

## Research Motivation

- Deep neural networks, including GNNs, can suffer significant performance degradation due to prediction errors when real-world data changes, resulting in critical misclassifications.
- Current model editing techniques focus primarily on computer vision and NLP, with limited exploration of editable training for GNNs.
- **Key question:** Can we develop an effective method to edit GNNs that ensures corrections for erroneous predictions while maintaining model stability across unaffected nodes? If so, how?

## Why Gradient Rewiring?

Preliminary experiments show that direct fine-tuning of GNNs for model editing can lead to a significant increase in training loss, indicating performance degradation.

- **Statement:** There is a considerable gradient discrepancy between the target and training data, causing higher degradation for GNNs compared to MLPs.
- **Insight:** A method is needed to maintain training performance during model editing, motivating the development of a gradient rewiring approach.

## Gradient Rewiring Method

- **Problem Formulation:** Model editing aims to fix prediction errors at the target node while preserving performance on training nodes: (1) the training loss should not exceed its value prior to model editing (see Eq. (2)); and (2) the differences in model predictions after editing should remain within a predefined range (see Eq. (3)).

$$\min_{\theta'} \mathcal{L}_{tg}(f_{\theta'}(\mathbf{x}_{tg}), y_{tg}) \quad (1)$$

$$\text{s.t. } \mathcal{L}_{train}(f_{\theta'}, \mathcal{V}_{train}) \leq \mathcal{L}_{train}(f_{\theta_0}, \mathcal{V}_{train}) \quad (2)$$

$$\left\| \frac{1}{|\mathcal{V}_{train}|} \sum_{i \in \mathcal{V}_{train}} f_{\theta'}(\mathbf{x}_i) - f_{\theta_0}(\mathbf{x}_i) \right\|^2 \leq \delta', \quad (3)$$

- **Problem Solver:** (1) Approximation: Use Taylor expansion to estimate the influence of the model's parameters for both the target prediction and the training performance. (2) Transforming into Gradient Optimization (3) Solution via Dual Optimization: Solve the gradient adjustment problem more efficiently by converting it into a simpler form in the dual space.

### Algorithm 1 Gradient Rewiring Editable (GRE) Graph Neural Networks Training

- 1: **Input:** Target samples  $(\mathbf{x}_{tg}, y_{tg})$ , hyperparameter  $\lambda$ , well-trained GNN model  $f_{\theta}(\cdot)$ , and its corresponding gradient for the training subgraph.
- 2: **Output:** Updated GNN model  $f_{\theta'}(\cdot)$ .
- 3: **while**  $f_{\theta}(\mathbf{x}_{tg}) \neq y_{tg}$  **do**
- 4:     Compute the model gradient  $g_{tg}$  for the target loss  $\mathcal{L}_{tg}$ .
- 5:     Rewire the target loss gradient  $g_{tg}$  by reducing the projection component on  $g_{train}$ , then scale with  $(1 + \lambda)^{-1}$ :
- 6:      $g^* = (1 + \lambda)^{-1} (g_{tg} - v^* g_{train})$ .
- 7:     Replace  $g_{tg}$  with  $g^*$  and update the model parameters using the optimizer to obtain  $\theta'$ .
- 8: **end while**

## Experiment Results

- **Experimental Results in the Independent Editing Setting** (a) Our proposed GRE and GRE+ notably surpass both GD and ENN in terms of test drawdown; (b) Our proposed GRE and GRE+ are compatible with EGNN and further improve the performance.

Editor	Cora			A-computers			A-photo			Coauthor-CS			
	Acc $\uparrow$	DD $\downarrow$	SR $\uparrow$	Acc $\uparrow$	DD $\downarrow$	SR $\uparrow$	Acc $\uparrow$	DD $\downarrow$	SR $\uparrow$	Acc $\uparrow$	DD $\downarrow$	SR $\uparrow$	
MLP	GD	68.15 $\pm$ 0.33	3.85 $\pm$ 0.33	0.98	73.22 $\pm$ 0.48	6.78 $\pm$ 0.48	1.00	83.19 $\pm$ 0.91	6.81 $\pm$ 0.91	1.00	93.59 $\pm$ 0.05	0.41 $\pm$ 0.05	1.00
	ENN	37.16 $\pm$ 3.80	52.24 $\pm$ 4.76	1.00	15.51 $\pm$ 10.99	72.36 $\pm$ 10.87	1.00	16.71 $\pm$ 14.81	77.07 $\pm$ 15.20	1.00	4.94 $\pm$ 3.78	89.43 $\pm$ 3.34	1.00
	GRE	69.41 $\pm$ 0.44	2.59 $\pm$ 0.44	0.96	61.21 $\pm$ 1.26	18.79 $\pm$ 1.26	1.00	73.56 $\pm$ 1.41	16.44 $\pm$ 1.41	1.00	93.27 $\pm$ 0.09	0.73 $\pm$ 0.09	1.00
	GRE+	71.19 $\pm$ 0.28	0.61 $\pm$ 0.28	0.96	61.27 $\pm$ 1.15	18.73 $\pm$ 1.15	1.00	78.26 $\pm$ 1.15	11.74 $\pm$ 1.15	1.00	93.73 $\pm$ 0.07	0.27 $\pm$ 0.07	1.00
GCN	GD	84.37 $\pm$ 5.84	5.03 $\pm$ 6.40	1.00	44.78 $\pm$ 22.41	43.09 $\pm$ 22.32	1.00	28.70 $\pm$ 21.26	65.08 $\pm$ 20.13	1.00	91.07 $\pm$ 3.23	3.30 $\pm$ 2.22	1.00
	ENN	37.16 $\pm$ 3.80	52.24 $\pm$ 4.76	1.00	15.51 $\pm$ 10.99	72.36 $\pm$ 10.87	1.00	16.71 $\pm$ 14.81	77.07 $\pm$ 15.20	1.00	4.94 $\pm$ 3.78	89.43 $\pm$ 3.34	1.00
	GRE	84.98 $\pm$ 0.47	4.02 $\pm$ 0.47	0.96	46.28 $\pm$ 3.47	51.72 $\pm$ 3.47	0.98	35.88 $\pm$ 2.26	58.12 $\pm$ 2.26	0.99	89.46 $\pm$ 0.29	4.54 $\pm$ 0.29	1.00
	GRE+	88.84 $\pm$ 0.35	0.56 $\pm$ 0.35	0.98	47.75 $\pm$ 0.45	40.25 $\pm$ 0.45	1.00	50.13 $\pm$ 1.36	43.87 $\pm$ 1.36	1.00	91.99 $\pm$ 0.30	2.01 $\pm$ 0.30	1.00
Graph-SAGE	GD	82.06 $\pm$ 4.33	4.54 $\pm$ 5.32	1.00	21.68 $\pm$ 20.98	61.15 $\pm$ 20.33	1.00	38.98 $\pm$ 30.24	55.32 $\pm$ 29.35	1.00	90.15 $\pm$ 5.58	5.01 $\pm$ 5.32	1.00
	ENN	33.16 $\pm$ 1.45	53.44 $\pm$ 2.23	1.00	16.89 $\pm$ 16.98	65.94 $\pm$ 16.75	1.00	15.06 $\pm$ 11.92	79.24 $\pm$ 11.25	1.00	13.71 $\pm$ 2.73	81.45 $\pm$ 2.11	1.00
	GRE	83.64 $\pm$ 0.20	3.36 $\pm$ 0.20	1.00	20.11 $\pm$ 2.30	62.89 $\pm$ 2.30	0.96	41.96 $\pm$ 1.57	52.04 $\pm$ 1.57	0.98	91.07 $\pm$ 0.44	3.93 $\pm$ 0.44	1.00
	GRE+	86.59 $\pm$ 0.07	0.41 $\pm$ 0.07	1.00	22.23 $\pm$ 1.60	60.77 $\pm$ 1.60	0.97	44.05 $\pm$ 0.83	50.32 $\pm$ 0.83	1.00	91.75 $\pm$ 0.43	3.25 $\pm$ 0.43	1.00
EGNN-GCN	GD	87.58 $\pm$ 0.31	1.42 $\pm$ 0.31	1.00	87.27 $\pm$ 0.14	0.73 $\pm$ 0.14	0.78	93.24 $\pm$ 0.59	0.76 $\pm$ 0.59	0.77	93.99 $\pm$ 0.02	0.01 $\pm$ 0.02	0.91
	GRE	87.47 $\pm$ 0.41	1.53 $\pm$ 0.41	1.00	83.38 $\pm$ 1.20	4.62 $\pm$ 1.20	0.87	88.01 $\pm$ 1.20	5.99 $\pm$ 1.20	0.86	93.92 $\pm$ 0.07	0.08 $\pm$ 0.07	0.94
	GRE+	88.99 $\pm$ 0.21	0.05 $\pm$ 0.21	1.00	88.10 $\pm$ 1.21	0.51 $\pm$ 1.21	1.00	94.22 $\pm$ 0.98	-0.21 $\pm$ 0.98	1.00	94.32 $\pm$ 0.06	-0.32 $\pm$ 0.06	1.00
	GD	85.05 $\pm$ 0.11	0.95 $\pm$ 0.11	1.00	85.93 $\pm$ 0.08	0.07 $\pm$ 0.08	0.90	93.87 $\pm$ 0.20	0.13 $\pm$ 0.20	0.81	95.0 $\pm$ 0.01	0.00 $\pm$ 0.01	0.99
EGNN-SAGE	GRE	84.79 $\pm$ 0.19	1.21 $\pm$ 0.19	1.00	81.94 $\pm$ 1.71	4.06 $\pm$ 1.71	0.96	88.55 $\pm$ 1.19	5.45 $\pm$ 1.19	0.95	94.85 $\pm$ 0.05	0.15 $\pm$ 0.05	1.00
	GRE+	86.24 $\pm$ 1.43	-0.24 $\pm$ 1.43	1.00	85.97 $\pm$ 0.83	-0.16 $\pm$ 0.83	1.00	94.07 $\pm$ 0.03	-0.07 $\pm$ 0.03	0.98	95.07 $\pm$ 0.03	-0.07 $\pm$ 0.03	1.00

- **Experimental Results in the Sequential Editing Setting.** (a) The proposed GRE and GRE+ consistently outperform GD in the sequential setting. (b) The improvement of GRE+ over GRE is quite limited in the sequential setting.

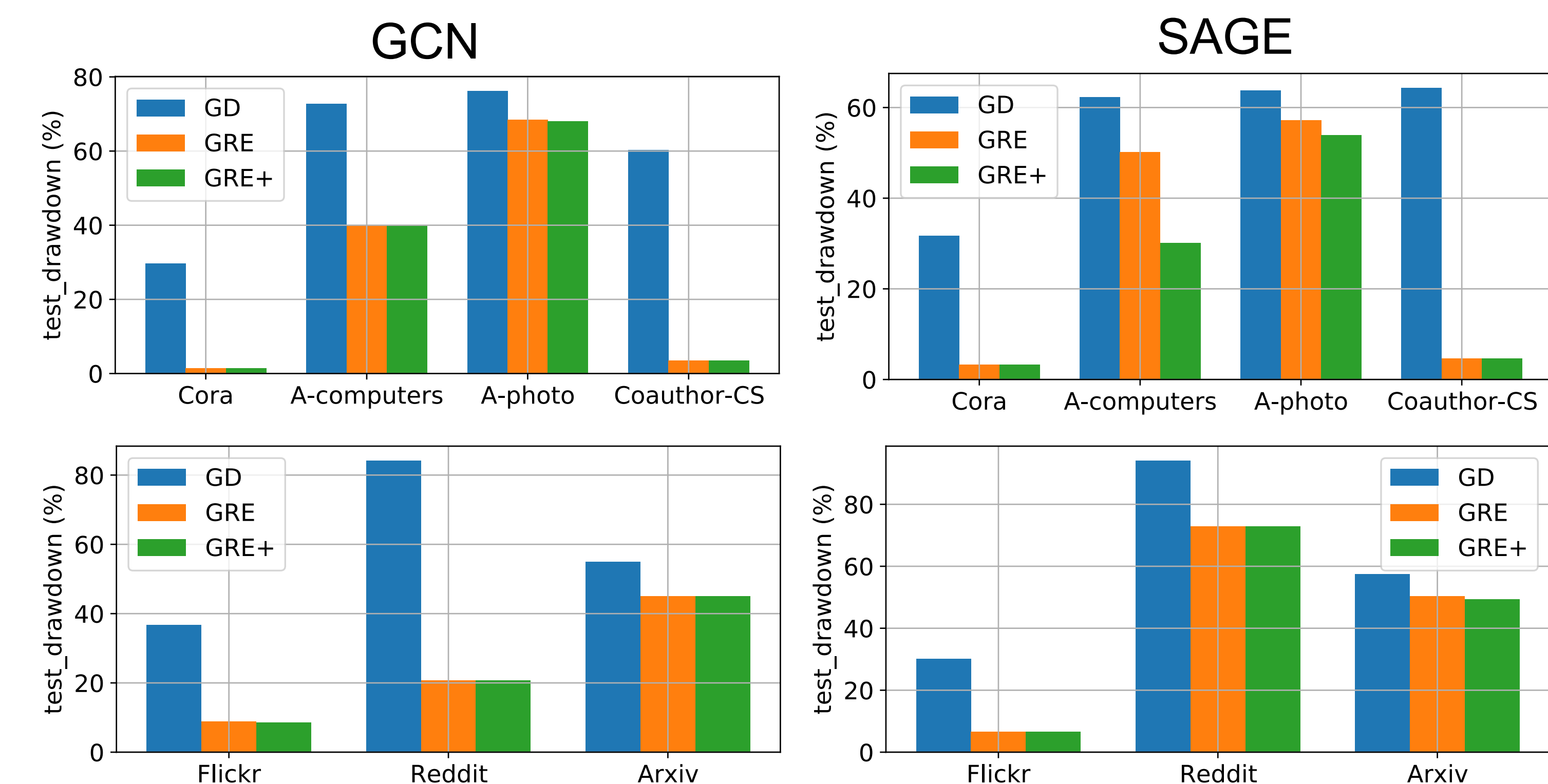


Figure: The test accuracy drawdown in sequential editing setting for GCN and GraphSAGE on various datasets. The units for y-axis are percentages (%).

**Acknowledgement:** The work is in part supported by NSF grants NSF IIS-2310260, IIS-2224843, IIS-2450662, IIS-2431515 and IIS-2239257.