

# When Witnesses Defend: A Witness Graph Topological Layer for Adversarial Graph Learning.

**Naheed Anjum Arafat**, Debabrota Basu, Yulia Gel, Yuzhou Chen

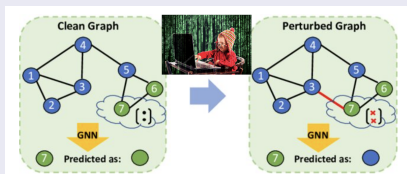
Post-doc  
Nanyang Technological University, Singapore (2021-2024)  
PhD  
National University of Singapore (2020)

February 11, 2025

# Adversarial attack: a Modern Challenge to GNNs

## Adversarial attack on Graph learning algorithms.

Attacker misleads a learning algorithm (e.g. GNN) into making incorrect predictions or classifications by deliberately perturbing a small number of edges (e.g. remove/add edges) or node features.



Adversarial perturbation (around target node 7) causes misclassification.

## Contributions

- 1 We introduced a novel topological adversarial defense, namely, the *Witness Graph Topological Layer (WGTL)*.
- 2 WGTL integrates local and global higher-order graph characteristics and controls their potential defense role via a topological regularizer.

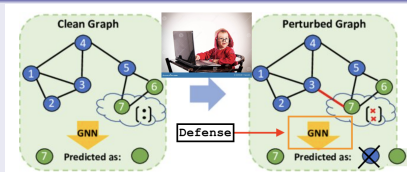


# Adversarial attack: a Modern Challenge to GNNs

## Adversarial attack on Graph learning algorithms.

Attacker misleads a learning algorithm (e.g. GNN) into making incorrect predictions or classifications by deliberately perturbing a small number of edges (e.g. remove/add edges) or node features.

## Problem

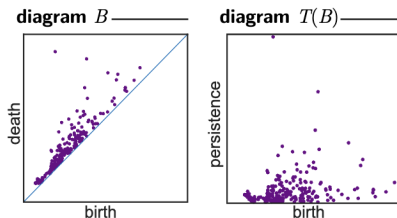


Design a defense algorithm that mitigates the effect of adversarial attack

## Contributions

- 1 We introduced a novel topological adversarial defense, namely, the *Witness Graph Topological Layer (WGTL)*.
- 2 WGTL integrates local and global higher-order graph characteristics and controls their potential defense role via a topological regularizer.

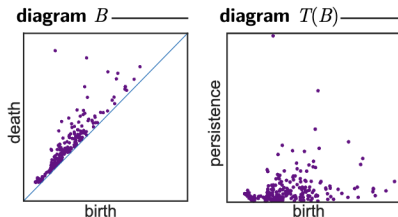
# Topological Features



A persistence diagram is transformed using function  $T : (x, y) \rightarrow (0, y - x)$ .

## Why Topological features?

**Stability theorem:** Small change in the data (graph) only result in small changes in the persistence diagram.



A persistence diagram is transformed using function  $T : (x, y) \rightarrow (0, y - x)$ .

## Why Topological features?

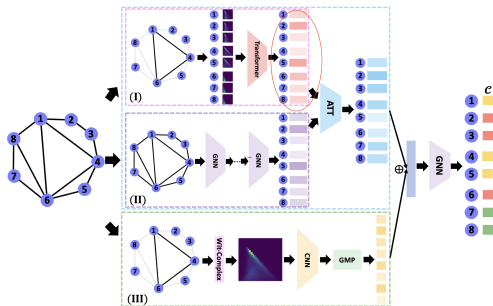
**Stability theorem:** Small change in the data (graph) only result in small changes in the persistence diagram.

- 1 This is the first work that shows that topological features can make GNNs robust against adversarial perturbations.
- 2 Effective against a wide variety of attacks, for instance,
  - Targetted poisoning attack (Graybox, modify the neighbors of a target node and their features)
  - Global poisoning attack (Graybox, instead of targetting specific neighborhood modify whichever edges minimizes the model accuracy)
  - Adaptive attacks (White-box, the model architecture, parameters and defense mechanisms are known to the attacker)
  - Node feature attack
- 3 WGTl improves existing defenses such as Pro-GNN, GNNGuard, and SimP-GCN respectively by 5%, 15%, and 5%.

- 1 This is the first work that shows that topological features can make GNNs robust against adversarial perturbations.
- 2 Effective against a wide variety of attacks, for instance,
  - Targetted poisoning attack (Graybox, modify the neighbors of a target node and their features)
  - Global poisoning attack (Graybox, instead of targetting specific neighborhood modify whichever edges minimizes the model accuracy)
  - Adaptive attacks (White-box, the model architecture, parameters and defense mechanisms are known to the attacker)
  - Node feature attack
- 3 WGTL improves existing defenses such as Pro-GNN, GNNGuard, and SimP-GCN respectively by 5%, 15%, and 5%.

- 1 This is the first work that shows that topological features can make GNNs robust against adversarial perturbations.
- 2 Effective against a wide variety of attacks, for instance,
  - Targetted poisoning attack (Graybox, modify the neighbors of a target node and their features)
  - Global poisoning attack (Graybox, instead of targetting specific neighborhood modify whichever edges minimizes the model accuracy)
  - Adaptive attacks (White-box, the model architecture, parameters and defense mechanisms are known to the attacker)
  - Node feature attack
- 3 WGTL improves existing defenses such as Pro-GNN, GNNGuard, and SimP-GCN respectively by 5%, 15%, and 5%.

# WGTL: Topological Encodings



Architecture of Witness Graph Topological Layer.

- Ⓘ **Local Topology Encoding:** Encodes local topological features of every node. ( $Z_{T_L}$ )
- Ⓙ **Node Representation Learning.** Learns node representations using any backbone GNN. ( $Z_G$ )
- Ⓜ **Global Topology Encoding.** Encodes topological feature of the entire graph. ( $Z_{T_G}$ )
- Ⓝ **Aggregated Topological Encoding.** Encodes local and global topological priors. ( $Z_{WGTL}$ )

$$z = [Z_{T_L}, Z_G]$$

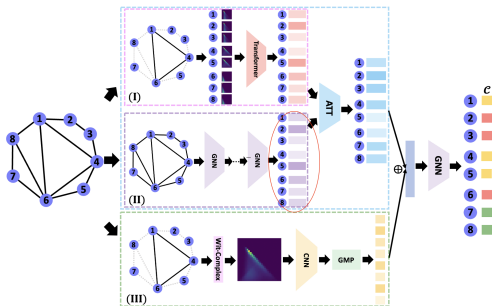
$$\text{Attention coefficients, } \alpha_i = \text{Softmax}(W_2 \cdot \tanh(W_1 z_i + b_1))$$

$$\text{Additive attention, } Z_{AGG} = \alpha_1 \times Z_{T_L} + \alpha_2 \times Z_G$$

$$Z_{WGTL} = Z_{AGG} Z_{T_G}$$



# WGTL: Topological Encodings



Architecture of Witness Graph Topological Layer.

- I** Local Topology Encoding: Encodes local topological features of every node. ( $Z_{T_L}$ )
- II** Node Representation Learning. Learns node representations using any backbone GNN. ( $Z_G$ )
- III** Global Topology Encoding. Encodes topological feature of the entire graph. ( $Z_{T_G}$ )
- IV** Aggregated Topological Encoding. Encodes local and global topological priors. ( $Z_{WGTL}$ )

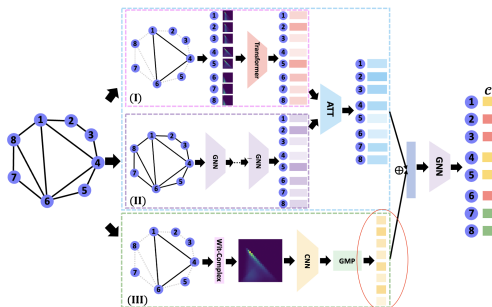
$$z = [Z_{T_L}, Z_G]$$

Attention coefficients,  $\alpha_i = \text{Softmax}(W_2 \cdot \tanh(W_1 z_i + b_1))$

Additive attention,  $Z_{AGG} = \alpha_1 \times Z_{T_L} + \alpha_2 \times Z_G$

$$Z_{WGTL} = Z_{AGG} Z_{T_G}$$

# WGTL: Topological Encodings



Architecture of Witness Graph Topological Layer.

- I** Local Topology Encoding: Encodes local topological features of every node. ( $Z_{T_L}$ )
- II** Node Representation Learning. Learns node representations using any backbone GNN. ( $Z_G$ )
- III** Global Topology Encoding. Encodes topological feature of the entire graph. ( $Z_{T_G}$ )
- IV** Aggregated Topological Encoding. Encodes local and global topological priors. ( $Z_{WGTL}$ )

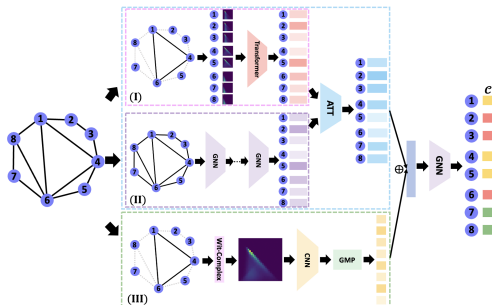
$$z = [Z_{T_L}, Z_G]$$

$$\text{Attention coefficients, } \alpha_i = \text{Softmax}(W_2 \cdot \tanh(W_1 z_i + b_1))$$

$$\text{Additive attention, } Z_{AGG} = \alpha_1 \times Z_{T_L} + \alpha_2 \times Z_G$$

$$Z_{WGTL} = Z_{AGG} Z_{T_G}$$

# WGTL: Topological Encodings



Architecture of Witness Graph Topological Layer.

- ❶ Local Topology Encoding: Encodes local topological features of every node. ( $Z_{T_L}$ )
- ❷ Node Representation Learning. Learns node representations using any backbone GNN. ( $Z_G$ )
- ❸ Global Topology Encoding. Encodes topological feature of the entire graph. ( $Z_{T_G}$ )
- ❹ Aggregated Topological Encoding. Encodes local and global topological priors. ( $Z_{WGTL}$ )

$$z = [Z_{T_L}, Z_G]$$

Attention coefficients,  $\alpha_i = \text{Softmax}(W_2 \cdot \tanh(W_1 z_i + b_1))$

Additive attention,  $Z_{AGG} = \alpha_1 \times Z_{T_L} + \alpha_2 \times Z_G$

$$Z_{WGTL} = Z_{AGG} Z_{T_G}$$

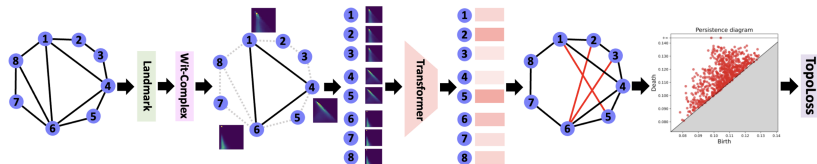


Illustration of Witness Complex-based topological regularizer  $L_{Topo}$ .

$$L_{topo}(\mathbb{T}(\mathcal{G})) \triangleq \sum_{i=1}^m (d_i - b_i)^2 \left( \frac{d_i + b_i}{2} \right)^2, \quad (1)$$

- A localized attack (perturbing certain nodes or edges) appears as topological noise in the final persistent diagram, and exhibit lower persistence.
- And minimising  $L_{topo}$  forces the Transformer to learn local topology encodings ( $Z_{T_L}$ ) which produces PD with small persistence, i.e.,  $(d_i - b_i)$ .

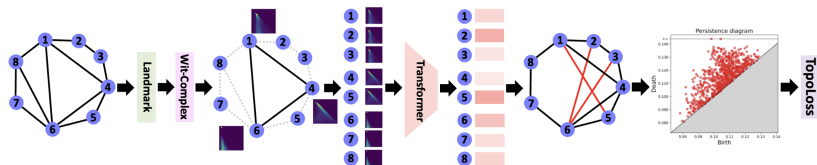


Illustration of Witness Complex-based topological regularizer  $L_{Topo}$ .

$$L_{topo}(T(\mathcal{G})) \triangleq \sum_{i=1}^m (d_i - b_i)^2 \left( \frac{d_i + b_i}{2} \right)^2, \quad (1)$$

- A localized attack (perturbing certain nodes or edges) appears as topological noise in the final persistent diagram, and exhibit lower persistence.
- And minimising  $L_{topo}$  forces the Transformer to learn local topology encodings ( $Z_{T_L}$ ) which produces PD with small persistence, i.e.,  $(d_i - b_i)$ .

Table 1: Comparison of performances (avg. accuracy $\pm$ std.) with existing defenses under metttack.

Dataset	Models	Perturbation Rate		
		0%	5%	10%
Cora-ML	Pro-GNN	82.98 $\pm$ 0.23	80.14 $\pm$ 1.34	71.59 $\pm$ 1.33
	Pro-GNN+WGTL	<b>83.85<math>\pm</math>0.38</b>	<b>81.90<math>\pm</math>0.73</b>	<b>72.51<math>\pm</math>0.76</b>
	GCN+GNNGuard	83.21 $\pm$ 0.34	76.57 $\pm$ 0.50	69.13 $\pm$ 0.77
	GCN+GNNGuard+WGTL	<b>*84.78<math>\pm</math>0.43</b>	<b>*83.23<math>\pm</math>0.82</b>	<b>*79.96<math>\pm</math>0.49</b>
	SimP-GCN	79.52 $\pm$ 1.81	74.75 $\pm$ 1.40	70.87 $\pm$ 1.70
	SimP-GCN+WGTL	<b>81.49<math>\pm</math>0.52</b>	<b>76.65<math>\pm</math>0.65</b>	<b>72.88<math>\pm</math>0.83</b>
Citeseer	ProGNN	72.34 $\pm$ 0.99	68.96 $\pm$ 0.67	67.36 $\pm$ 1.12
	ProGNN+WGTL	<b>72.83<math>\pm</math>0.94</b>	<b>71.85<math>\pm</math>0.74</b>	<b>70.70<math>\pm</math>0.57</b>
	GCN+GNNGuard	71.82 $\pm$ 0.43	70.79 $\pm$ 0.22	66.86 $\pm$ 0.54
	GCN+GNNGuard+WGTL	<b>73.37<math>\pm</math>0.63</b>	<b>72.57<math>\pm</math>0.17</b>	<b>66.93<math>\pm</math>0.21</b>
	SimP-GCN	73.73 $\pm$ 1.54	73.06 $\pm$ 2.09	72.51 $\pm$ 1.25
	SimP-GCN+WGTL	<b>*74.32<math>\pm</math>0.19</b>	<b>*74.05<math>\pm</math>0.71</b>	<b>*73.09<math>\pm</math>0.50</b>
Pubmed	Pro-GNN	87.33 $\pm$ 0.18	87.25 $\pm$ 0.09	87.20 $\pm$ 0.12
	Pro-GNN + WGTL (ours)	<b>87.90<math>\pm</math>0.30</b>	<b>*87.77<math>\pm</math>0.08</b>	<b>*87.67<math>\pm</math>0.22</b>
	GCN+GNNGuard	83.63 $\pm$ 0.08	79.02 $\pm$ 0.14	76.58 $\pm$ 0.16
	GCN+GNNGuard+WGTL	OOM	OOM	OOM
	SimP-GCN	*88.11 $\pm$ 0.10	86.98 $\pm$ 0.19	86.30 $\pm$ 0.28
	SimP-GCN+WGTL	OOM	OOM	OOM
Polblogs	GCN+GNNGuard	95.03 $\pm$ 0.25	73.25 $\pm$ 0.16	72.76 $\pm$ 0.75
	GCN+GNNGuard+WGTL	<b>*96.22<math>\pm</math>0.25</b>	<b>*73.62<math>\pm</math>0.22</b>	<b>*73.72<math>\pm</math>1.00</b>
	SimP-GCN	89.78 $\pm$ 6.47	65.75 $\pm$ 5.03	61.53 $\pm$ 6.41
	SimP-GCN+WGTL	<b>94.56<math>\pm</math>0.24</b>	<b>69.78<math>\pm</math>4.10</b>	<b>69.55<math>\pm</math>4.42</b>

Table 2: Efficiency of WGTl. All the times are in seconds.

Datasets/ (# Landmarks)	Landmark selection time	Local feat. comput. time	Global feat. comput. time
Cora-ML/124	$0.01 \pm 0.01$	$0.12 \pm 0.03$	$5.11 \pm 0.13$
Citeseer/105	$0.01 \pm 0.01$	$0.16 \pm 0.02$	$5.23 \pm 1.22$
Polblogs/61	$0.01 \pm 0.00$	$0.07 \pm 0.01$	$4.64 \pm 0.2$
Snap-patents/91	$0.03 \pm 0.02$	$0.64 \pm 0.00$	$7.54 \pm 1.15$
Pubmed/394	$0.07 \pm 0.01$	$0.51 \pm 0.03$	$27.83 \pm 0.47$
OGBN-arXiv/84	$1.02 \pm 0.00$	$12.79 \pm 0.31$	$83.04 \pm 2.19$