# On the Feasibility of Fidelity− for Graph Pruning

Yong-Min Shin, Won-Yong Shin
Yonsei Univeristy
Seoul, South Korea

Paper    Homepage

## Introduction

### XAI for graph neural networks (GNNs)

- Many XAI methods for GNNs **highlight local edges** that are **highly relevant** to the output.

- **Edge attribution**: How much can we attribute the model's output to each input edge?

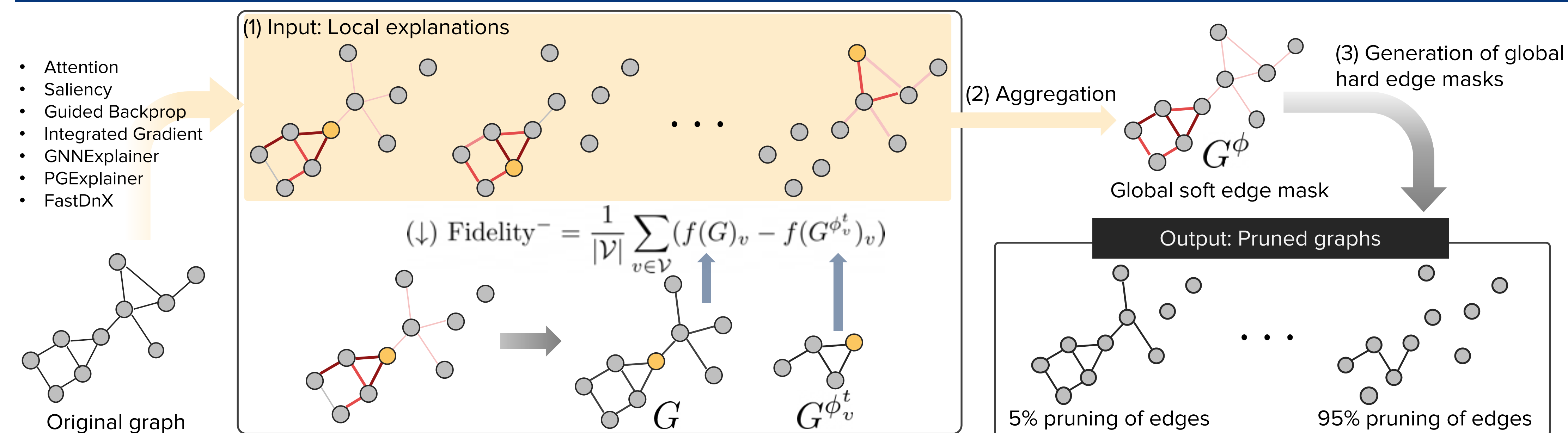### Quantitative measurement: Fidelity-

- Output **difference** between two instances when the **original input graph** is used and when the **unimportant edges are removed**.

- If the explanation is valid, then underline{edges deemed less important} should have underline{less impact to the model's output} after **removal** from the input.

## Research Question

**Can we improve the efficiency of the GNN model by performing graph pruning by performing graph pruning based on GNN explanations?**
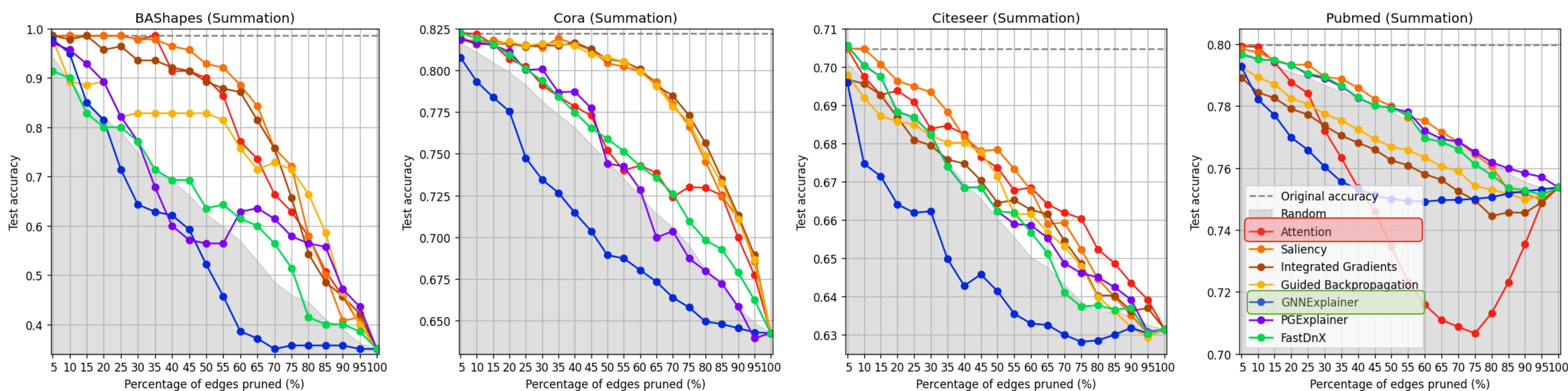
- **Intuition**: If an edge is frequently removed in Fidelity-, it may simply be removed from the original graph.

## Overview: FiP (Fidelity-inspired Pruning)



- Attention
- Saliency
- Guided Backprop
- Integrated Gradient
- GNNExplainer
- PGExplainer
- FastDnX

Original graph

(1) Input: Local explanations

(2) Aggregation

(3) Generation of global hard edge masks

Global soft edge mask    $G^\phi$

$$(\downarrow)\ \text{Fidelity}^- = \frac{1}{|\mathcal{V}|}\sum_{v\in\mathcal{V}}(f(G)_v - f(G^{\phi_v^t})_v)$$

$G$    $G^{\phi_v^t}$

Output: Pruned graphs

5% pruning of edges    95% pruning of edges

## Observations and Discussion

### Obs.1: Can Explanation be used for Graph Pruning?



### Obs.2: Does Graph Pruning and Fidelity- Scores Translate?

| Method | BAShapes | Cora | Citeseer | Pubmed |
|--------|----------|------|----------|--------|
| Att | $4.06 \times 10^{-2}$ | $3.67 \times 10^{-2}$ | $2.23 \times 10^{-2}$ | $2.46 \times 10^{0}$ |
| SA | $3.54 \times 10^{-7}$ | $2.21 \times 10^{-7}$ | $\mathbf{8.90 \times 10^{-8}}$ | $2.46 \times 10^{0}$ |
| IG | $6.25 \times 10^{0}$ | $1.26 \times 10^{0}$ | $5.68 \times 10^{-1}$ | $\mathbf{2.25 \times 10^{0}}$ |
| GB | $3.77 \times 10^{0}$ | $1.42 \times 10^{0}$ | $7.04 \times 10^{-1}$ | $2.40 \times 10^{0}$ |
| GNNEx | $\mathbf{3.44 \times 10^{-7}}$ | $\mathbf{2.14 \times 10^{-7}}$ | $3.52 \times 10^{-1}$ | $2.46 \times 10^{0}$ |
| PGEx | $3.83 \times 10^{-7}$ | $2.04 \times 10^{-2}$ | $7.11 \times 10^{-3}$ | $2.46 \times 10^{0}$ |
| FDnX | $1.41 \times 10^{-1}$ | $1.77 \times 10^{-2}$ | $7.05 \times 10^{-3}$ | $2.46 \times 10^{0}$ |

### Discussions on Results

1) **Explanations can be used for graph pruning.**
2) Avoid GNN-tailored methods.
3) Graph pruning naturally improve efficiency for GNN.
4) Fidelity does not translate to graph pruning.
5) The problem likely resides in the **practical problem of aggregating local explanation**:
   - Scale of attribution score across nodes
   - Num. of edges for each explanation
6) **Inherent limitation** of graph pruning: Local explanation is unique for each node, **creating a single (pruned) graph is always a lossy compression of information**.