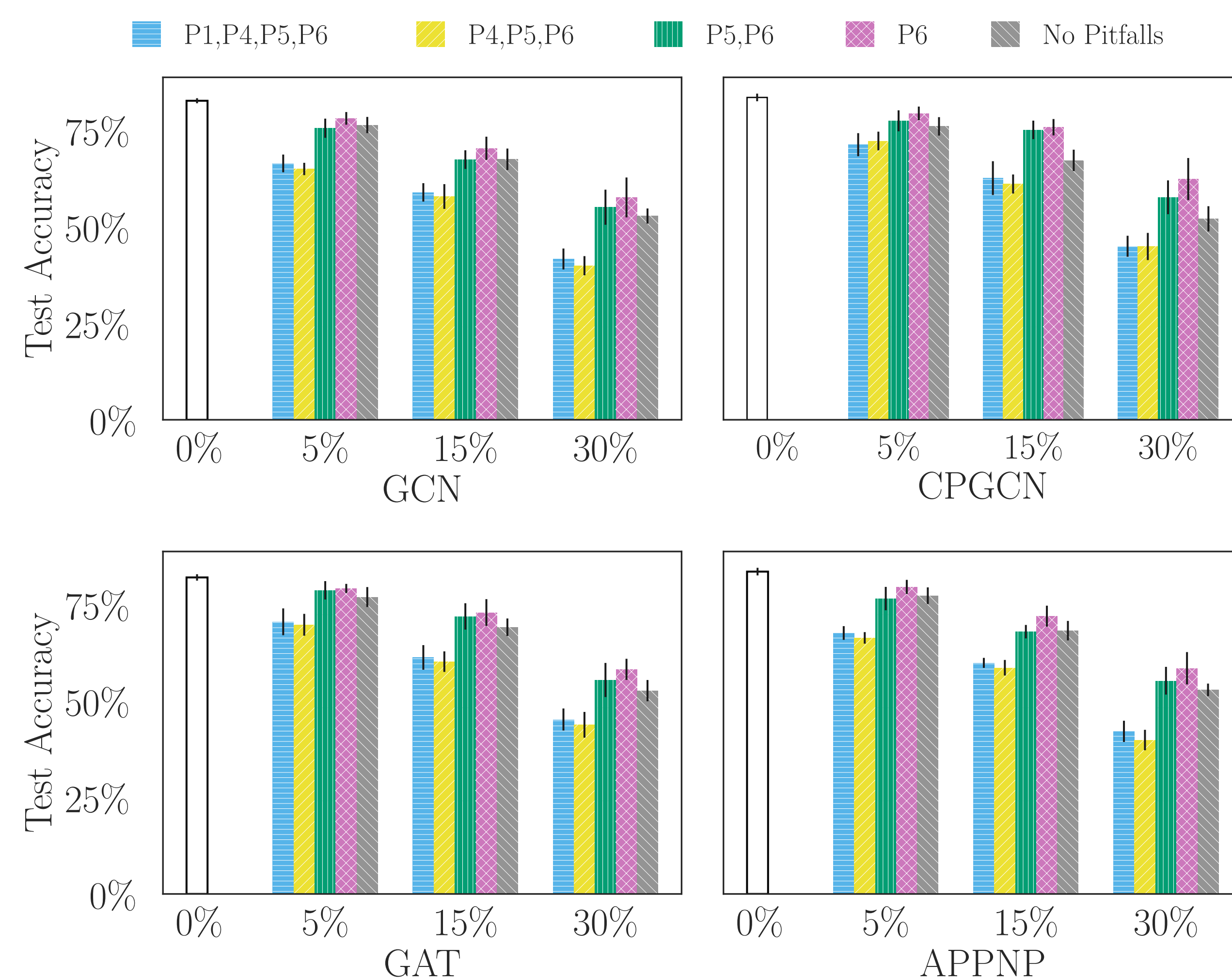
Motivation

GNNs have wide range of applications including critical ones.

Label poisoning poses a distinct threat as training data can be compromised.

Existing attacks are not effective; do better attacks exist?

Existing attacks are not as powerful as claimed



- P1:** Large Validation Set
- P2:** Missing stdev
- P3:** Eval. on undefended models
- P4:** Class equalised splits
- P5:** Hyper-parameter tuning
- P6:** Clean Validation set

Fixing the above pitfalls leads to a massive reduction in LafAK's performance (previous strongest attack).

Threat Model and Baselines

Flip a small fraction of labels to decrease test acc.

Results in a difficult bi-level optimization problem for which we propose different relaxations.

We used two family of attacks as the baselines:

- Heuristic-based:** Random (**RND**), Degree (**DEG**)
- Learning-based:** LP, LafAK (**LFK**), MG

Linear surrogate attacks

Linearize the classifier and compute the optimal weights in closed-form

$$\begin{aligned} \min_{\mathbf{H} \in \{0,1\}^{L \times C}} \mathcal{L}(\mathbf{Y}_u, \widehat{\mathbf{Y}}_u) \\ \|\mathbf{H} - \mathbf{Y}_l\|_0 \leq 2\epsilon L \\ \widehat{\mathbf{Y}}_u = \widehat{\mathbf{X}}_u \widetilde{\mathbf{X}}_l \mathbf{H} \\ \mathbf{H} \mathbf{1}_C = \mathbf{1}_L \end{aligned}$$

where: $\widetilde{\mathbf{X}} = (\widehat{\mathbf{X}}^T \widehat{\mathbf{X}} + \lambda \mathbf{I})^{-1} \widehat{\mathbf{X}}$ is the closed form solution of LR.

Variant-1: SGC surrogate $\widehat{\mathbf{X}} = \widehat{\mathbf{A}}^2 \mathbf{X}$

Variant-2: NTK surrogate $\widehat{\mathbf{X}} = \text{NTK}$ - Kernel

Proposition: LSA closed form solution

Given fixed target labels $\widetilde{\mathbf{Y}}_l$, the optimal nodes to poison are the subset of nodes corresponding to the smallest $\lfloor \epsilon L \rfloor$ negative elements of an L -dimensional vector \mathbf{c} , where the l -th element of \mathbf{c} is computed as $c_l = \sum_{ij} Q_{il} P_{lj} R_{ij}$ where $\mathbf{Q} = \widehat{\mathbf{X}}_u \widetilde{\mathbf{X}}_l$, $\mathbf{P} = \widetilde{\mathbf{Y}}_l - \mathbf{Y}_l$, and $\mathbf{R} = \mathbf{Y}_u$.

Meta attacks

Meta gradients w.r.t. labels by backpropagating through the unrolled inner optimization. The poisoned labels are constructed as follows:

$$\mathbf{H} = \text{diag}(\mathbf{b}) \widetilde{\mathbf{Y}} + \text{diag}(\mathbf{1}_L - \mathbf{b}) \mathbf{Y}_l$$

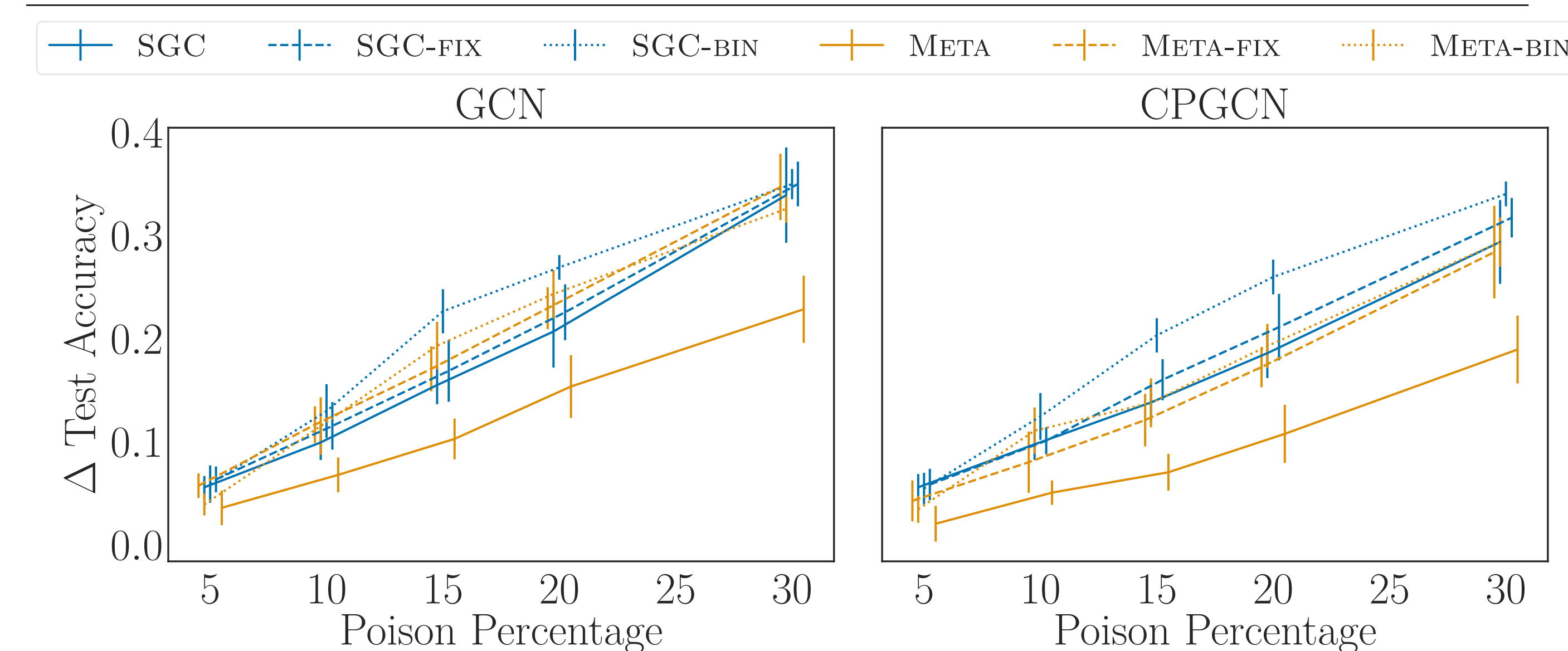
$$\begin{aligned} \text{where: } \widetilde{\mathbf{Y}} &= \text{GumbleSoftmax}(\widetilde{\mathbf{Y}}_{\log}); \widetilde{\mathbf{Y}}_{\log} \in \mathbb{R}^{N \times C} \\ \mathbf{b} &= \text{top}_k(\widetilde{\mathbf{b}}); \mathbf{b} \in \mathbb{R}^N \end{aligned}$$

Note: since topk is not differentiable, we apply soft-top-k followed by k-subset selection.

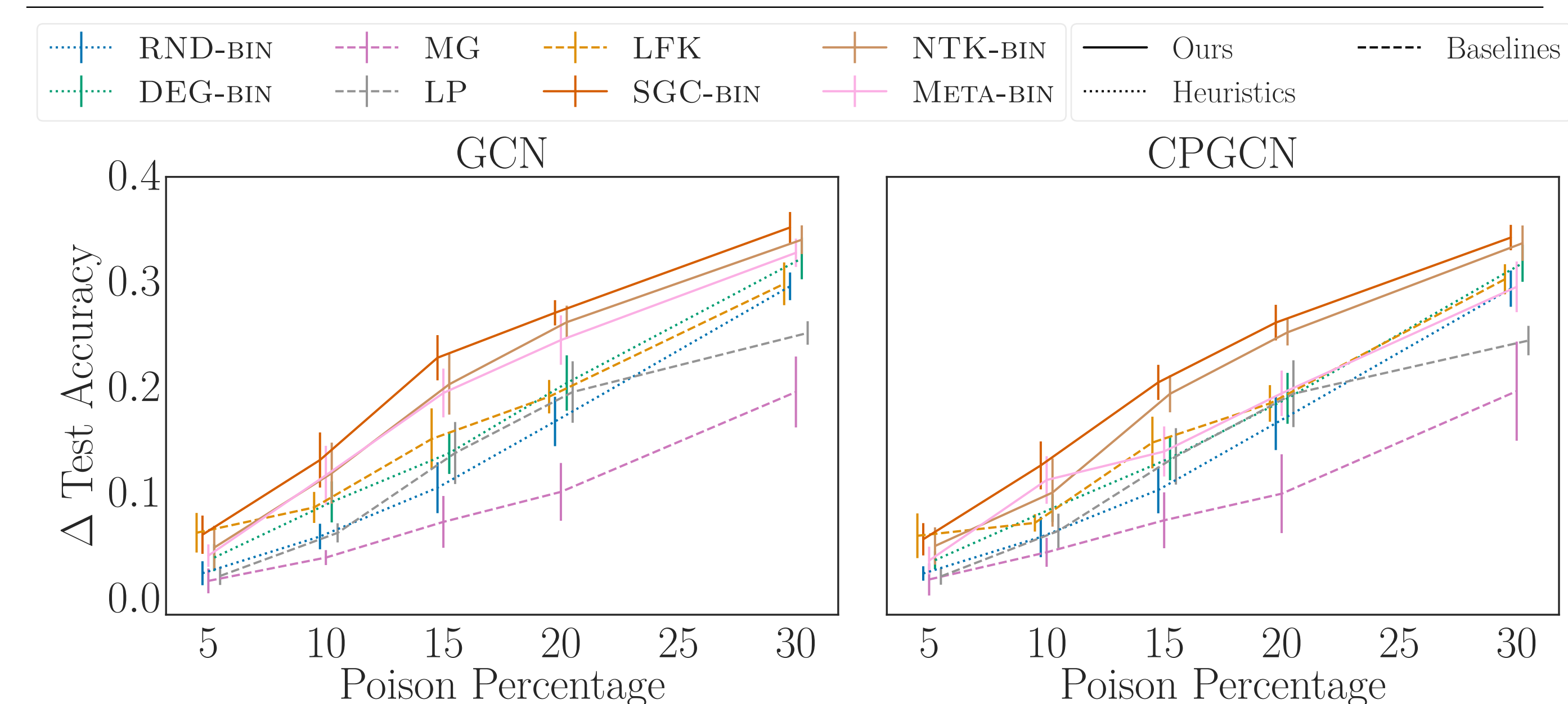
Proposition: Optimality of binary random attack

Let the adversary flip label p to label $q \neq p$ with probability $\frac{\epsilon}{s} \cdot t_{pq}$ and retain label p with probability $1 - \epsilon$, where ϵ is the poisoning budget, $t_{pq} \in \{0, 1\}$ indicates whether the adversary is allowed to flip p to q , and $s = \sum_{q \neq p} t_{pq}$ is the number of allowed classes. The test accuracy of the Bayes optimal classifier trained on randomly flipped labels is minimized for $s = 1$ (binary flips).

LSA outperforms meta & Binary outperforms multi-label



Our proposed attacks significantly outperform baselines



Key takeaways

- Faithfully simulating the defender is crucial to evaluate the efficacy of an attack.
- Simple label poisoning attacks are surprisingly powerful.
- Our findings highlight the need to further study label poisoning attacks as well as develop defences.