

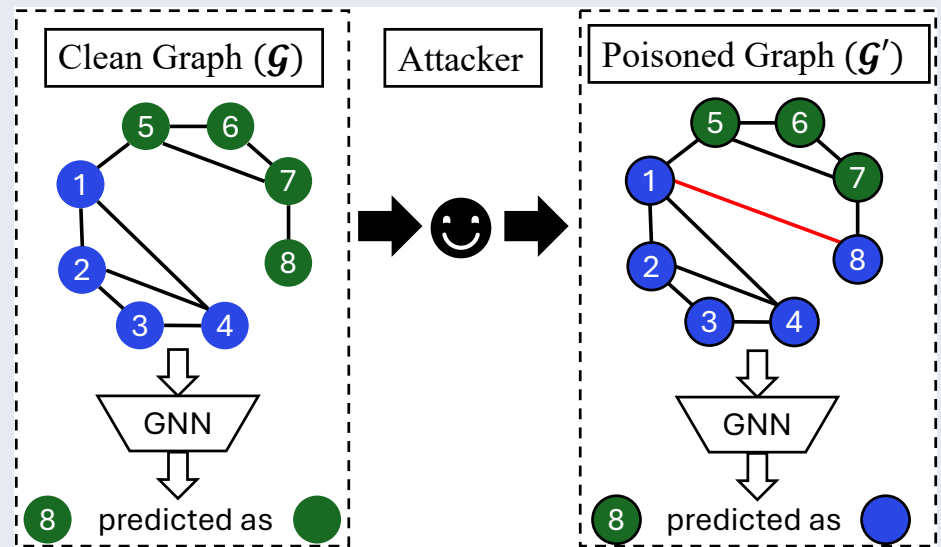
# When Witnesses Defend: A Witness Graph Topological Layer for Adversarial Graph Learning

Naheed Anjum Arafat\*, Debabrota Basu<sup>+</sup>, Yulia R. Gel<sup>#</sup>, Yuzhou Chen<sup>\$</sup>

\*Nanyang Technological Univ. (Singapore), <sup>+</sup>Univ. Lille, Inria (France), <sup>#</sup>Virginia Tech (USA), <sup>\$</sup> UC Riverside (USA)

## Adversarial Attack on GNNs

The attacker **misleads** GNNs into making incorrect predictions by deliberately perturbing a **small number of edges** (e.g. remove/add edges) or node features.



**Goal:** Design a robust representation  $R$  such that

$$\|R(\mathcal{G}) - R(\mathcal{G}')\|_1 = \mathcal{O}(\delta)$$

when  $\|\mathcal{G} - \mathcal{G}'\|_1 = \delta$

## Our Contributions

- This is the first work that shows that **topological features can make GNNs robust** against adversarial attacks.
- Our approach integrates **local** and **global** higher-order graph characteristics (**using Persistence homology**) and controls their potential defense role via a **topological regularizer**.
- Effective against a wide **variety of attacks**:
  - Targetted (perturbs neighbors of a set of target nodes)
  - Global (perturbs whichever edges minimizes the attackers loss)
  - Adaptive (White-box, the model architecture, parameters and the defense mechanisms are known to the attacker)
  - Node feature attack

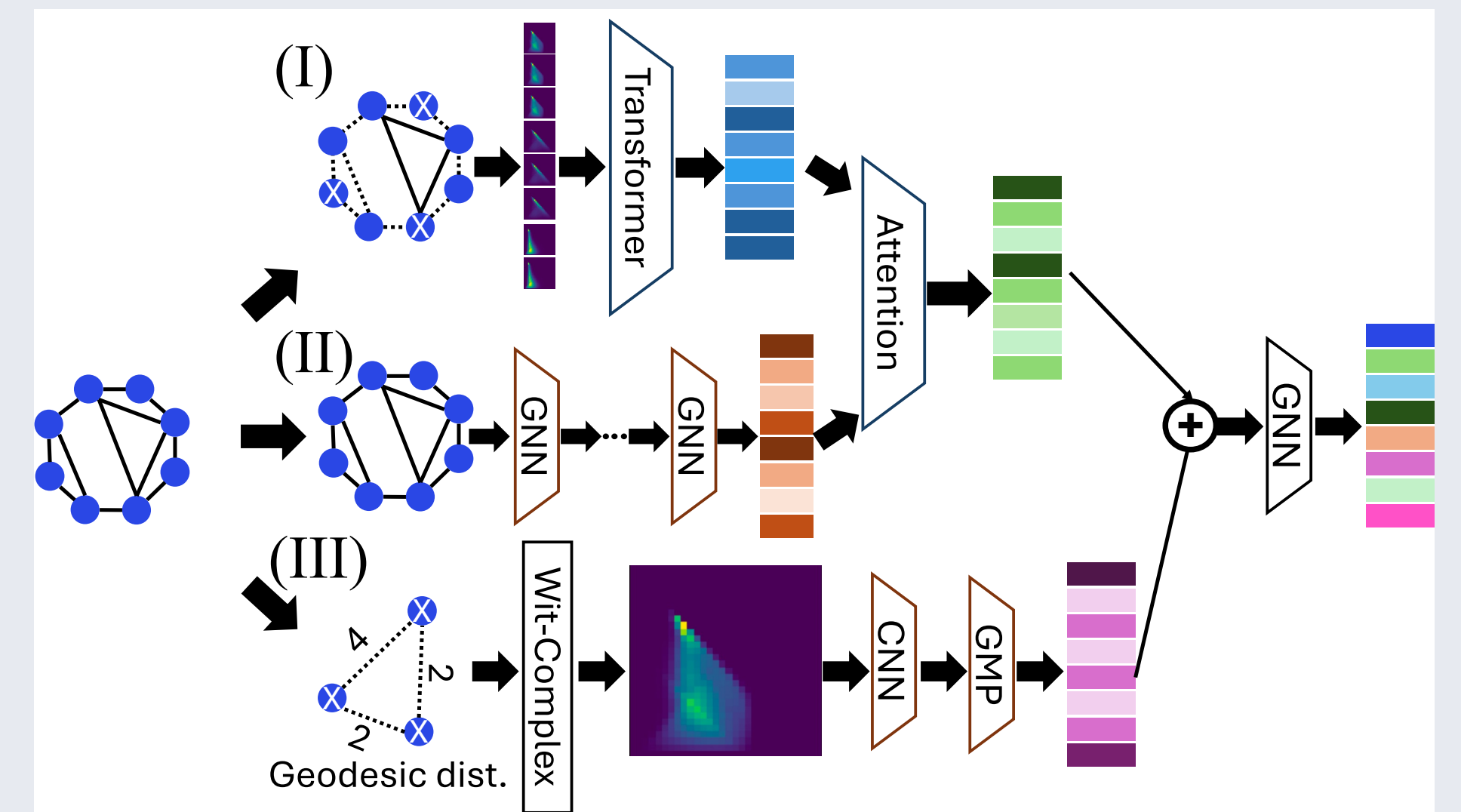
## Performance (Mettack)

Dataset	Models	Perturbation Rate	
		0%	10%
Cora-ML	GCN	82.87±0.83	70.39±1.28
	GCN+WGTLp(ours)	<b>83.83±0.55</b>	<b>71.31±0.85</b>
	GAT	84.25±0.67	72.63±1.56
	GAT+WGTLp(ours)	<b>*86.07±2.10</b>	<b>*73.74±1.92</b>
	GraphSAGE	81.00±0.27	70.92±1.18
	GraphSAGE+WGTLp(ours)	<b>83.63±0.35</b>	<b>73.57±0.73</b>
	ProGNN	82.98±0.23	71.59±1.33
	ProGNN+WGTLp(ours)	<b>83.85±0.38</b>	<b>72.71±1.26</b>
	GCN+GNNGuard	83.21±0.34	69.13±0.77
	GCN+GNNGuard+WGTLp(ours)	<b>84.78±0.43</b>	<b>70.15±0.89</b>
Polblogs	GCN	94.40±1.47	69.16±1.86
	GCN+WGTLp(ours)	<b>95.95±0.15</b>	<b>74.52±0.28</b>
	GAT	95.28±0.51	73.11±1.20
	GAT+WGTLp(ours)	<b>95.87±0.26</b>	<b>74.21±0.74</b>
	GraphSAGE	94.54±0.27	74.66±0.85
	GraphSAGE+WGTLp(ours)	<b>95.58±0.50</b>	<b>*74.93±0.81</b>
	GCN+GNNGuard	95.03±0.25	72.76±0.75
	GCN+GNNGuard+WGTLp(ours)	<b>*96.22±0.25</b>	<b>73.72±1.00</b>

## Efficiency comparison

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

## Witness Graph Topological Layer (WGTL)



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

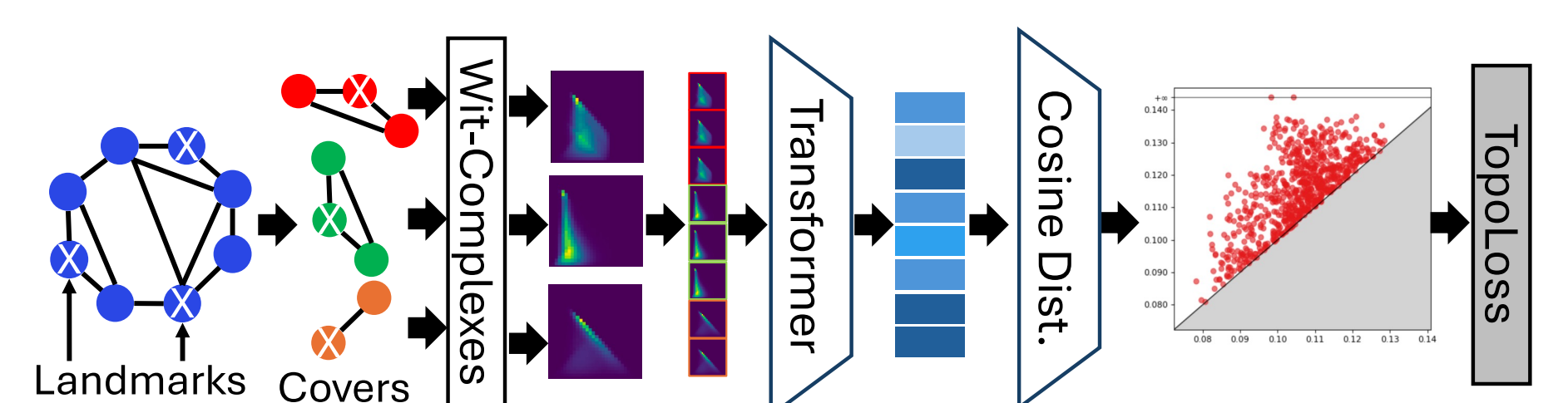
$$\mathbf{z} = [\mathbf{Z}_{T_L}, \mathbf{Z}_G]$$

Attention coefficients,  $\alpha_i = \text{Softmax}(\mathbf{W}_2 \cdot \tanh(\mathbf{W}_1 \mathbf{z}_i + \mathbf{b}_1))$

Additive attention,  $\mathbf{Z}_{AGG} = \alpha_1 \times \mathbf{Z}_{T_L} + \alpha_2 \times \mathbf{Z}_G$

$$\mathbf{Z}_{WGTL} = \mathbf{Z}_{AGG} \mathbf{Z}_{T_G}$$

## Topological Regularizer



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

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

## Conclusion

- Proposed the first topological defense against adversarial attacks on graphs.
- We have derived theoretical properties of WGTL, both at the local and global levels.
- WGTL improves robustness against a wide variety of attacks.

**Acknowledgements:** NSF grant TIP-2333703, ONR grant N00014-21-1-2530, ANR JCJC project REPUBLIC (ANR-22-CE23-0003-01), PEPR project FOUNDRY (ANR23-PEIA-0003), and the CHIST-ERA project CausalXRL (ANR-21-CHR4-0007).

arXiv Link

