

CONTRACTIVE SYSTEMS IMPROVE GRAPH NEURAL NETWORKS AGAINST ADVERSARIAL ATTACKS

Anonymous authors

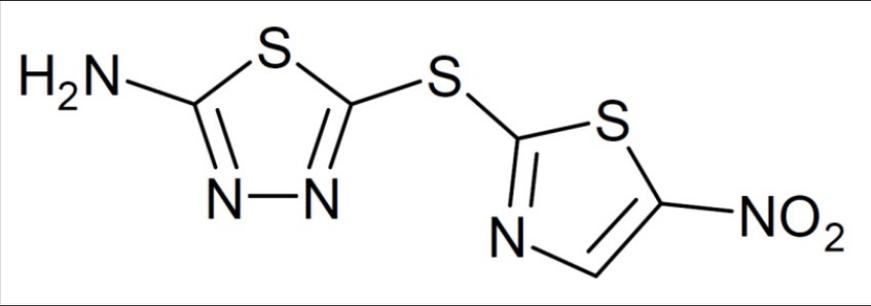
Paper under double-blind review

ABSTRACT

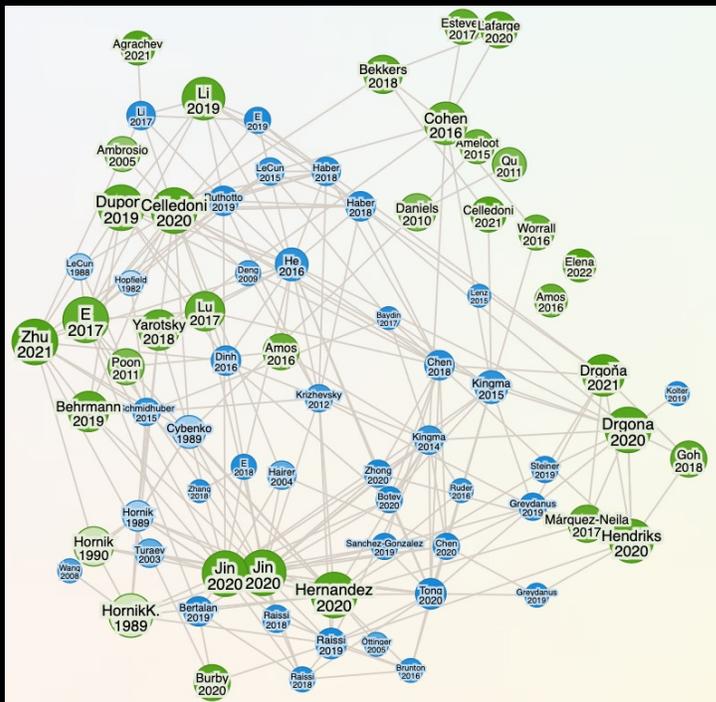
Graph Neural Networks (GNNs) have established themselves as a key component in addressing diverse graph-based tasks. Despite their notable successes, GNNs remain susceptible to input perturbations in the form of adversarial attacks. This paper introduces an innovative approach to fortify GNNs against adversarial perturbations through the lens of contractive dynamical systems. Our method introduces graph neural layers based on differential equations with contractive properties, which, as we show, improve the robustness of GNNs. A distinctive feature of the proposed approach is the simultaneous learned evolution of both the node features and the adjacency matrix, yielding an intrinsic enhancement of model robustness to perturbations in the input features and the connectivity of the graph. We mathematically derive the underpinnings of our novel architecture and provide theoretical insights to reason about its expected behavior. We demonstrate the efficacy of our method through numerous real-world benchmarks, reading on par or improved performance compared to existing methods.

With Moshe Eliasof, Ferdia Sherry and Carola Schönlieb

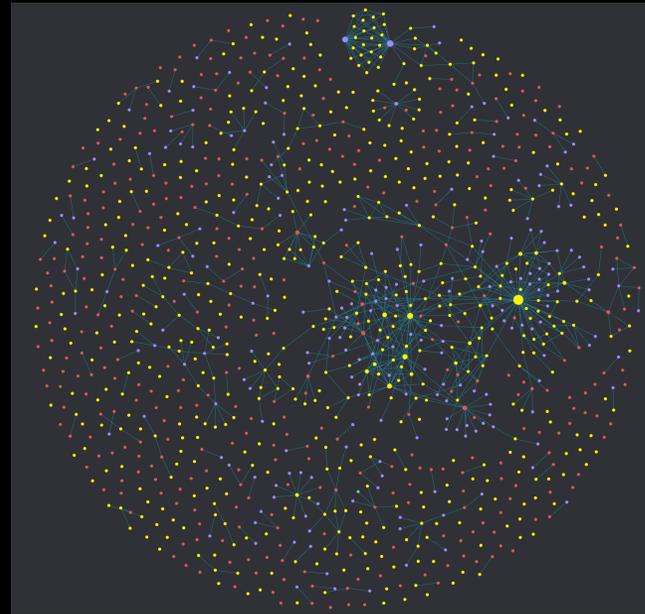
Graphs are everywhere



Molecule structure



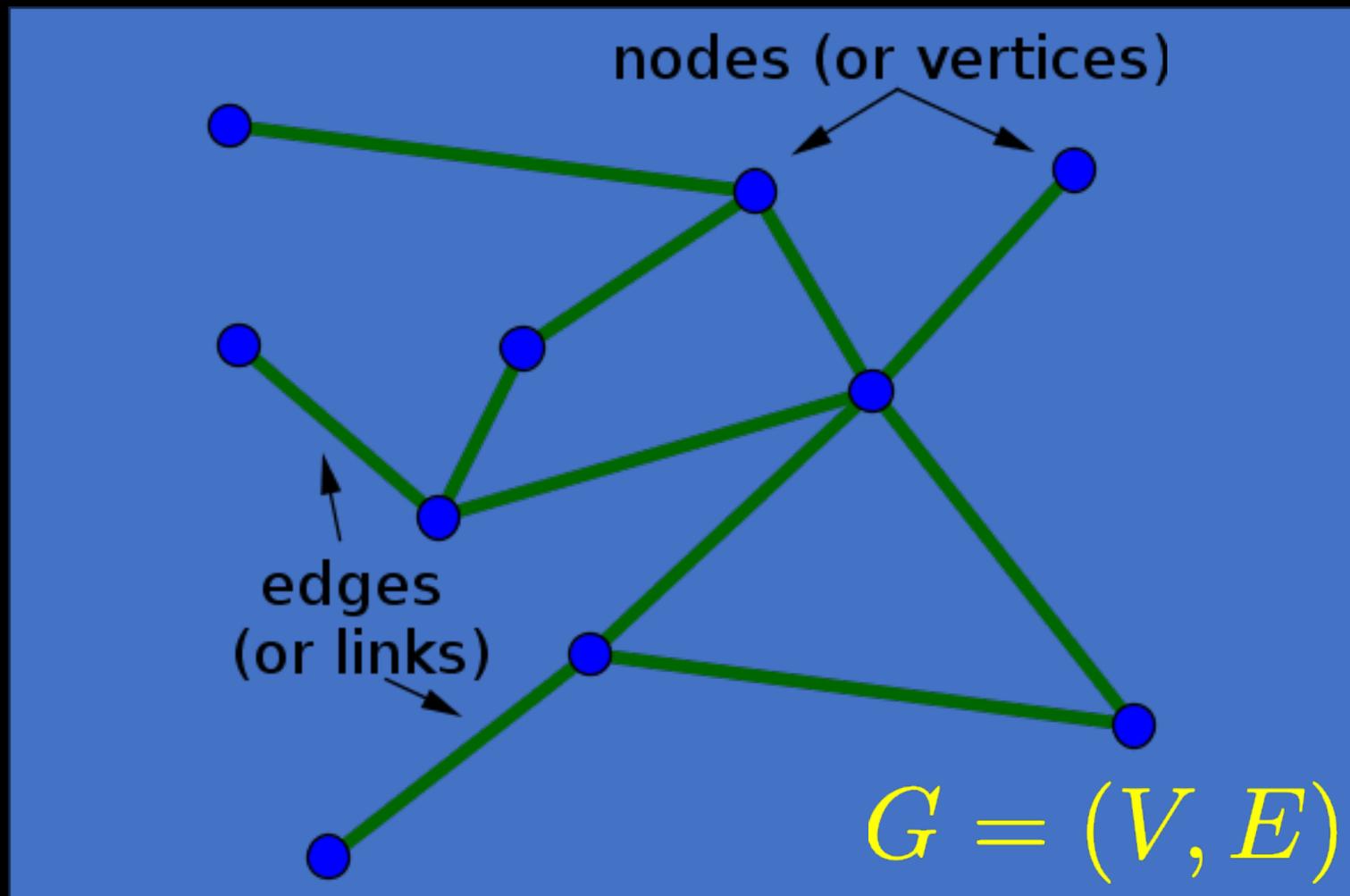
Citation graph of papers I saved



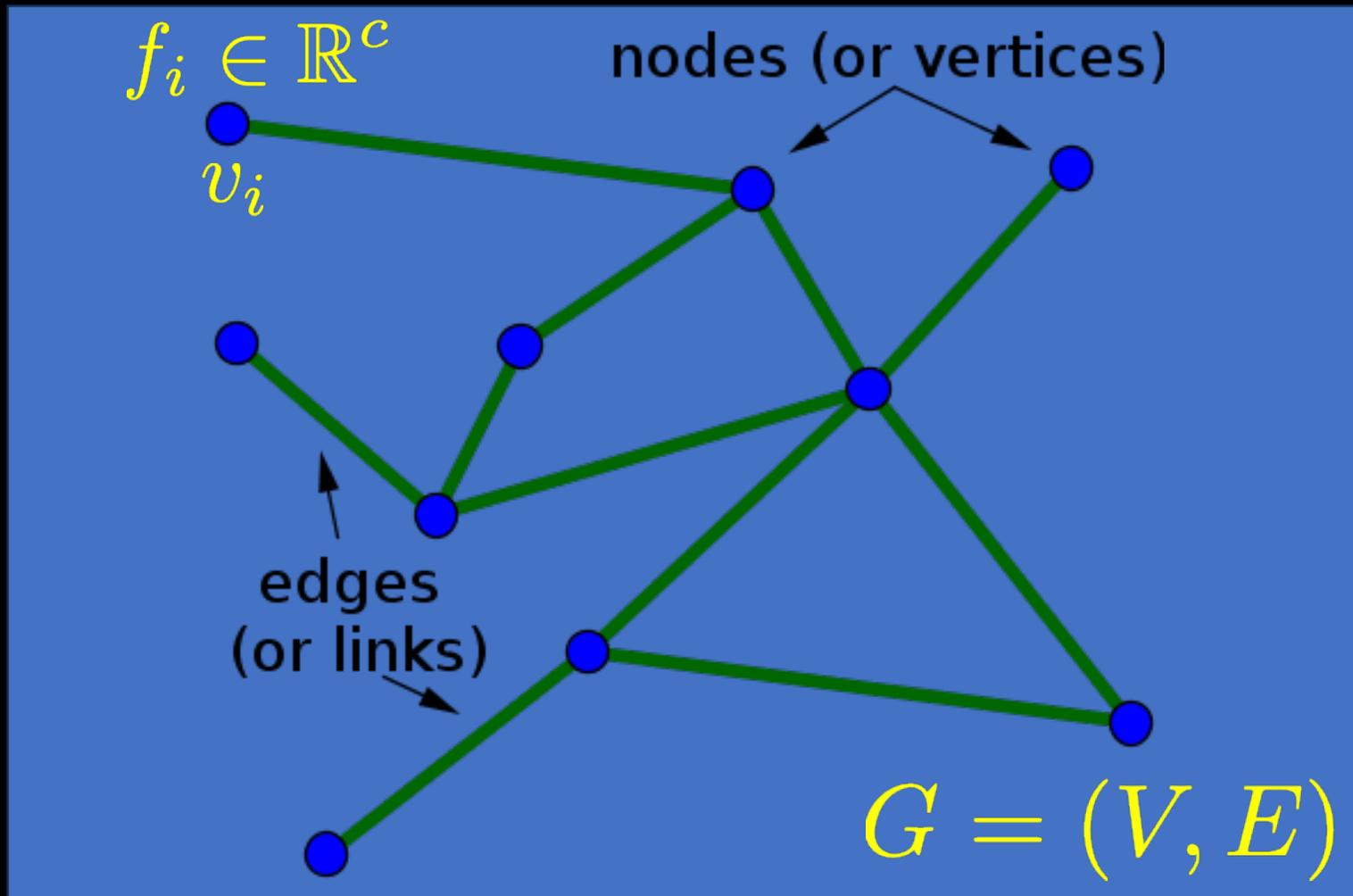
Reference graph some notes
I wrote in Obsidian



Connectivity of neurons
in the brain



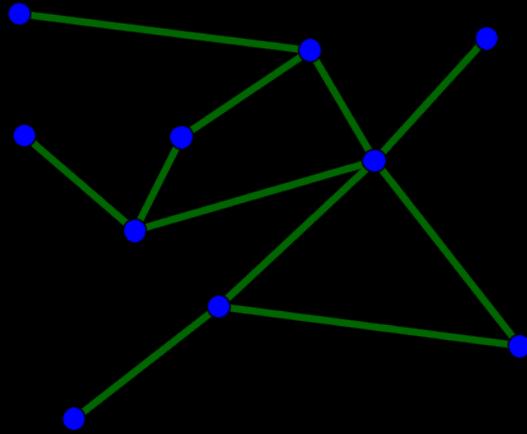
$$A \in \{0, 1\}^{n \times n}, \quad A^T = A$$



$$F \in \mathbb{R}^{n \times c}, e_i^T F = f_i$$

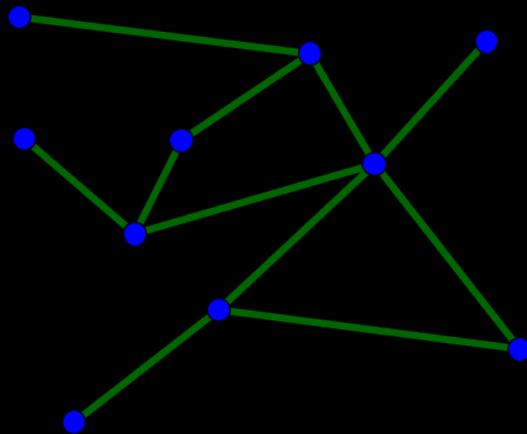
Classical tasks solved with GNNs

Graph classification



It is a protein

Node classification



Red
Blue
Blue
...
Green

Usual structure of GNNs

$$F^{(0)} = F$$

$$F^{(l+1)} = T_l \left(F^{(l)}, A \right), l = 0, \dots, L - 1$$

$$F^{(O)} = \text{MLP} \left(F^{(L)} \right) =: \text{GNN}(F, A)$$

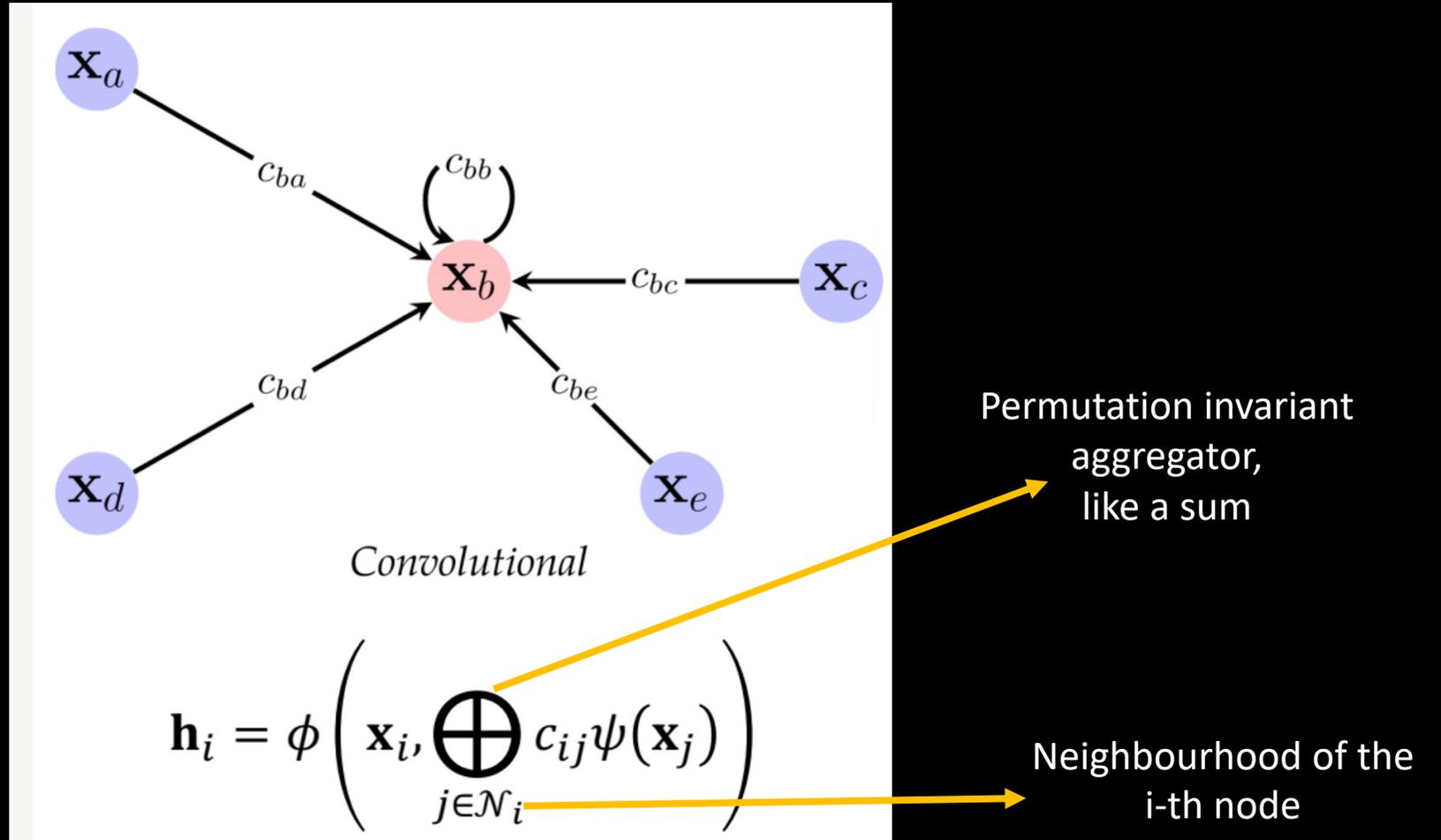
Invariant

$$\text{GNN}(F, A) = \text{GNN}(PF, PAP^T)$$

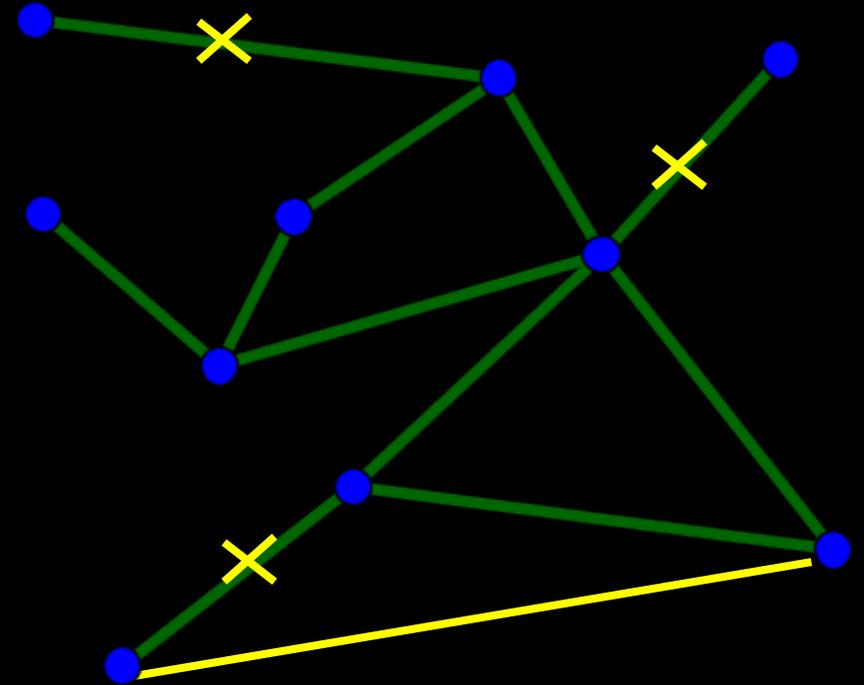
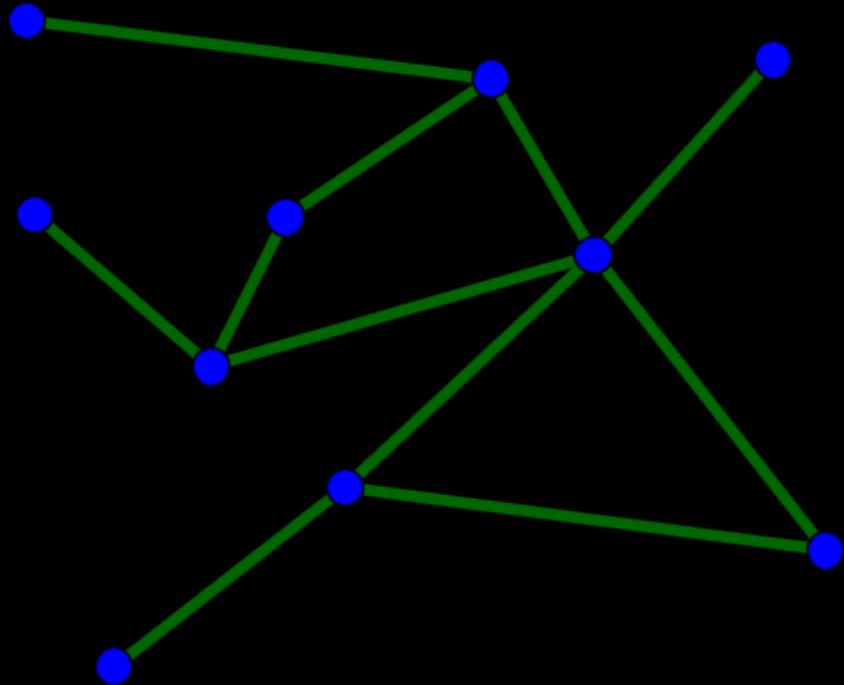
$$PGNN(F, A) = \text{GNN}(PF, PAP^T)$$

Equivariant

Usual structure of GNNs



Adversarial attacks



e.g. Add/remove a friendship
on Facebook

Adversarial attacks

$$A_* = A + \delta A, \quad \|\delta A\|_0 \leq \varepsilon_1$$

$$F_* = F + \delta F, \quad \|\delta F\|_F \leq \varepsilon_2$$

Attacks do not break the properties of symmetry generally

Goal:

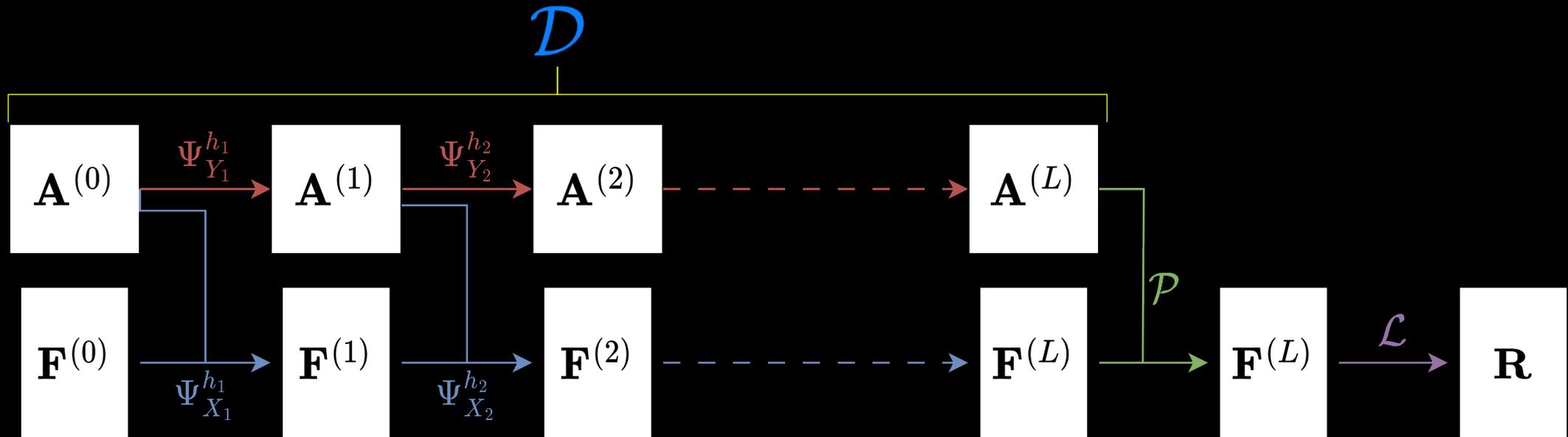
$$\text{GNN}(F, A) \approx \text{GNN}(F_*, A_*)$$

Remark on Nuclear Norm

$$A \in \{0, 1\}^{n \times n} \implies \|A\|_0 = \#\{i, j \in \{1, \dots, n\} : A_{ij} \neq 0\} = \|\text{vec}(A)\|_{\ell^1}$$

The 1-norm of the vectorisation is better suited for what we do, and we will use such norm instead of the nuclear norm.

Our proposed architecture: CSGNN



$$(\mathbf{F}^{(0)}, \mathbf{A}^{(0)}) := (\mathcal{K}(\mathbf{F}_*), \mathbf{A}_*)$$

$$\Psi_{X_i}^{h_i}(\mathbf{F}, \mathbf{A}) = \mathbf{F} - h_i \mathbf{G}(\mathbf{A})^T \sigma(\mathbf{G}(\mathbf{A}) \mathbf{F} \mathbf{W}_i) \mathbf{W}_i^T \left(\frac{\mathbf{K}_i + \mathbf{K}_i^T}{2} \right)$$

$$\Psi_{Y_i}^{h_i}(\mathbf{A}) = \mathbf{A} + h_i \sigma(\mathbf{M}_i(\mathbf{A}))$$

Our proposed architecture: CSGNN

$$\begin{aligned} M(A) = & k_1 A + k_2 \text{diag}(\text{diag}(A)) + \frac{k_3}{2n} (A \mathbf{1}_n \mathbf{1}_n^\top + \mathbf{1}_n \mathbf{1}_n^\top A) + k_4 \text{diag}(A \mathbf{1}_n) \\ & + \frac{k_5}{n^2} (\mathbf{1}_n^\top A \mathbf{1}_n) \mathbf{1}_n \mathbf{1}_n^\top + \frac{k_6}{n} (\mathbf{1}_n^\top A \mathbf{1}_n) I_n + \frac{k_7}{n^2} (\mathbf{1}_n^\top \text{diag}(A)) \mathbf{1}_n \mathbf{1}_n^\top \\ & + \frac{k_8}{n} (\mathbf{1}_n^\top \text{diag}(A)) I_n + \frac{k_9}{2n} (\text{diag}(A) \mathbf{1}_n^\top + \mathbf{1}_n (\text{diag}(A))^\top) \end{aligned}$$

$$M(PAP^T) = PM(A)P^T, \quad (M(A))^T = M(A)$$

Robustness of the network

Theorem 2 (Equation (8) can induce contractive node dynamics). *There are choices of $(\mathbf{W}_l, \mathbf{K}_l) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{c \times c}$, for which the explicit Euler step in Equation (8) is contractive for a small enough $h_l > 0$, i.e.*

$$\|\Psi_{X_l}^{h_l}(\mathbf{F} + \delta\mathbf{F}, \mathbf{A}) - \Psi_{X_l}^{h_l}(\mathbf{F}, \mathbf{A})\|_F \leq \|\delta\mathbf{F}\|_F, \quad \delta\mathbf{F} \in \mathbb{R}^{n \times c}. \quad (10)$$

Theorem 3 (Equation (14) defines contractive adjacency dynamics). *Let $\alpha \leq 0$, $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be a Lipschitz continuous function, with $\sigma'(s) \in [0, 1]$ almost everywhere. If $0 \leq h_l \leq \hat{h}_l^{\text{adj}} := \frac{2}{(2 \sum_{i=2}^9 |k_i|) - \alpha}$, then the explicit Euler step*

$$\mathbf{A}^{(l)} = \Psi_{Y_l}^{h_l}(\mathbf{A}^{(l-1)}) := \mathbf{A}^{(l-1)} + h_l \sigma \left(M(\mathbf{A}^{(l-1)}) \right), \quad (15)$$

where $k_1 = \left(\alpha - \sum_{i=2}^9 |k_i| \right)$, is contractive in the vectorized ℓ^1 norm.

$$\begin{aligned} d(\mathcal{D}(\mathbf{F}^{(0)}, \mathbf{A}^{(0)}), \mathcal{D}(\mathbf{F}_*^{(0)}, \mathbf{A}_*^{(0)})) &:= \|\text{vec}(\mathbf{A}^{(L)}) - \text{vec}(\mathbf{A}_*^{(L)})\|_1 + \|\mathbf{F}^{(L)} - \mathbf{F}_*^{(L)}\|_F \\ &\leq \varepsilon_1 + \varepsilon_2 \left(1 + \sum_{i=1}^L \text{Lip}(X_{i, \mathbf{F}^{(i-1)}}) h_i \right) \\ &=: \varepsilon_1 + c(h_1, \dots, h_L) \varepsilon_2. \end{aligned}$$

Experimental setup

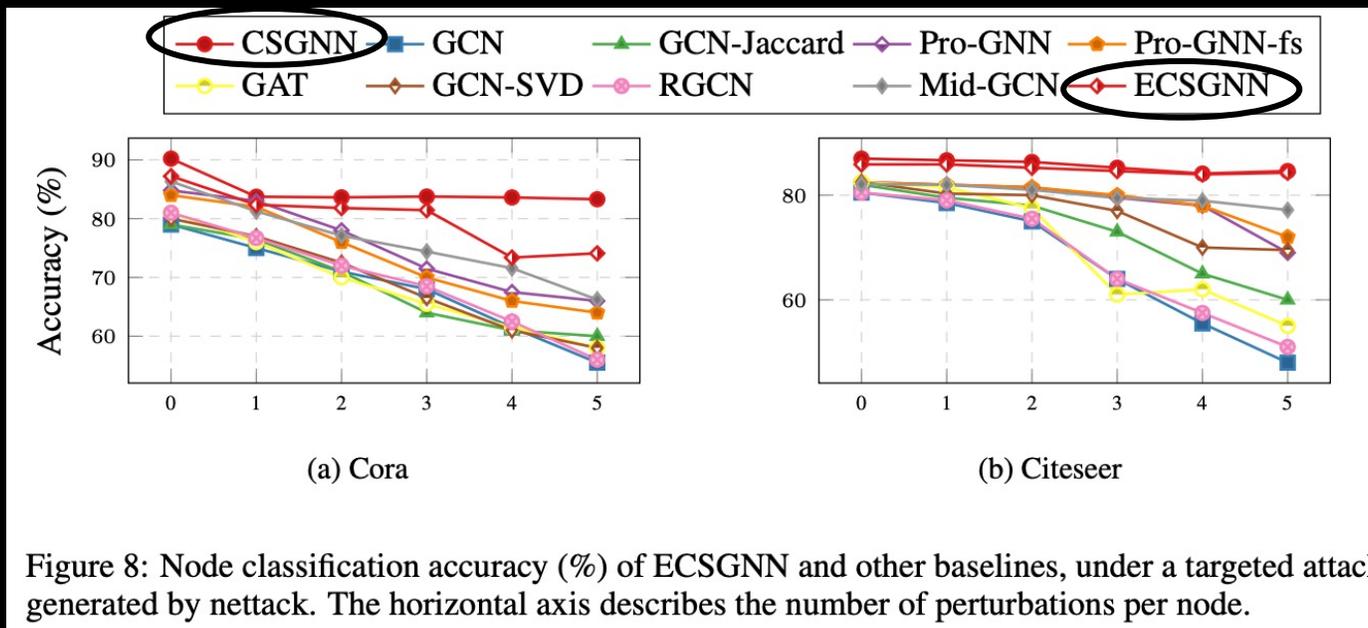
Table 3: Hyperparameter ranges

Hyperparameter	Range	Distribution
input/output embedding learning rate	$[10^{-5}, 10^{-2}]$	uniform
node dynamics learning rate	$[10^{-5}, 10^{-2}]$	uniform
adjacency dynamics learning rate	$[10^{-5}, 10^{-2}]$	uniform
input/output embedding weight decay	$[5 \cdot 10^{-8}, 5 \cdot 10^{-2}]$	log uniform
node dynamics weight decay	$[5 \cdot 10^{-8}, 5 \cdot 10^{-2}]$	log uniform
adjacency dynamics weight decay	$[5 \cdot 10^{-8}, 5 \cdot 10^{-2}]$	log uniform
input/output embedding dropout	$[0, 0.6]$	uniform
node dynamics dropout	$[0, 0.6]$	uniform
share weights between time steps	{yes, no}	discrete uniform
step size h	$[10^{-2}, 1]$	log uniform
adjacency contractivity parameter α	$[-2, 0]$	uniform
#layers L	{2, 3, 4, 5}	discrete uniform
#channels c	{8, 16, 32, 64, 128}	discrete uniform

Some experimental results

Method	Cora			Citeseer		
	nettack	metattack	random	nettack	metattack	random
CSGNN _{noAdj}	81.90	70.25	77.19	82.20	70.17	71.28
CSGNN	83.29	74.46	78.38	84.60	72.94	72.70

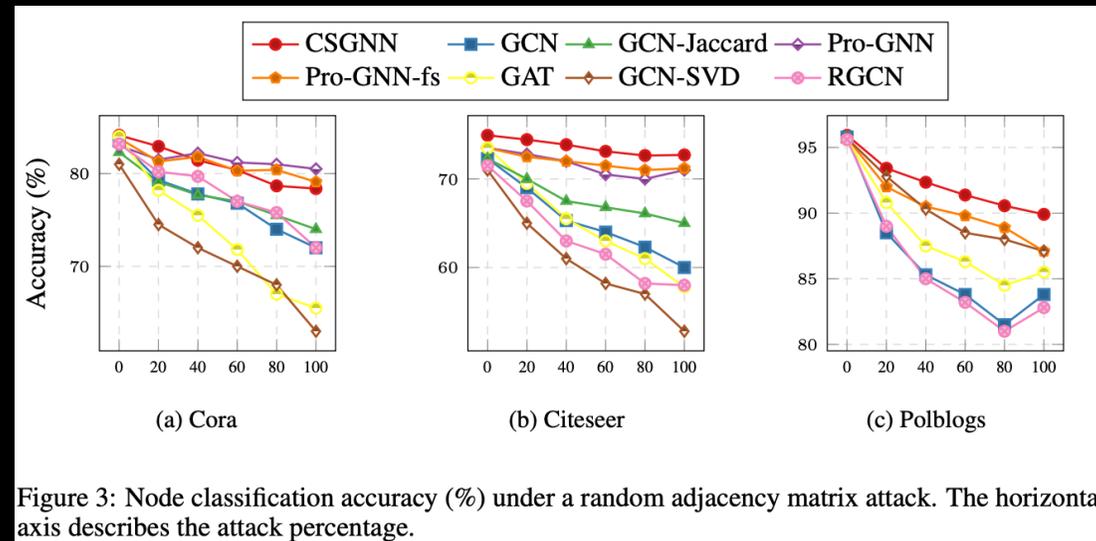
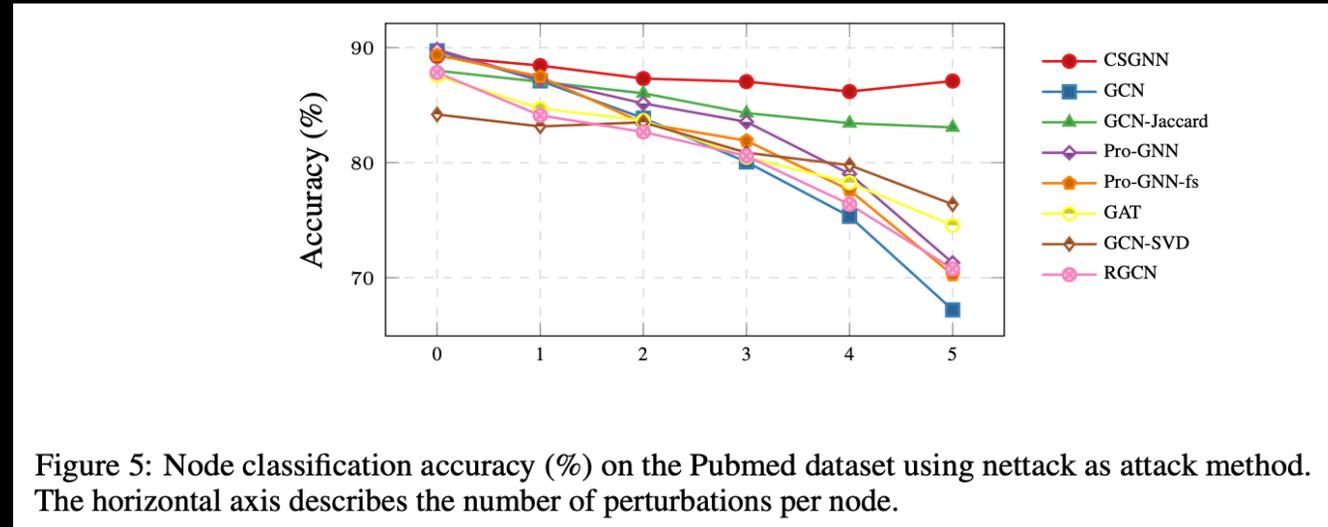
Table 6: The influence of learning the adjacency dynamical system $\Psi_{Y_i}^{h_i}$. The results show the node classification accuracy (%) with and without learning the adjacency dynamical system.



We target the nodes with degree at least 10 and flip few of their incident edges

Figure 8: Node classification accuracy (%) of ECSGNN and other baselines, under a targeted attack generated by nettack. The horizontal axis describes the number of perturbations per node.

Some experimental results



The adjacency matrix is attacked by adding random fake edges, from 0% to 100% of the number of edges in the true adjacency matrix