



DShield: Defending against Backdoor Attacks on Graph Neural Networks via Discrepancy Learning

Hao Yu¹, Chuan Ma², Xinhang Wan¹, Jun Wang¹, Tao Xiang², Meng Shen³, and Xinwang Liu¹ National University of Defense Technology Chongqing University Beijing Institute of Technology

Speaker: Zhongming Wang







Feb. 27, 2025@San Diego, USA



- Many types of data can be represented as graph structures, e.g., social networks and protein molecules
- GNNs propagate information across graph structures, effectively capturing complex relationships between nodes and edges
- >Tasks: node classification, graph classification, link prediction









Social Networks

Drug Discovery

Smart City

Bank Risk Management



Backdoor attacks

- Inducing undesirable behaviors in backdoored models by injecting triggers into original graphs
- Two variants
 - Dirty-label backdoor attacks (**DLBAs**): Altering label information
 - Clean-label backdoor attacks (CLBAs): Altering attributes of normal nodes





>Two phenomena

- Semantic Drift of DLBAs: between node's attributes and structures, and their labels
- Attribute over-emphasis of CLBAs: high similarities of important attributes



Semantic Drift

Attribute Over-emphasis



➢ Pipeline

- Auxiliary Model Training: Preparing to identify poisoned nodes in DLBAs
- Discrepancy Matrix Construction: Detecting poisoned nodes in both DLBAs and CLBAs
- Backdoor-free Model Training: Training a model without backdoors





Backdoored Model

$$\mathcal{L} = rac{1}{|\mathcal{V}^{train}|} \sum_{v_i \in \mathcal{V}^{train}} \operatorname{CE} \left(\widetilde{f}(oldsymbol{A}, oldsymbol{X})_i, oldsymbol{y}_i
ight)$$

Self-supervised Model

- View Augmentation: Sampling two stochastic augmentation functions
- View Encoding: Extracting latent representations of nodes
- Contrast and Reconstruction: Distinguishing between representations of two nodes





View Augmentation

$$P\left\{(v_i, v_j) \in \hat{\mathcal{E}}
ight\} = 1 - p_{ij}^s; \quad p_{ij}^s = rac{\exp\left(-s\left(\widetilde{oldsymbol{h}}_i, \widetilde{oldsymbol{h}}_j
ight)
ight)}{\sum_{(v_k, v_t) \in \widetilde{\mathcal{E}}} \exp\left(-s\left(\widetilde{oldsymbol{h}}_k, \widetilde{oldsymbol{h}}_t
ight)
ight)}$$

Structure-level

$$oldsymbol{g}_i = rac{\partial ext{CE}(f(\widetilde{oldsymbol{A}},\widetilde{oldsymbol{X}};\phi)_i,oldsymbol{y}_i)}{\partialoldsymbol{x}_i}; \quad \widehat{oldsymbol{A}} = [oldsymbol{x}_1 \odot \widehat{oldsymbol{m}}_1; \cdots; oldsymbol{x}_N \odot \widehat{oldsymbol{m}}_N]^{ op}$$

Attribute-level

≻View Encoding

$$\widehat{oldsymbol{H}}_1 = g_{ ext{enc}}(\widehat{oldsymbol{A}}_1, \widehat{oldsymbol{X}}_1), \qquad \widehat{oldsymbol{H}}_2 = g_{ ext{enc}}(\widehat{oldsymbol{A}}_2, \widehat{oldsymbol{X}}_2)$$

Contrast and Reconstruction

Contrast

$$\mathcal{L}_{c_1} = \sum_{v_i \in \mathcal{V}^{ ext{train}}} - rac{1}{|\mathcal{P}(v_i)|} \sum_{v_j \in \mathcal{P}(v_i)} \log rac{\exp(s(\widehat{oldsymbol{h}}_{1,i}, \widehat{oldsymbol{h}}_{2,j})/ au)}{\sum_{v_k \in \mathcal{N}(v_i)} \exp(s(\widehat{oldsymbol{h}}_{1,i}, \widehat{oldsymbol{h}}_{2,k})/ au)};$$

 $\widehat{\boldsymbol{H}}_{1}^{\prime} = [\widehat{\boldsymbol{h}}_{1} \odot \widehat{\boldsymbol{m}}_{1}^{\prime}; \widehat{\boldsymbol{h}}_{2} \odot \widehat{\boldsymbol{m}}_{2}^{\prime}; \cdots; \widehat{\boldsymbol{h}}_{N} \odot \widehat{\boldsymbol{m}}_{N}^{\prime}]^{\mathsf{T}} \qquad \widehat{\boldsymbol{H}}_{1}^{\prime} = g_{\mathrm{dec}}(\widehat{\boldsymbol{A}}_{1}, \widehat{\boldsymbol{X}}_{1}^{\prime}), \quad \widehat{\boldsymbol{H}}_{2}^{\prime} = g_{\mathrm{dec}}(\widehat{\boldsymbol{A}}_{2}, \widehat{\boldsymbol{X}}_{2}^{\prime})$

$${\mathcal{L}}_r \!=\! rac{1}{2} \sum_{i=1}^N (1-\!s(\stackrel{\sim}{m{x}}_i, \stackrel{\wedge}{m{x}}_{1,i}'))^2 \!+\! rac{1}{2} \sum_{i=1}^N (1-\!s(\stackrel{\sim}{m{x}}_i, \stackrel{\wedge}{m{x}}_{2,i}'))^2$$



Semantic Discrepancy Matrix

$$\widetilde{\boldsymbol{H}}_{sl} = \widetilde{f}(\widetilde{\boldsymbol{A}}, \widetilde{\boldsymbol{X}}; \phi_{1}), \quad \widetilde{\boldsymbol{H}}_{ssl} = g_{enc}(\widetilde{\boldsymbol{A}}, \widetilde{\boldsymbol{X}})$$

$$\boldsymbol{D}_{sl}^{y} = \boldsymbol{1} \operatorname{diag}(\widetilde{\boldsymbol{H}}_{sl}^{y} \widetilde{\boldsymbol{H}}_{sl}^{y^{\top}})^{\top} - 2\widetilde{\boldsymbol{H}}_{sl}^{y} \widetilde{\boldsymbol{H}}_{sl}^{y^{\top}} + \operatorname{diag}(\widetilde{\boldsymbol{H}}_{sl}^{y} \widetilde{\boldsymbol{H}}_{sl}^{y^{\top}}) \boldsymbol{1}^{\top};$$

$$\boldsymbol{D}_{ssl}^{y} = \boldsymbol{1} \operatorname{diag}(\widetilde{\boldsymbol{H}}_{ssl}^{y} \widetilde{\boldsymbol{H}}_{ssl}^{y^{\top}})^{\top} - 2\widetilde{\boldsymbol{H}}_{ssl}^{y} \widetilde{\boldsymbol{H}}_{ssl}^{y^{\top}} + \operatorname{diag}(\widetilde{\boldsymbol{H}}_{ssl}^{y} \widetilde{\boldsymbol{H}}_{ssl}^{y^{\top}}) \boldsymbol{1}^{\top},$$

$$\boldsymbol{D}_{ssl}^{y} = \boldsymbol{1} \operatorname{diag}(\widetilde{\boldsymbol{H}}_{ssl}^{y} \widetilde{\boldsymbol{H}}_{ssl}^{y^{\top}})^{\top} - 2\widetilde{\boldsymbol{H}}_{ssl}^{y} \widetilde{\boldsymbol{H}}_{ssl}^{y^{\top}} + \operatorname{diag}(\widetilde{\boldsymbol{H}}_{ssl}^{y} \widetilde{\boldsymbol{H}}_{ssl}^{y^{\top}}) \boldsymbol{1}^{\top},$$

Attribute Importance Discrepancy Matrix

$$\begin{split} \boldsymbol{g}_{i} &= \frac{\partial \mathrm{CE}(f(\widetilde{\boldsymbol{A}}, \widetilde{\boldsymbol{X}}; \phi)_{i}, \boldsymbol{y}_{i})}{\partial \boldsymbol{x}_{i}}, \quad w_{ij} = \mathbb{1}_{g_{ij} > 0} \qquad \widetilde{\boldsymbol{X}}^{y} = \widetilde{\boldsymbol{X}}^{y} \circ \boldsymbol{W}^{y} \\ \overline{\boldsymbol{X}}^{y} &= \mathrm{UMAP}(\widetilde{\boldsymbol{X}}^{y}, d) \\ \boldsymbol{D}_{a}^{y} &= \boldsymbol{1} \operatorname{diag}(\overline{\boldsymbol{X}}^{y} \, \overline{\boldsymbol{X}}^{y^{\top}})^{\top} - 2\overline{\boldsymbol{X}}^{y} \, \overline{\boldsymbol{X}}^{y^{\top}} + \operatorname{diag}(\overline{\boldsymbol{X}}^{y} \, \overline{\boldsymbol{X}}^{y^{\top}}) \boldsymbol{1}^{\top} \end{split}$$

Clustering

 $oldsymbol{D}_y = rac{\mathrm{std}(oldsymbol{D}_s^y)}{\mathrm{std}(oldsymbol{D}_s^y) + \mathrm{std}(oldsymbol{D}_a^y)} oldsymbol{D}_s^y + rac{\mathrm{std}(oldsymbol{D}_a^y)}{\mathrm{std}(oldsymbol{D}_s^y) + \mathrm{std}(oldsymbol{D}_a^y)} oldsymbol{D}_a^y; \hspace{0.3cm} \mathcal{V}_1, \mathcal{V}_2 = \mathrm{HDBSCAN}(oldsymbol{D}_y)$





Importance of UMAP



Target Label Discovery

• Model Traning

$$\begin{split} \widetilde{A}_{ij}^{y,\text{pos}} & \begin{cases} = \widetilde{A}_{ij}, & \text{if } v_i \notin \mathcal{V}_2^y \text{ and } v_j \notin \mathcal{V}_2^y \\ = 0, & \text{if } v_i \in \mathcal{V}_2^y \text{ or } v_j \in \mathcal{V}_2^y \end{cases} & \qquad \widetilde{A}_{ij}^{y,\text{ neg}} \begin{cases} = \widetilde{A}_{ij}, & \text{if } v_i \in \mathcal{V}_2^y \text{ or } v_j \in \mathcal{V}_2^y \\ = 0, & \text{if } v_i \notin \mathcal{V}_2^y \text{ and } v_j \notin \mathcal{V}_2^y \end{cases} \\ \mathcal{L}_{\text{pos}} = \frac{1}{|\mathcal{V}_1^y|} \sum_{v_i \in \mathcal{V}_1^y} \text{CE}(f'(\widetilde{\boldsymbol{A}}^{y,\text{pos}}, \widetilde{\boldsymbol{X}})_i, \boldsymbol{y}_i) & \qquad \mathcal{L}_{\text{neg}} = \frac{1}{|\mathcal{V}_2^y|} \sum_{v_i \in \mathcal{V}_2^y} \text{CE}(f'(\widetilde{\boldsymbol{A}}^{y,\text{neg}}, \widetilde{\boldsymbol{X}})_i, \boldsymbol{y}_i). & \qquad \mathcal{L} = \mathcal{L}_{\text{pos}} - \log\left(\mathcal{L}_{\text{neg}}\right) \end{split}$$

• Target Label

$$egin{aligned} egin{aligned} egin{aligned} eta^y &= rac{1}{|\mathcal{V}_1^y|} \sum_{v_i \in \mathcal{V}_1^y} f'(\widetilde{oldsymbol{A}}^{y, ext{pos}}, \widetilde{oldsymbol{X}}; \phi_1')_i & s^y &= rac{1}{|\mathcal{V}_2^y|} \sum_{v_i \in \mathcal{V}_2^y} \left\| f'(\widetilde{oldsymbol{A}}^{y, ext{neg}}, \widetilde{oldsymbol{X}}; \phi_1')_i - oldsymbol{h}^y
ight\|_2^2 \ &y &= \left\{ y \,|\, y \!\in\! [1,C] \wedge rac{|s^y - ext{median}(oldsymbol{s})|}{ ext{mad}(oldsymbol{s})} \leq eta
ight\} \end{aligned}$$

Backdoor-free Training

$$\mathcal{L}_{\text{clf}} = \frac{1}{|\mathcal{V}^{\dagger}|} \sum_{v_i \in \mathcal{V}^{\dagger}} \text{CE}(f(\widetilde{\boldsymbol{A}}, \widetilde{\boldsymbol{X}})_i, \boldsymbol{y}_i) \qquad \qquad \mathcal{L}_{\text{penalty}} = \frac{1}{|\mathcal{V}^{\dagger}|} \sum_{v_i \in \mathcal{V}^{\dagger}} \text{CE}(f(\widetilde{\boldsymbol{A}}', \widetilde{\boldsymbol{X}})_i, \boldsymbol{y}_i) \qquad \qquad \mathcal{L} = \mathcal{L}_{\text{clf}} - \gamma \log \mathcal{L}_{\text{penalty}}$$



Elements of Discrepancy Matrices

• Differences in elements of two discrepancy matrices between normal and malicious nodes in the Cora dataset.





Seven DLBAs & Four Datasets

Datasets	Defenses	SBA [50]		GTA [38]		EBA	EBA [43]		GB-FGSM [4]		LGCB [4]		UGBA [5]		TRAP [44]	
	Derenses	ASR \downarrow	ACC ↑	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC ↑									
Cora	no-defense	$53.14_{\pm 2.10}$	84.07 _{±1.02}	$96.22_{\pm 4.41}$	82.74 _{±2.38}	71.81 _{±9.44}	83.48 ±1.32	$97.34_{\pm 1.64}$	82.52 ± 1.53	97.79 _{±1.01}	83.70 _{±0.52}	97.51 _{±2.07}	83.03 _{±1.32}	$19.86_{\pm 2.42}$	84.15 ±1.09	
	Prune [5]	$30.26_{\pm 2.78}$	83.26 ± 0.96	$16.42_{\pm 2.35}$	84.00 ±1.80	$100.0_{\pm 0.00}$	$80.74_{\pm 1.57}$	$56.98_{\pm 1.00}$	$81.92_{\pm 3.36}$	$56.98_{\pm 0.93}$	81.92 _{±1.95}	$54.47_{\pm 2.19}$	$83.41_{\pm 1.52}$	$11.36_{\pm 2.08}$	$81.48_{\pm 1.87}$	
	Isolate [5]	$04.06_{\pm 2.42}$	$37.19_{\pm 1.07}$	$03.54_{\pm 1.85}$	$77.78_{\pm 1.48}$	$96.43_{\pm 1.66}$	80.67 _{±1.29}	$56.09_{\pm 0.78}$	$81.85_{\pm 1.08}$	$56.46_{\pm 0.78}$	$81.63_{\pm 1.97}$	$56.02_{\pm 0.80}$	$80.74_{\pm 3.21}$	10.15 $_{\pm 1.94}$	$67.26_{\pm 1.63}$	
	PXGBD [9]	$46.37_{\pm 2.03}$	84.00 ± 0.85	$97.05_{\pm 2.46}$	$83.40_{\pm 1.74}$	$71.04_{\pm 9.65}$	$82.74_{\pm 1.40}$	$97.86_{\pm 1.44}$	82.74 _{±0.93}	$98.38_{\pm1.64}$	$82.37_{\pm 1.13}$	$97.86_{\pm 2.17}$	$\textbf{84.00}_{\pm 0.55}$	$22.14_{\pm 2.31}$	$81.55_{\pm 1.37}$	
	GXGBD [9]	$56.28_{\pm 1.31}$	84.17 $_{\pm 1.30}$	$83.77_{\pm 1.56}$	$83.42_{\pm 1.10}$	$88.19_{\pm 2.81}$	$73.52_{\pm 2.75}$	$98.16_{\pm 1.16}$	84.07 ±1.09	$98.43_{\pm1.01}$	84.26 ±0.98	$80.81_{\pm 2.30}$	$83.15_{\pm 2.31}$	$19.70_{\pm 3.84}$	$81.18_{\pm 0.80}$	
	ABL [21]	$54.86_{\pm 2.37}$	$82.81_{\pm 1.40}$	$97.79_{\pm 2.09}$	$82.59_{\pm 2.38}$	$81.43_{\pm 3.98}$	83.26 ± 1.37	$97.56_{\pm 1.34}$	$83.11_{\pm 0.81}$	$98.08_{\pm0.84}$	83.70 ± 0.59	$97.42_{\pm 1.04}$	$83.89_{\pm 2.19}$	$18.38_{\pm 2.52}$	$83.33_{\pm 0.64}$	
	RS [49]	$08.12_{\pm 1.51}$	$80.81_{\pm 1.03}$	$94.83_{\pm 2.93}$	80.67 _{±0.99}	$58.92_{\pm 2.72}$	$81.11_{\pm 1.05}$	97.27 _{±0.85}	$80.15_{\pm 1.24}$	96.97 _{±1.09}	$80.15_{\pm 1.84}$	$90.41_{\pm 1.56}$	$80.81_{\pm 1.03}$	$20.81_{\pm 1.78}$	$79.85_{\pm 1.13}$	
	DShield	$01.33_{\pm 2.56}$	$81.78_{\pm 1.93}$	$00.74_{\pm 0.74}$	$81.92_{\pm 1.98}$	$02.95_{\pm 2.32}$	$81.50_{\pm 1.08}$	$02.51_{\pm 2.83}$	$82.29_{\pm 0.71}$	$00.89_{\pm 1.15}$	$82.57_{\pm 0.51}$	$01.33_{\pm 1.21}$	$82.15_{\pm 0.80}$	$13.38_{\pm 2.12}$	$82.07_{\pm 1.65}$	
	no-defense	$62.84_{\pm 5.05}$	85.17 _{±0.20}	96.54 _{±2.28}	85.09 _{±0.26}	84.79 _{±2.10}	85.23 _{±0.17}	99.67 _{±0.66}	$85.03_{\pm 0.15}$	97.14 _{±1.44}	85.16 _{±0.22}	93.27 _{±3.23}	85.20 _{±0.25}	48.34 _{±0.93}	85.21 _{±0.38}	
	Prune [5]	$53.01_{\pm 2.59}$	$85.33_{\pm 0.37}$	$46.15_{\pm 0.40}$	85.30 _{±0.21}	99.95 _{±0.00}	85.44 ±0.18	$70.28_{\pm 0.40}$	85.31 ±0.25	$69.36_{\pm 0.56}$	85.36 ±0.18	$66.85_{\pm 1.47}$	85.38 ±0.29	41.77 ±0.28	85.33 ±0.25	
PubMed	Isolate [5]	40.69 ± 0.47	81.25 ± 0.28	$45.46_{\pm 0.22}$	$84.47_{\pm 0.21}$	99.36 _{±0.38}	$84.56_{\pm 0.22}$	$70.16_{\pm 0.31}$	85.14 _{±0.15}	69.80 _{±0.90}	85.11 _{±0.13}	$67.62_{\pm 1.34}$	$85.14_{\pm 0.20}$	$43.30_{\pm 3.41}$	$83.34_{\pm 1.59}$	
	PXGBD [9]	$55.26_{\pm 3.06}$	$83.55_{\pm 1.01}$	99.67 _{±0.54}	83.17 _{±0.92}	$86.26_{\pm 5.89}$	$82.64_{\pm 1.57}$	$100.0_{\pm 0.00}$	$82.32_{\pm 1.47}$	$100.0_{\pm 0.00}$	82.38 ± 0.99	$91.80_{\pm 0.87}$	$82.61_{\pm 1.35}$	$49.80_{\pm 1.70}$	$81.98_{\pm 1.45}$	
	GXGBD [9]	$56.20_{\pm 1.18}$	$83.93_{\pm 0.76}$	99.71±0.35	$82.57_{\pm 0.64}$	$90.21_{\pm 2.27}$	$75.45_{\pm 2.23}$	$100.0_{\pm 0.00}$	$82.31_{\pm 0.85}$	$100.0_{\pm 0.00}$	82.41 ± 0.82	$81.95_{\pm 0.00}$	$82.82_{\pm 0.91}$	49.36 _{±3.16}	$82.84_{\pm 1.47}$	
	ABL [21]	$65.16_{\pm 4.55}$	85.21 ± 0.38	99.98 _{±0.03}	$85.13_{\pm 0.21}$	$89.08_{\pm 2.61}$	$84.98_{\pm 0.16}$	99.98 _{±0.03}	84.89 _{±0.13}	$100.0_{\pm 0.00}$	85.04 _{±0.17}	$91.45_{\pm 5.44}$	85.00 ± 0.30	49.36 _{±0.68}	$84.64_{\pm 1.01}$	
	RS [49]	$51.10_{\pm 3.63}$	$83.13_{\pm 0.45}$	97.53 _{±1.33}	82.99 _{±0.37}	89.86 _{±1.20}	$83.68_{\pm 0.32}$	$100.0_{\pm 0.00}$	83.39 _{±0.58}	$100.0_{\pm 0.00}$	83.40 _{±0.57}	$59.68_{\pm 4.13}$	$83.17_{\pm 0.31}$	$48.48_{\pm 3.50}$	$82.83_{\pm 0.71}$	
	DShield	35.91 ±3.43	82.01 ± 0.47	00.73 ±0.33	85.28 ± 0.12	$02.78_{\pm 1.58}$	85.02 ± 0.29	00.66 ±0.72	$83.71_{\pm 0.24}$	03.54 ±1.97	$84.11_{\pm 0.56}$	$02.24_{\pm 0.97}$	$84.70_{\pm 0.42}$	$45.77_{\pm 0.76}$	$82.17_{\pm 0.36}$	



Two CLBAs & Four Datasets

	Cora				PubMed				Flickr				OGBN-arXiv			
Defenses	GCBA	A [49]	PerCB	BA [45]	GCBA	A [49]	PerCB	A [45]	GCBA	A [49]	PerCB	A [45]	GCBA	[49]	PerCB	A [45]
	ASR \downarrow	ACC ↑	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC ↑
no-defense	81.74 _{±3.93}	84.74 ±1.27	84.87 _{±3.13}	$83.18_{\pm 1.54}$	$ 78.86_{\pm 7.76} $	85.30 _{±0.13}	$100.0_{\pm 0.00}$	85.23 _{±0.13}	$60.48_{\pm 4.72}$	50.84 ±0.16	$07.05_{\pm 5.11}$	50.86 ±0.18	26.84 _{±2.79}	60.43 ±0.30	$18.06_{\pm 1.48}$	59.81 _{±0.15}
Prune [5]	$99.70_{\pm 0.31}$	$82.37_{\pm 0.42}$	$69.01_{\pm 1.31}$	82.66 ± 0.55	$99.83_{\pm 0.18}$	85.55 $_{\pm 0.16}$	$95.58_{\pm 4.50}$	$\textbf{85.34}_{\pm 0.30}$	$78.50_{\pm 3.63}$	$49.97_{\pm 0.27}$	$30.38_{\pm 2.70}$	$50.12_{\pm 0.15}$	$23.14_{\pm 0.30}$	$60.40_{\pm 0.35}$	$19.03_{\pm 1.15}$	$59.78_{\pm 0.17}$
Isolate [5]	$99.88_{\pm 0.21}$	$71.92_{\pm 2.05}$	$81.73_{\pm 2.35}$	$69.04_{\pm 2.07}$	99.10 _{±1.09}	$84.12_{\pm 0.14}$	$97.88_{\pm 1.05}$	$84.26_{\pm 1.01}$	$66.44_{\pm 5.40}$	$50.34_{\pm 0.08}$	$30.15_{\pm 3.17}$	$50.52_{\pm 0.18}$	22.24 _{±0.98}	60.43 ±0.08	19.00 ± 0.64	$59.75_{\pm 0.18}$
PXGBD [8]	$82.29_{\pm 3.66}$	$82.89_{\pm 0.88}$	$27.68_{\pm 2.58}$	$82.74_{\pm 1.19}$	$96.91_{\pm 3.08}$	$82.82_{\pm 1.34}$	$93.11_{\pm 0.36}$	$83.31_{\pm 1.54}$	26.50 ± 0.79	$47.74_{\pm 0.98}$	$00.16_{\pm 0.28}$	$47.29_{\pm 2.67}$	$11.87_{\pm 0.71}$	$60.42_{\pm 0.17}$	00.78 ± 0.49	60.42 ±0.89
GXGBD [8]	$04.61_{\pm 1.77}$	$82.96_{\pm 0.52}$	80.69 _{±2.98}	84.30 ±1.07	$97.30_{\pm 0.51}$	$82.43_{\pm 1.11}$	$72.90_{\pm 1.18}$	82.78 ± 0.32	$49.90_{\pm 2.39}$	$40.57_{\pm 0.30}$	$02.42_{\pm 4.19}$	$40.61_{\pm 0.44}$	09.21 _{±1.54} :	$55.94_{\pm 0.16}$	$00.34_{\pm 0.14}$	$48.51_{\pm 0.43}$
ABL [20]	$79.70_{\pm 3.89}$	$84.30_{\pm 1.35}$	$88.68_{\pm 4.25}$	$83.70_{\pm 0.91}$	$83.40_{\pm 1.46}$	84.98 ± 0.23	$82.50_{\pm 3.53}$	$85.13_{\pm 0.36}$	$66.88_{\pm 0.84}$	$50.49_{\pm 0.21}$	$01.35_{\pm 0.60}$	$50.68_{\pm 0.04}$	02.73 _{±3.84} :	$58.40_{\pm 1.61}$	22.04 ± 2.39	$60.23_{\pm 0.79}$
RS [49]	$42.44_{\pm 4.11}$	$81.26_{\pm 1.92}$	$85.22_{\pm 3.15}$	$80.74_{\pm 1.14}$	$74.05_{\pm 1.70}$	83.28 ± 0.58	$88.46_{\pm 3.41}$	$83.51_{\pm 0.27}$	00.02 ±0.03	40.48 ± 0.09	23.72 ± 0.59	40.50 ± 0.16	27.76 _{±0.08} :	$55.36_{\pm 0.10}$	43.48 ± 5.32	$55.38_{\pm 0.19}$
DShield	01.11 ±1.04	$81.18_{\pm 1.32}$	03.41 ±2.53	$81.63_{\pm 2.04}$	01.91 ±1.66	$84.24_{\pm 0.24}$	08.09 ±2.04	$83.69_{\pm 2.01}$	$01.66_{\pm 2.34}$	$50.42_{\pm 0.25}$	$\textbf{00.04}_{\pm 0.05}$	50.86 ±0.06	00.00 ±0.00	$59.96_{\pm 0.54}$	00.00 ±0.00	$59.92_{\pm 0.44}$

- DShield exhibits a **notable** performance against most backdoor attacks, surpassing the efficacy of conventional defenses
- Although DShield demonstrates superior defense performance in most scenarios, a slight decline in model performance on normal nodes is observed

Robustness Evaluation



➢Six Poisoning Rates



≻Six Trigger Sizes



Five Adaptive Attacks

Detenata	A tto olira	no-de	efense	no-p	oison	DShield		
Datasets	Attacks	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC \uparrow	
	UGBA+LGCB	$97.91_{\pm 0.93}$	$84.32_{\pm 1.19}$	$08.00_{\pm 0.21}$	$78.02_{\pm 5.00}$	$00.49_{\pm 4.16}$	$82.22_{\pm 6.72}$	
	UGBA+GCBA	$88.07_{\pm 1.89}$	$83.83_{\pm 2.04}$	$08.00_{\pm 0.43}$	78.15 ± 4.98	$03.94_{\pm 4.49}$	81.97 _{±5.46}	
Cora	GCBA+PerCBA	$74.78_{\pm 3.77}$	$83.82_{\pm 1.13}$	$08.24_{\pm 0.93}$	$77.53_{\pm 3.36}$	$08.49_{\pm 2.82}$	$80.99_{\pm 4.64}$	
	AdaDA	$84.51_{\pm 2.61}$	$84.07_{\pm 0.52}$	$07.75_{\pm 0.64}$	$84.07_{\pm 1.29}$	$04.43_{\pm 3.13}$	$76.86_{\pm0.26}$	
	AdaCA	$89.85_{\pm 6.52}$	$82.59_{\pm 0.52}$	$09.47_{\pm 5.18}$	$84.69_{\pm 0.94}$	$00.74_{\pm 0.98}$	$81.23_{\pm 2.46}$	
	UGBA+LGCB	$ 88.73_{\pm 5.16} $	$84.98_{\pm 0.26}$	$37.92_{\pm 1.71}$	82.23 _{±1.97}	$00.15_{\pm 0.27}$	$82.50_{\pm 0.74}$	
	UGBA+GCBA	91.89 ± 3.07	$85.42_{\pm 0.03}$	39.61 _{±4.17}	83.66 ± 0.43	$02.92_{\pm 1.17}$	$82.51_{\pm 1.19}$	
PubMed	GCBA+PerCBA	$89.01_{\pm 5.25}$	84.53 ± 0.62	$39.52_{\pm 3.38}$	83.49 ± 0.61	$13.39_{\pm 2.08}$	$81.28_{\pm 1.17}$	
	AdaDA	$88.72_{\pm 3.34}$	85.22 ± 0.04	$40.23 _{\pm 0.72}$	$83.17_{\pm 3.57}$	32.61 ± 8.89	$79.17_{\pm 1.47}$	
	AdaCA	$ 87.88_{\pm 5.12} $	$85.13_{\pm 0.31}$	$03.82_{\pm 1.68}$	$85.39_{\pm 0.13}$	$05.68_{\pm 0.79}$	82.60 ± 1.87	

• The efficacy of DShield stems from the execution of the selfsupervised learning framework and UMAP on the manipulated graph 13

Sensitivity Analysis







Defending against Attacks on Graph Classification Tasks

Dotogoto	Defenses	G-S	BA	G-E	EBA	G-GCBA		
Datasets	Defenses	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC \uparrow	ASR \downarrow	ACC ↑	
	no-defense	$100.0_{\pm 0.00}$	32.67 _{±0.94}	$100.0_{\pm 0.00}$	$27.50_{\pm 3.54}$	$100.0_{\pm 0.00}$	$23.33_{\pm 0.00}$	
ENZ I MES	G-DShield	$01.67_{\pm 2.35}$	$35.00_{\pm 2.36}$	$02.50_{\pm 3.54}$	$31.67_{\pm 2.35}$	$00.00_{\pm 0.00}$	$30.84_{\pm 1.18}$	
PROTEINS	no-defense	$100.0_{\pm 0.00}$	$66.67_{\pm 1.27}$	$99.56_{\pm 0.63}$	$71.62_{\pm 3.18}$	99.11 _{±0.00}	$73.43_{\pm 5.73}$	
INOTEINS	G-DShield	$00.00_{\pm 0.00}$	$74.63_{\pm 4.04}$	$04.29_{\pm 6.06}$	$74.33_{\pm 3.19}$	$00.98_{\pm 1.39}$	$71.17_{\pm 1.27}$	
MNIST	no-defense	$100.0_{\pm 0.00}$	$36.33_{\pm 0.67}$	$100.0_{\pm 0.00}$	$36.89_{\pm 0.02}$	$100.0_{\pm 0.00}$	$38.34_{\pm 0.66}$	
10110101	G-DShield	$00.00_{\pm 0.00}$	38.86 ± 0.37	$01.87_{\pm 2.64}$	39.77 _{±1.65}	$00.00_{\pm 0.00}$	$40.57_{\pm 1.89}$	

• Experimental results show that semantic drift and attribute over-emphasis also occur in graph classification tasks and the defense efficiency is not limited to node classification backdoor attacks

Thank You!

Hao Yu

csyuhao@gmail.com

National University of Defense Technology

Source Code: https://github.com/csyuhao/DShield