

## TL;DR

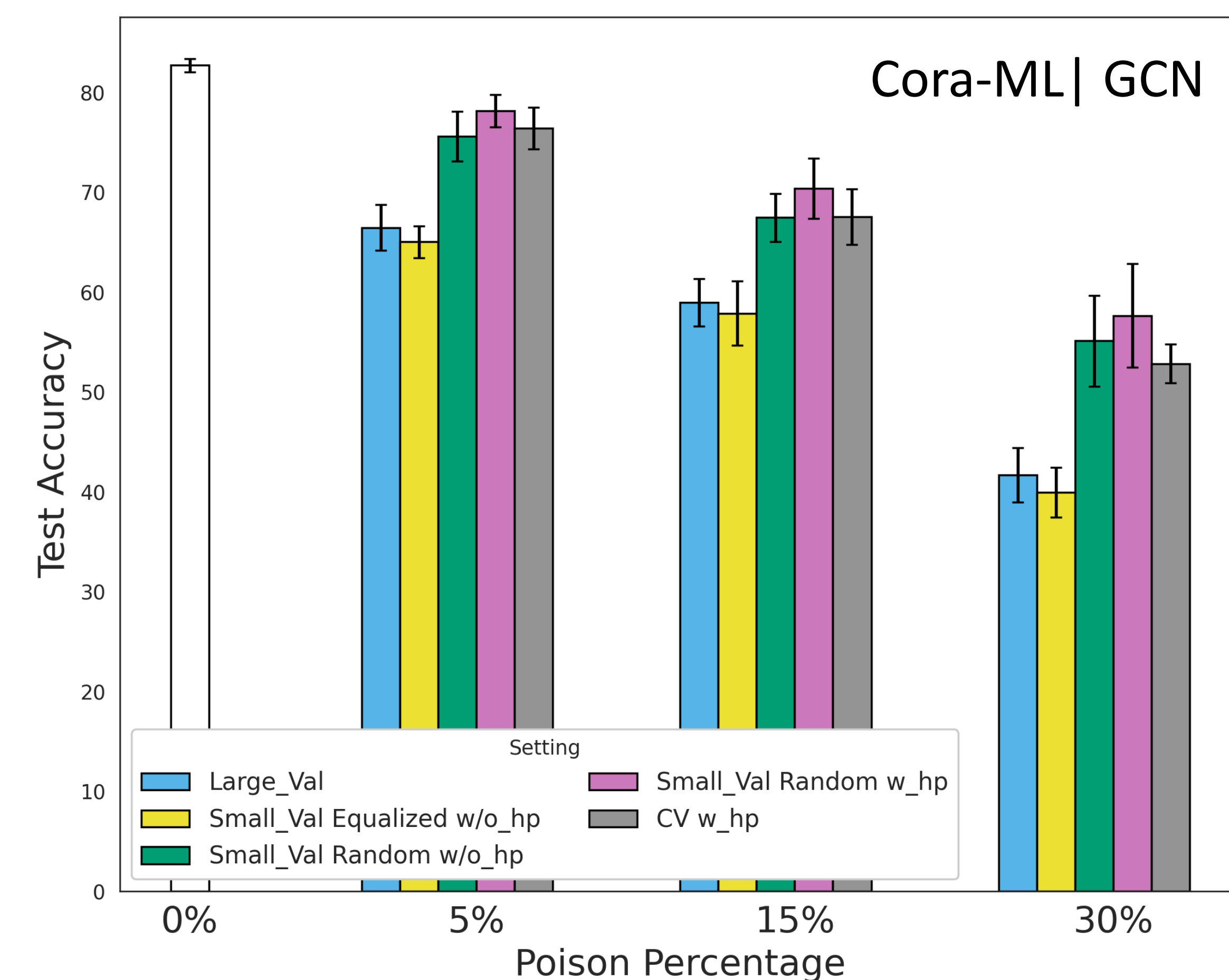
- Label poisoning for GNNs is plagued by serious evaluation pitfalls.
- Existing attacks render ineffective post fixing these fallacies.
- We introduce two new simple yet effective family of attacks that are significantly stronger (up to **~8%**) than previous strongest attacks.

## 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



- Large validation set
- Class equalised splits
- Hyper-param tuning
- Clean Validation set
- Missing stdev

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

## Threat Model

Flip a small fraction of labels to decrease test acc. Results in a difficult bi-level optimization problem for which we propose different relaxations.

## 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

$$\begin{aligned} \min_{\mathbf{H} \in \{0,1\}^{n \times C}, \mathbf{b} \in \{0,1\}^n} \mathcal{L}(\mathbf{Y}_u, \tilde{\mathbf{Y}}_u) \\ \hat{\mathbf{Y}}_l = \mathbf{b} \odot \mathbf{H} + (\mathbf{1}_N - \mathbf{b}) \odot \mathbf{Y}_l \quad \# \text{poisoned labels} \\ \mathbf{H} \neq \mathbf{Y}_l \\ \tilde{\mathbf{Y}}_u = \hat{\mathbf{X}}_u \tilde{\mathbf{X}}_l \hat{\mathbf{Y}}_l \quad \# \text{compute predictions} \\ \mathbf{H} \mathbf{1}_C = \mathbf{1}_N \quad \# \text{one-hot constraint} \\ \mathbf{b}^T \mathbf{1}_N \leq \epsilon N \quad \# \text{enforces budget} \end{aligned}$$

$$\tilde{\mathbf{X}} = (\hat{\mathbf{X}}^T \hat{\mathbf{X}} + \lambda I)^{-1} \hat{\mathbf{X}}^T \quad \# \text{Closed form solution of LR}$$

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

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

## Meta attacks

Meta gradients w.r.t. labels by backpropagating through the unrolled inner optimization.

Final poisoned labels are constructed as follows:

$$\begin{aligned} \hat{\mathbf{Y}}_l = \mathbf{b} \odot \mathbf{H} + (\mathbf{1} - \mathbf{b}) \odot \mathbf{Y}_l \\ \text{where } \mathbf{H} = \text{GumbelSoftmax}(\tilde{\mathbf{H}}); \quad \tilde{\mathbf{H}} \text{ in } \mathbb{R}^{N \times C} \\ \mathbf{b} = \text{top}_k(\tilde{\mathbf{b}}); \quad \tilde{\mathbf{b}} \in \mathbb{R}^N \end{aligned}$$

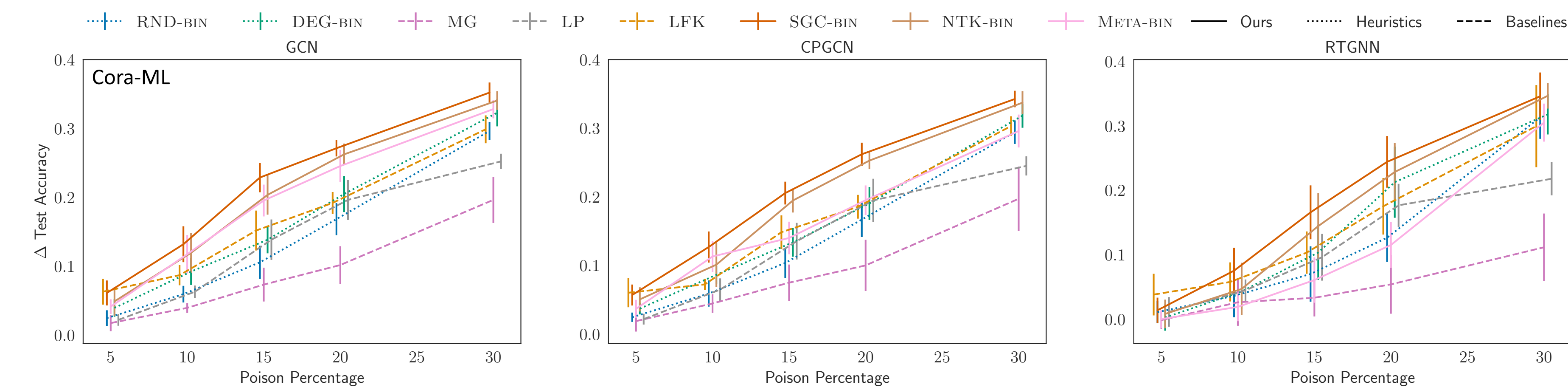
Note: since  $\text{top}_k$  is not differentiable, we apply soft-top-k followed by k-subset selection.

## Gumbel-softmax trick to approximate 0-1 loss

To make the 0-1 loss differentiable, we propose the following simple alternative:

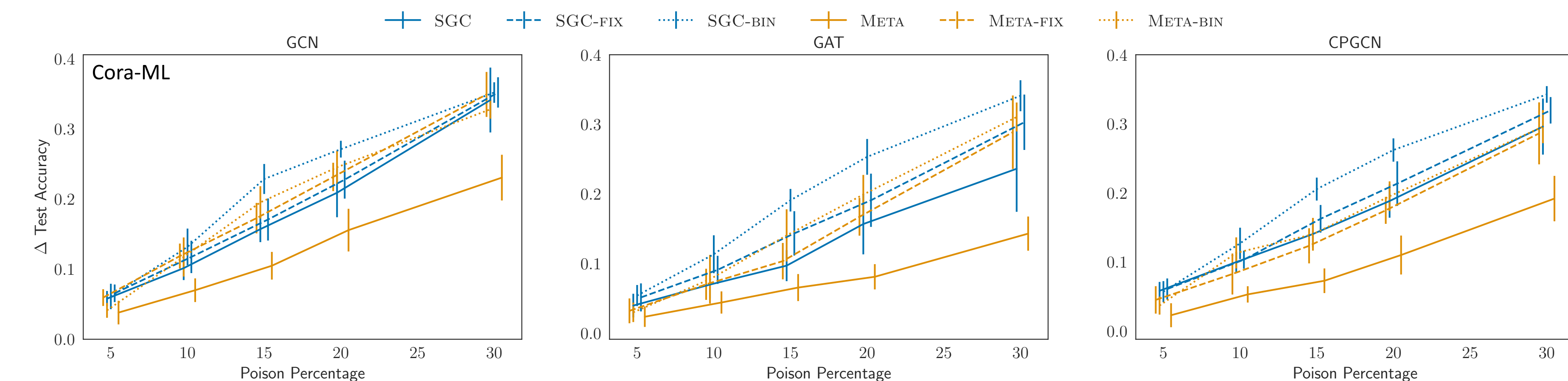
$$GS_{loss} = \frac{1}{|u|} \sum \text{GumbelSoftmax}(\hat{\mathbf{Y}}_u) \odot \mathbf{Y}_u$$

## Our proposed attacks significantly outperform baselines



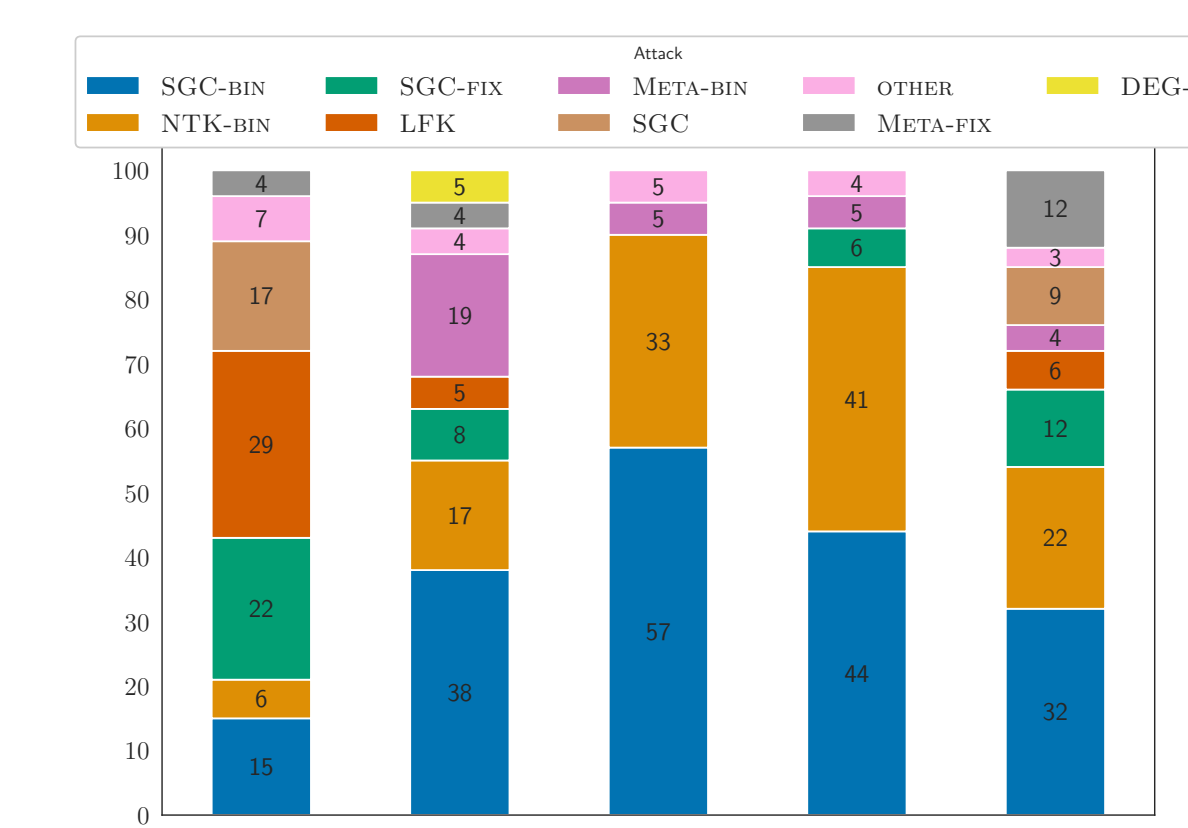
All our attacks (in solid), particularly SGC-BIN, outperform baseline attacks with max gains of up to **~13%**.

## Linear surrogate outperform meta attacks



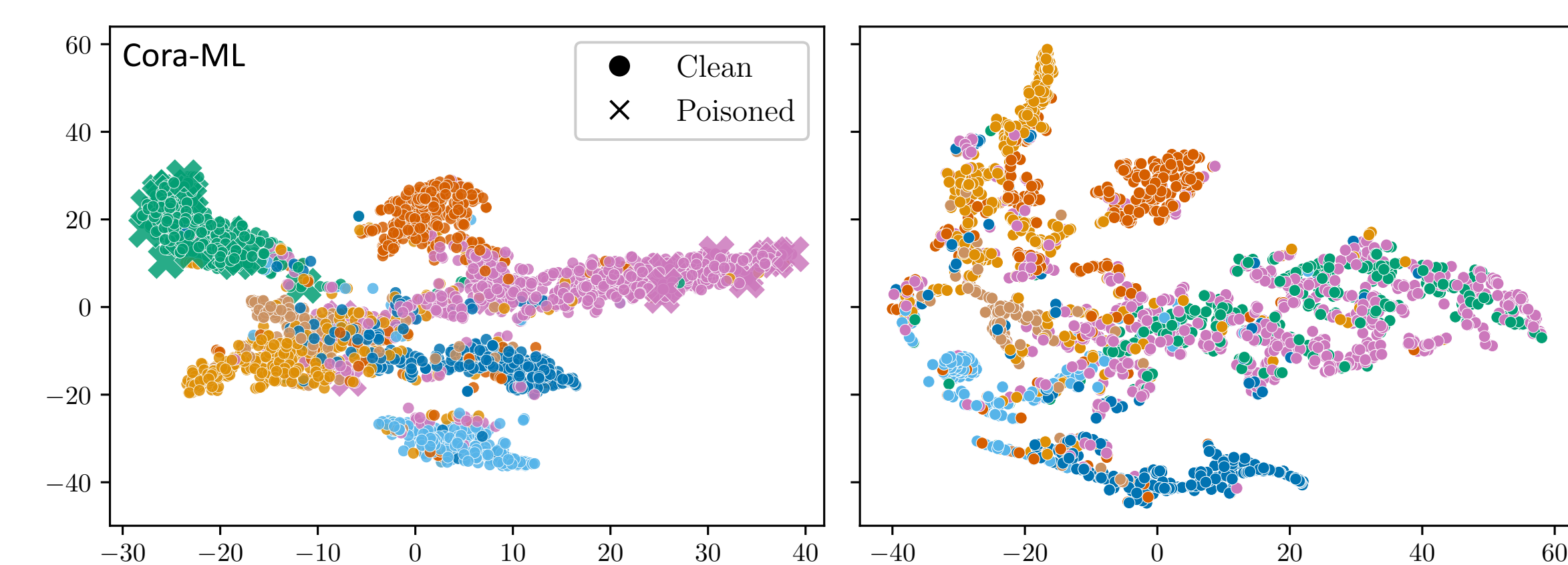
Different variants of the attacks that we propose. Binary variants are best on average.

## Strongest attack per split



Different variants of our linear attacks win most often.

## t-SNE of the logits before and after poisoning



Poisoning only a handful of training labels disrupts the learned representations compared to clean model representations.

## Key takeaways

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