

Sparse but Strong: Crafting Adversarially Robust Graph Lottery Tickets

Subhajit Dutta Chowdhury¹, Zhiyu Ni¹, Qingyuan Peng¹, Souvik Kundu².

Pierluigi Nuzzo¹

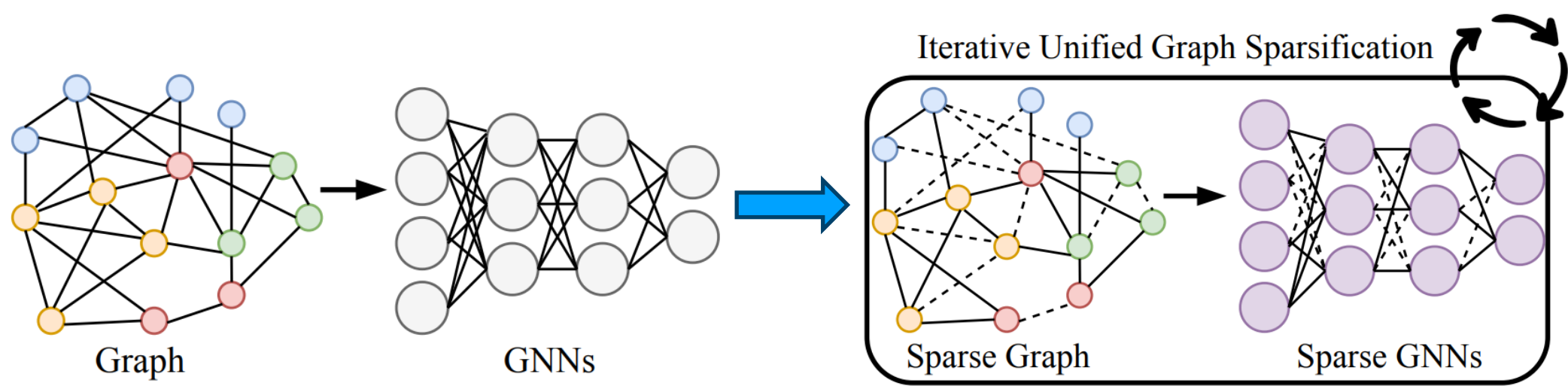
¹University of Southern California, Los Angeles, USA

²Intel Labs, San Diego, USA

Problem Statement

Are the existing graph lottery ticket (GLT) identification techniques capable of generating adversarially robust GLTs?

Graph Lottery Tickets in a Nutshell



Unified graph sparsification (UGS) iteratively removes edges and weights from the graph and GNN, respectively. Image courtesy: T. Chen et al., "A Unified Lottery Tickets Hypothesis for Graph Neural Networks," ICML 2021

Adversarially Robust Graph Sparsification

The proposed technique performs unified graph-GNN sparsification such that:

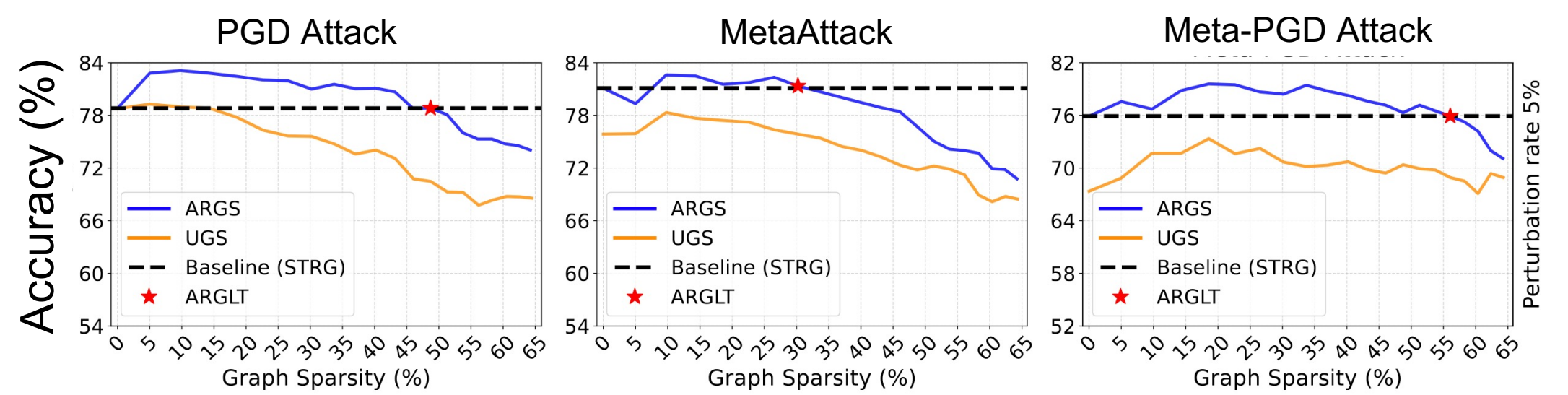
- Promotes feature smoothness
- Leverages self-training to remove edges from the locality of the test nodes

$$\mathcal{L}_{\text{ARGS}} = \mathcal{L}_0(f(\{m_g \odot A', X\}, m_\theta \odot \Theta)) + \beta \mathcal{L}_{f_s}(m_g \odot A', X) + \mathcal{L}_1(f(\{m_g \odot A', X\}, m_\theta \odot \Theta)) + \lambda_1 \|m_g\|_1 + \lambda_2 \|m_\theta\|_1$$

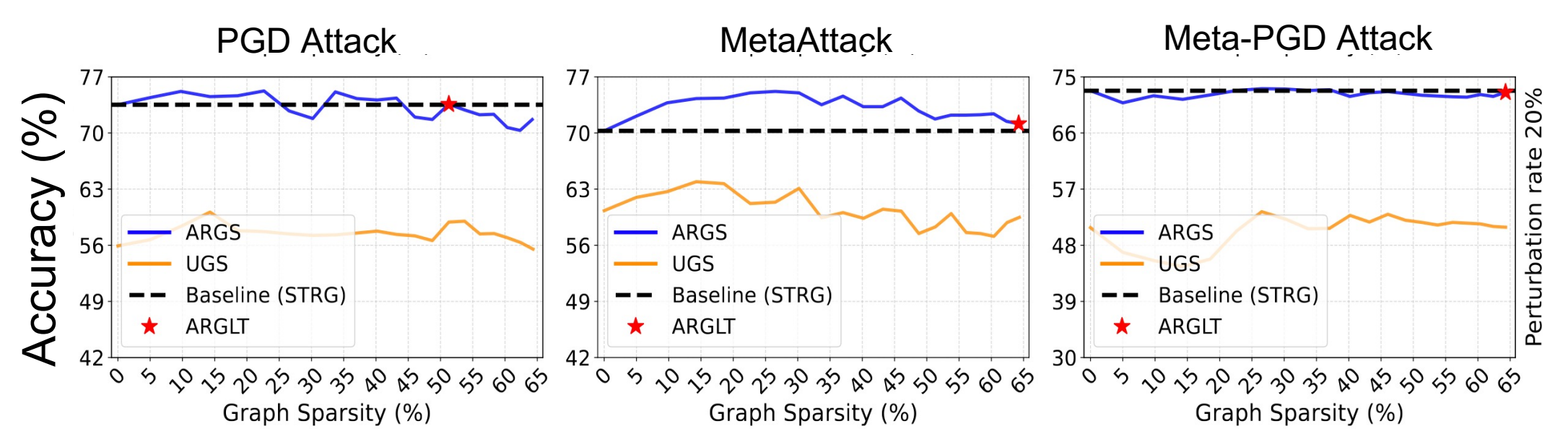
Annotations: CE loss for train nodes, Feature smoothness, CE loss for test nodes, l_1 regularizers for the masks.

Experimental Results

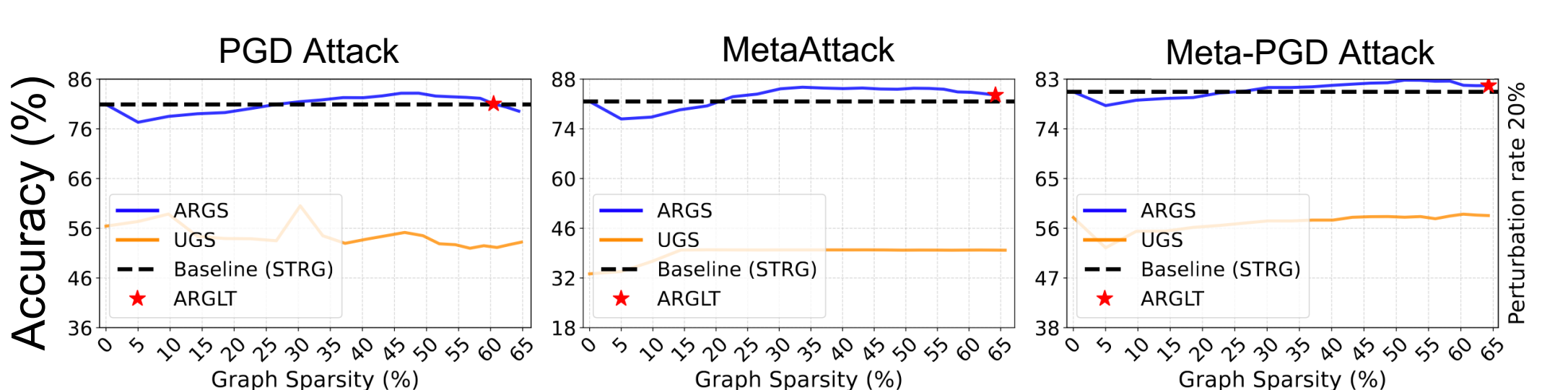
Cora (GCN)



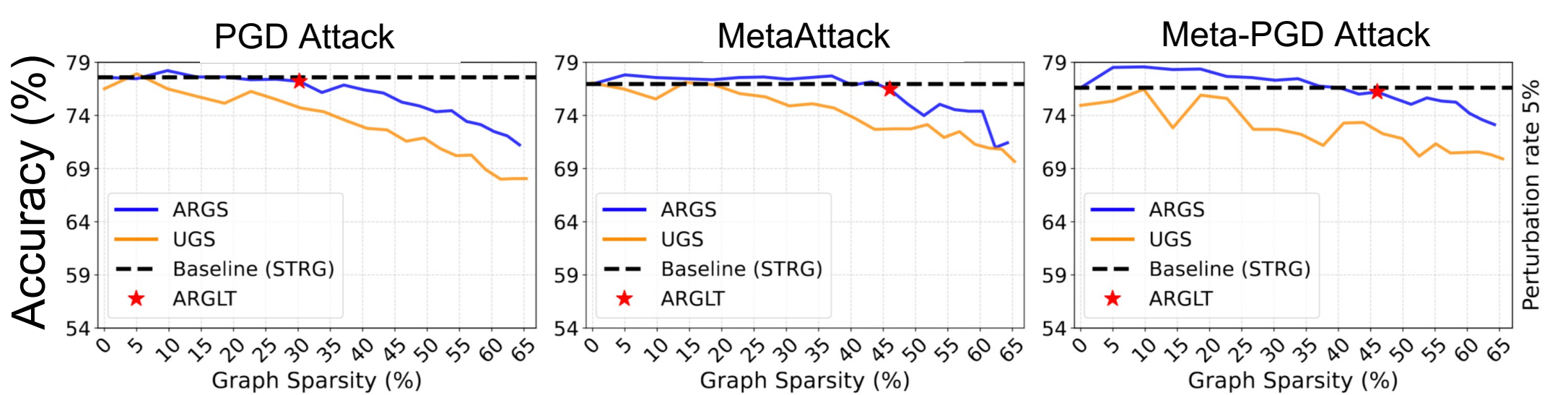
Citeseer (GCN)



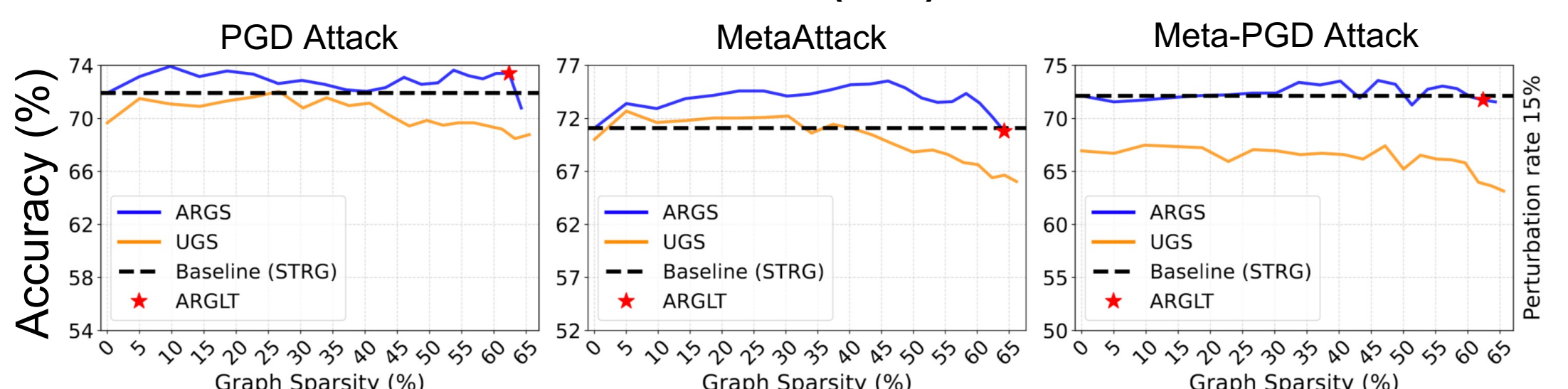
PubMed (GCN)



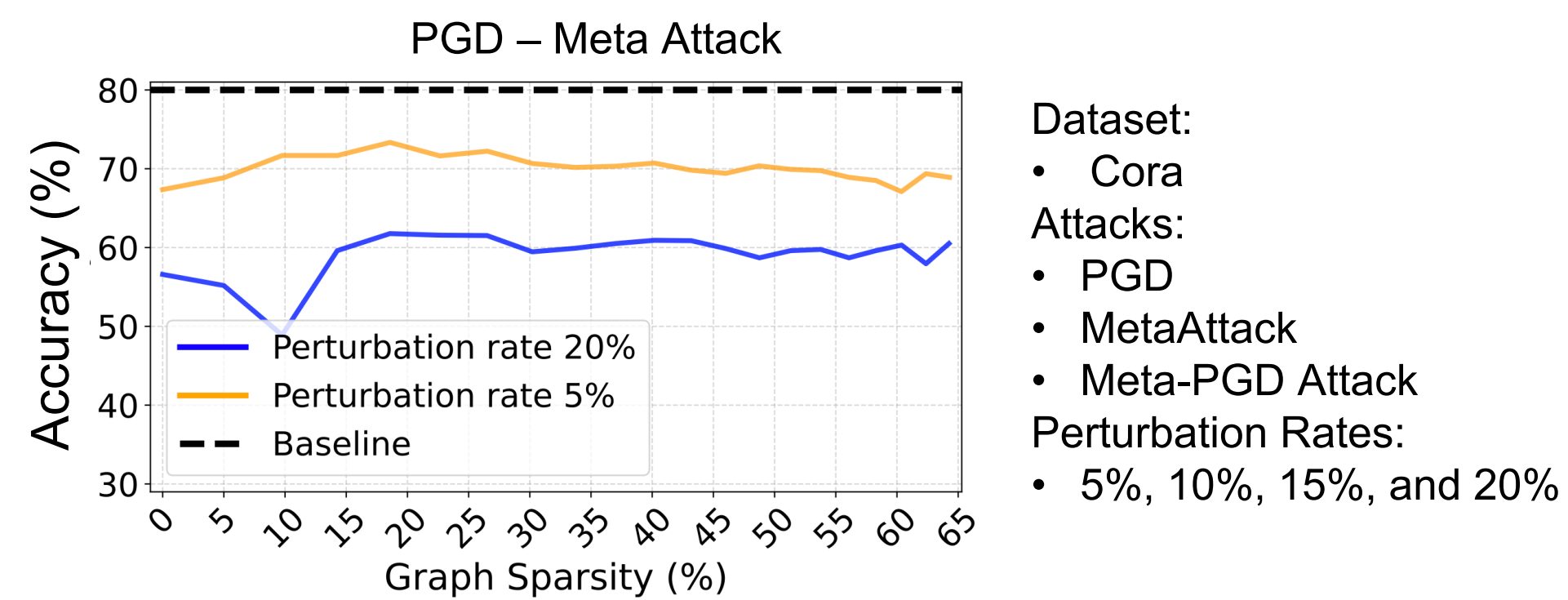
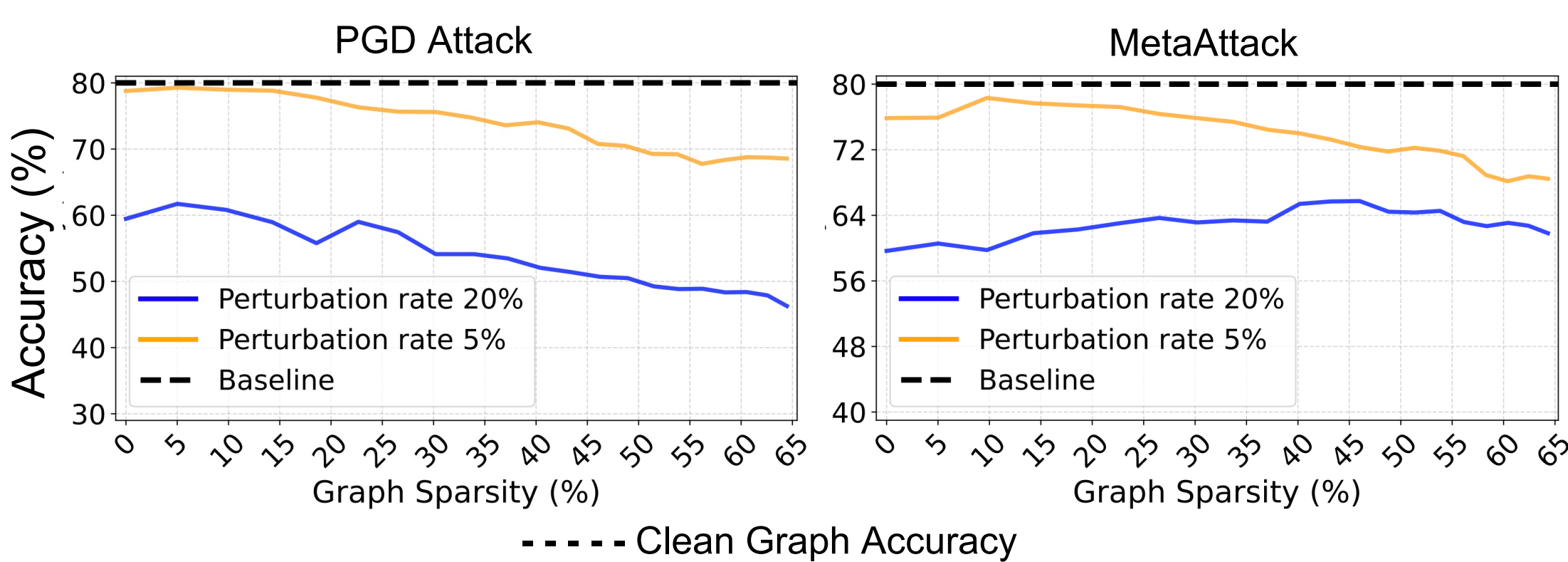
Cora (GIN)



Citeseer (GIN)



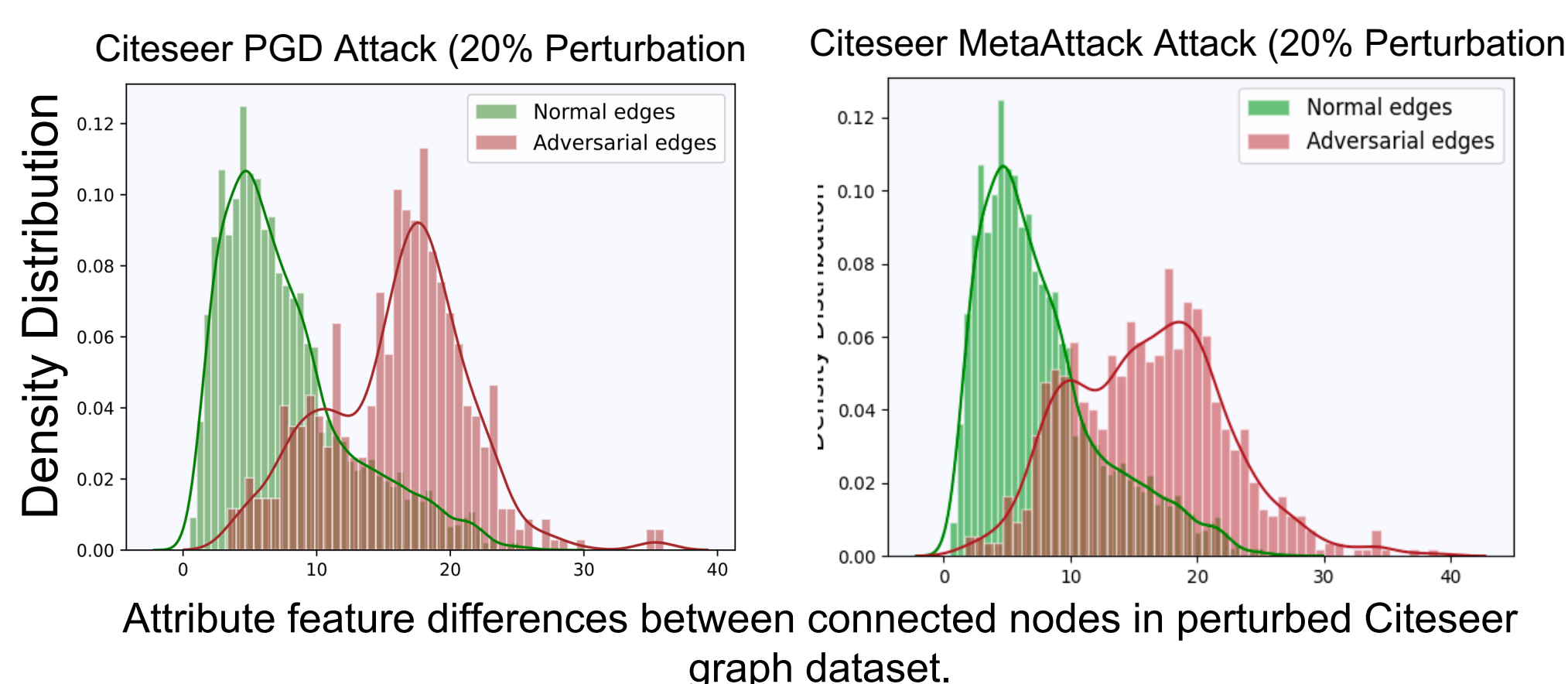
Analyzing the Robustness of GLTs to Adversarial Structure Perturbations



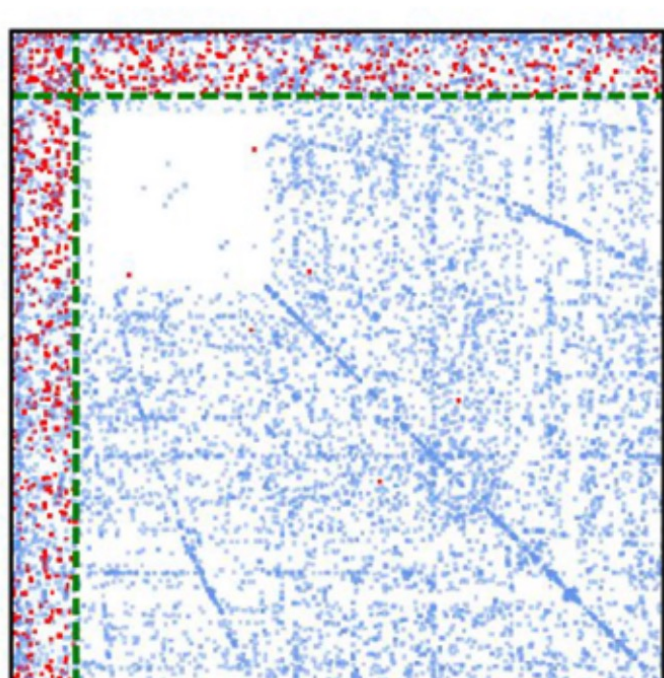
Dataset: Cora
Attacks: PGD, MetaAttack, Meta-PGD Attack
Perturbation Rates: 5%, 10%, 15%, and 20%

While UGS removes edges from the perturbed adjacency matrix, it **does not** effectively **remove** the **adversarial** edges.

Impact of Adversarial Attacks on Graphs



- Attacks tend to connect nodes with significant **attribute feature differences**.



- Edge perturbations are unevenly distributed on the graph
- Most modifications are around train nodes

Adjacency matrix of Cora attacked by MetaAttack (10%). Blue dots - Clean edges, Red dots - Adversarial edges, Green dotted line - Boundary of train and test nodes.

Summary

- Existing techniques for GLT identification are unable to generate adversarially robust GLTs
- The proposed technique ARGs can remove adversarial edges effectively and generate adversarially robust GLTs
- Future work includes extending ARGs to heterophilic graphs

