

# **Improving the robustness of Graph Neural Networks with coupled dynamical systems**

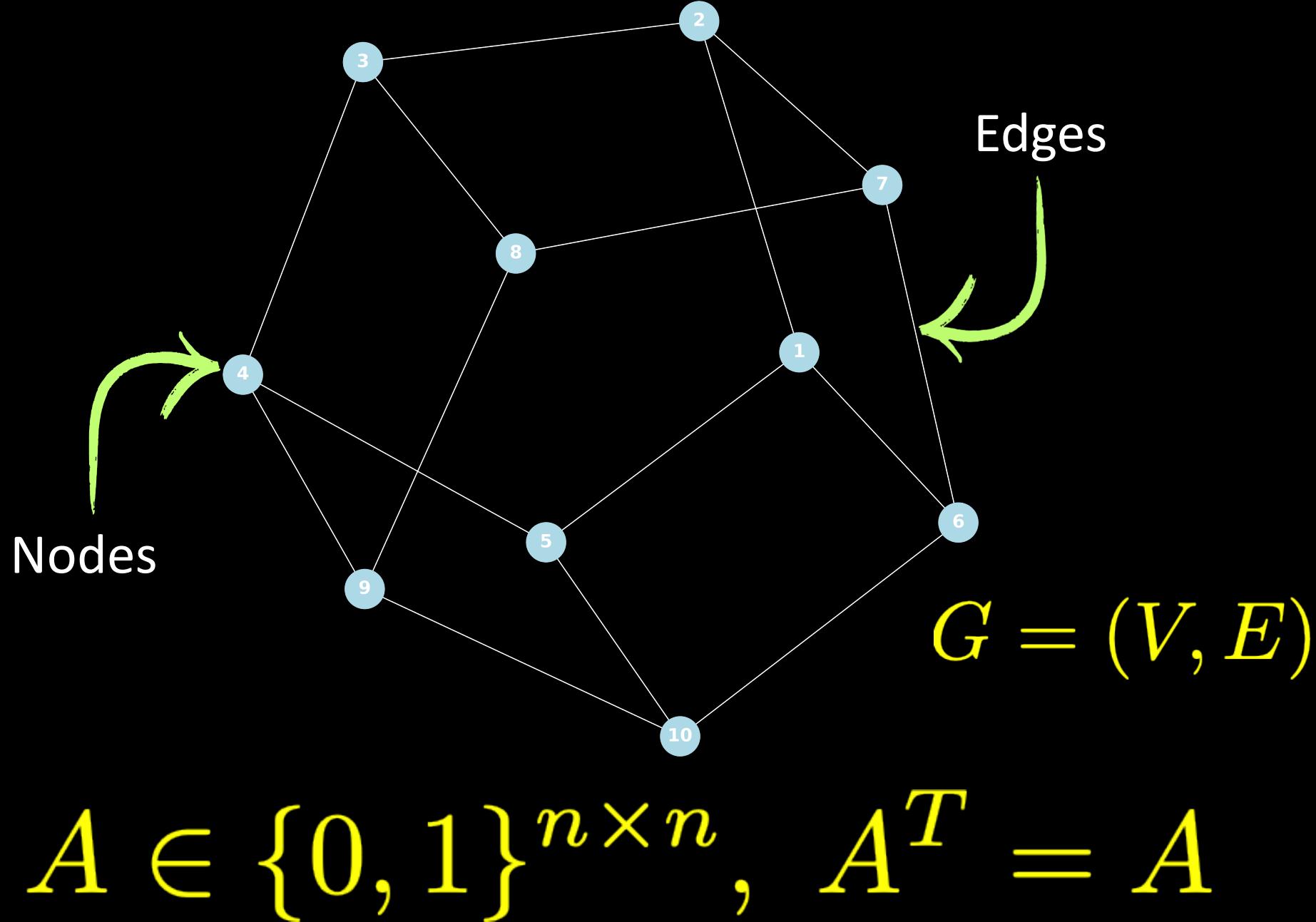
Davide Murari (NTNU)  
[davide.murari@ntnu.no](mailto:davide.murari@ntnu.no)

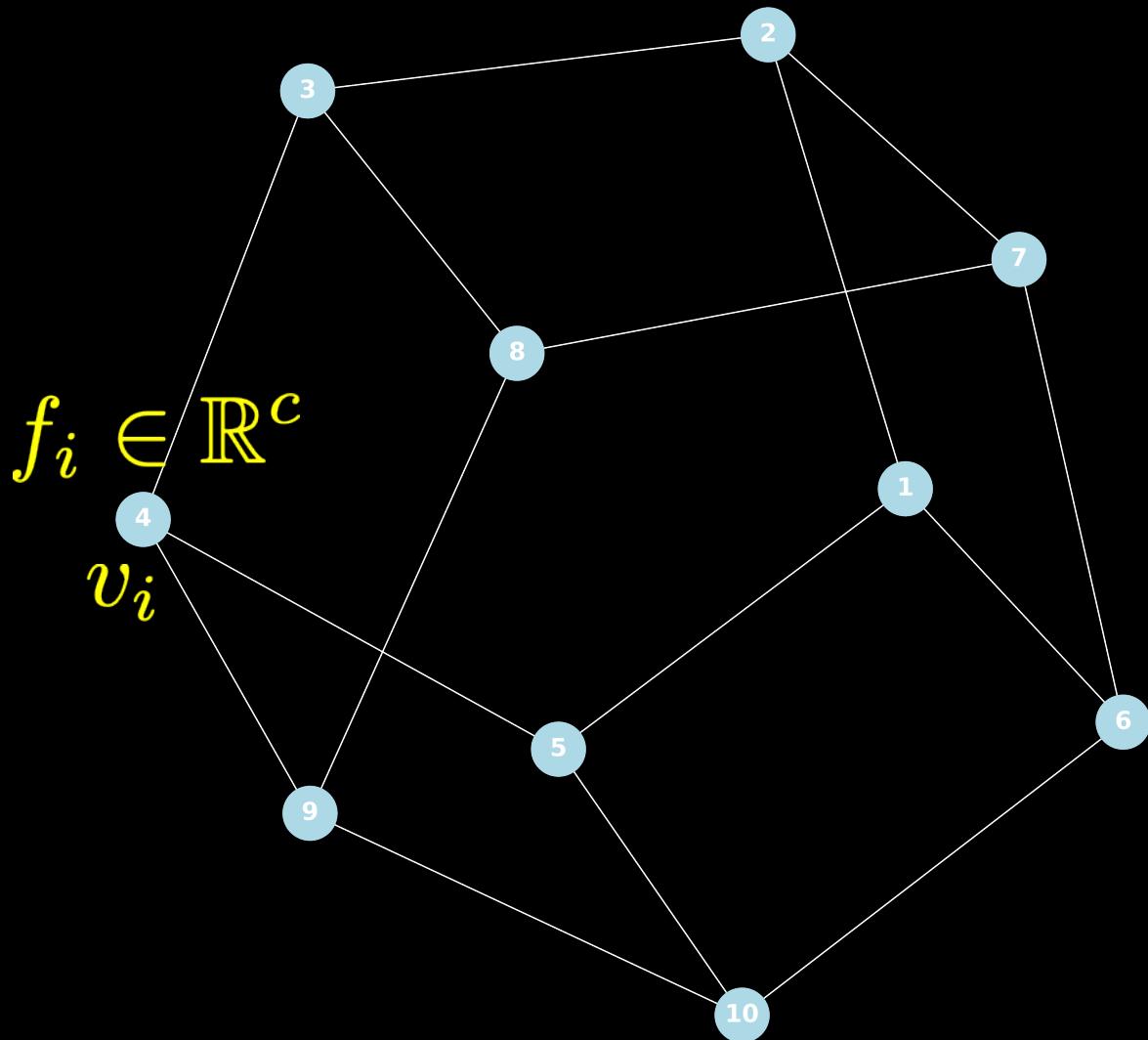
Structured Machine Learning and Time–Stepping for Dynamical Systems  
Banff, 20/02/2024

In collaboration with Moshe Eliasof, Ferdia Sherry and Carola Schönlieb

# Outline

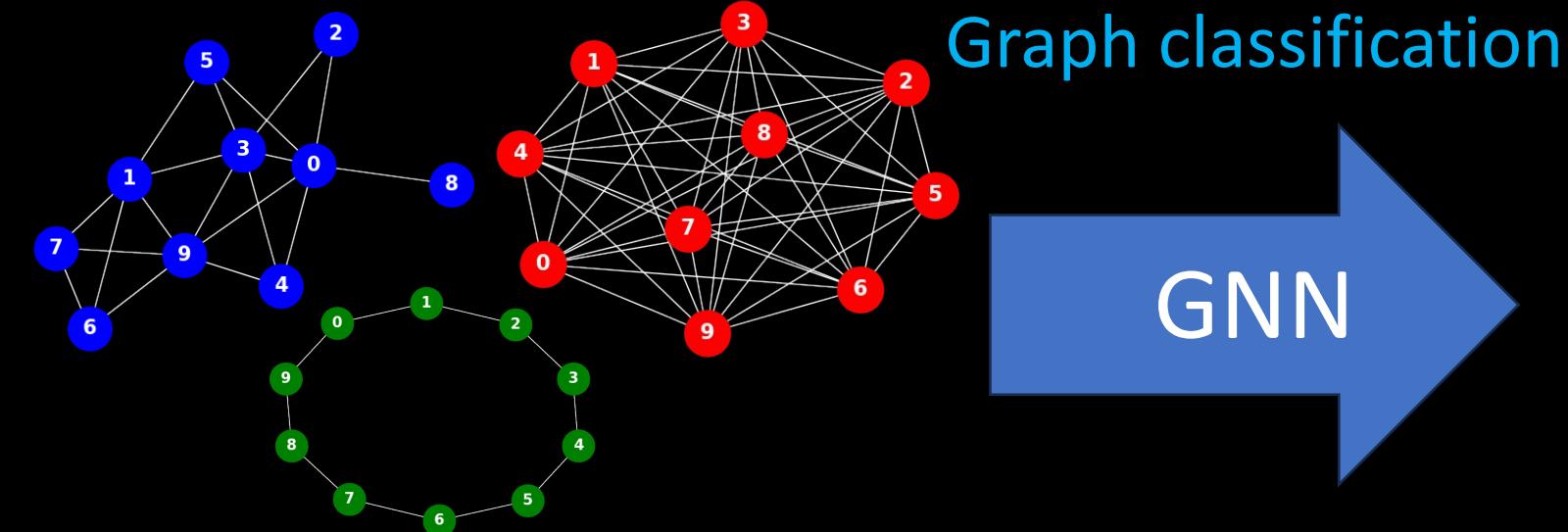
- 1 Brief introduction to graphs and Graph Neural Networks (GNNs)
- 2 The problem of adversarial robustness
- 3 Analysis of the GNN we propose
- 4 Numerical experiments



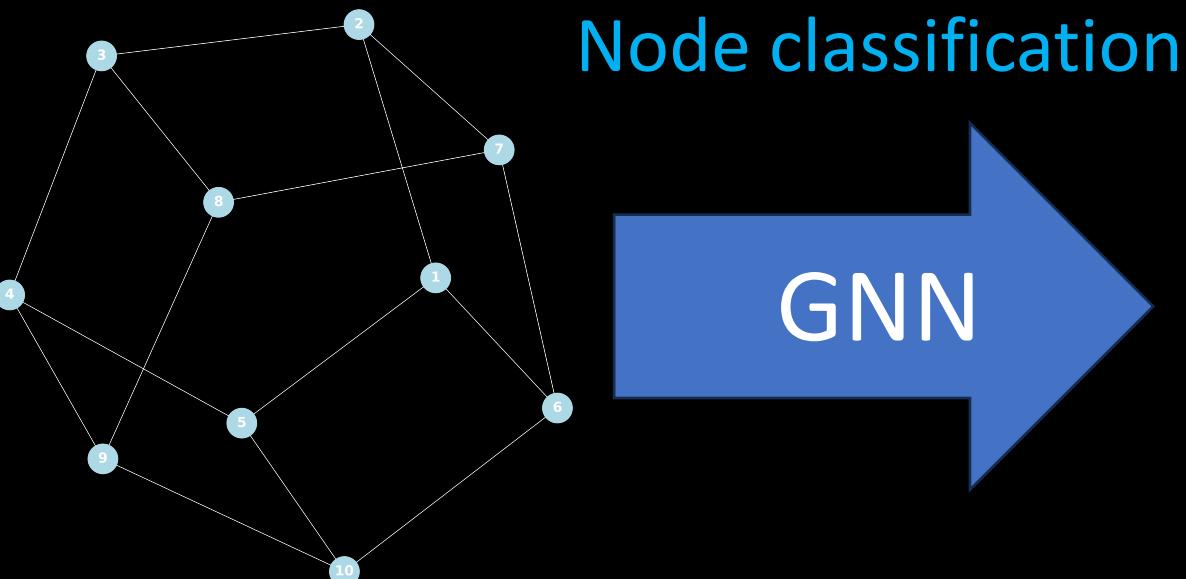


$$F \in \mathbb{R}^{n \times c}, e_i^T F = f_i$$

# Classical tasks solved with GNNs



It is a protein  
It is not a protein  
It is a protein



Red  
Blue  
Blue  
...  
Green

# Usual structure of GNNs

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

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

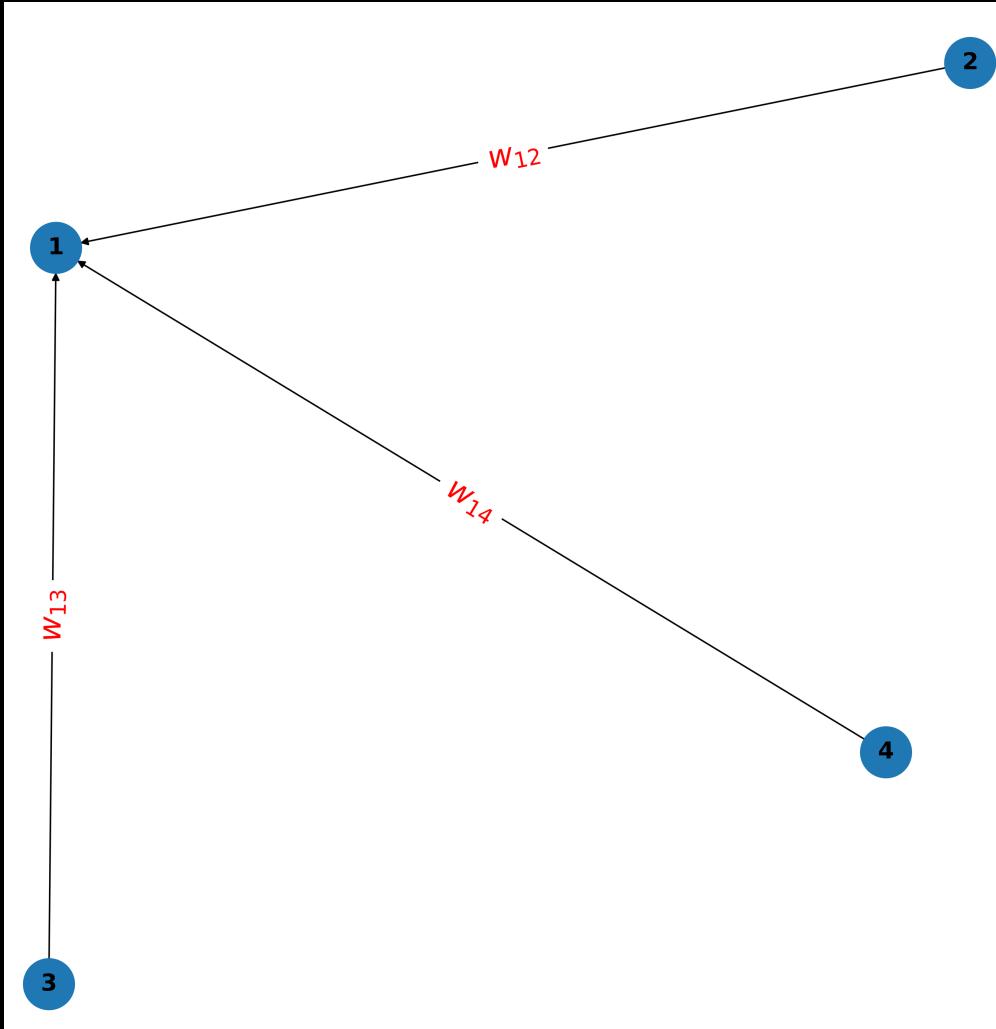
$$R = \text{MLP} \left( F^{(L)} \right) =: \text{GNN}(F, A)$$

Invariant

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

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

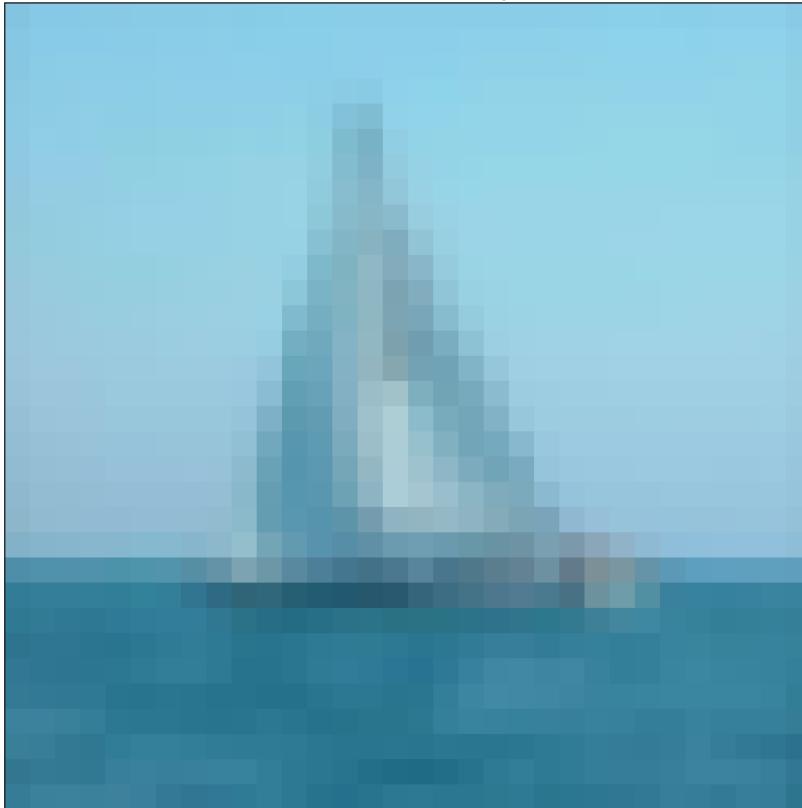
# Simple example of a GNN update



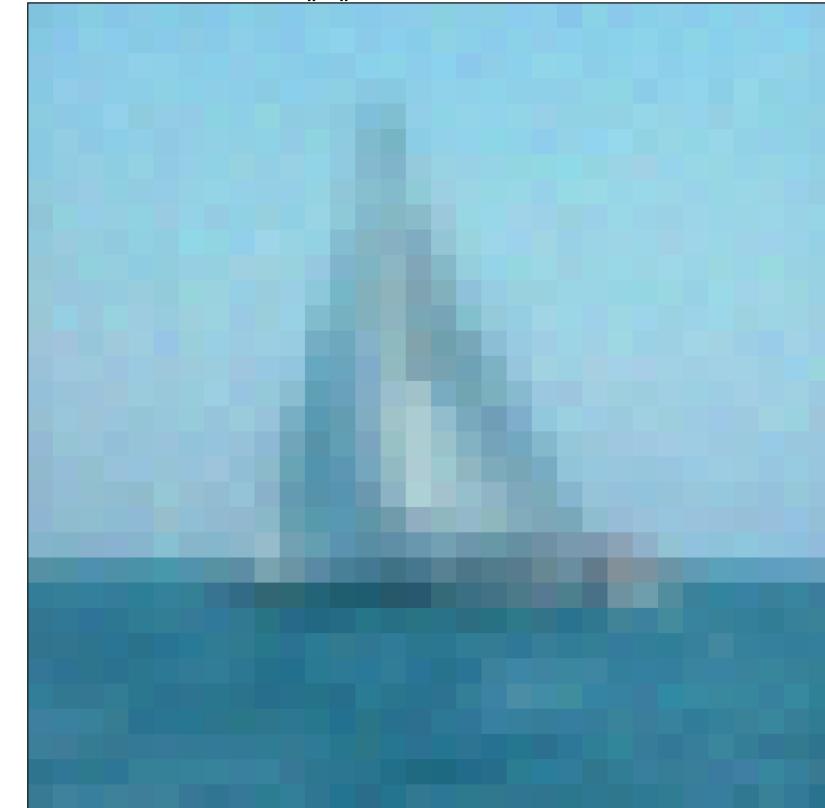
$$f_1 = \phi \left( f_1, \sum_{j=2}^4 w_{1j} \psi(f_j) \right)$$

# Adversarial Attacks for images

$X$ , Label : Ship



$X + \delta$ ,  $\|\delta\|_2 = 0.3$ , Label : Plane

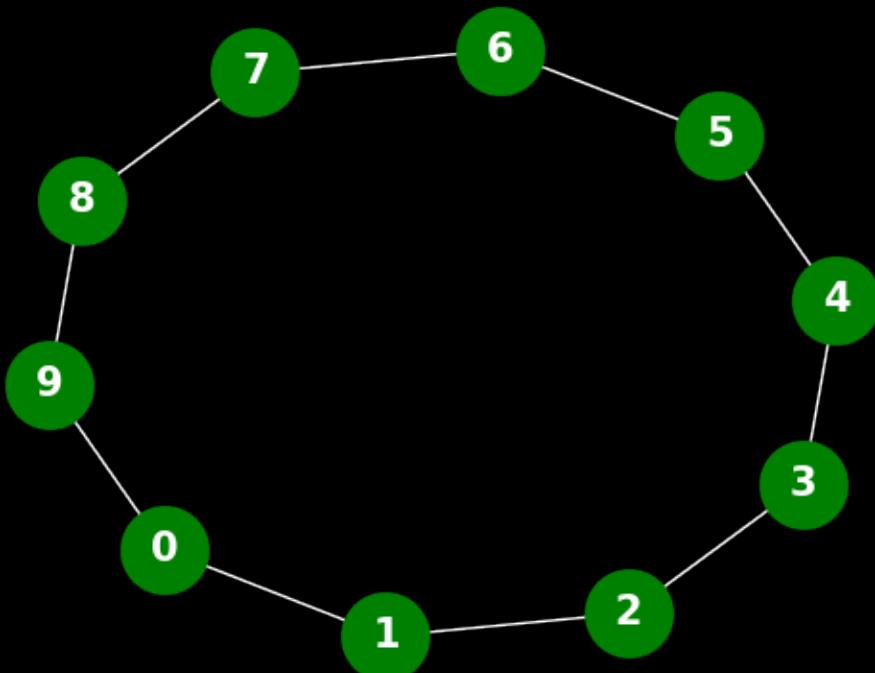


# Mathematically...

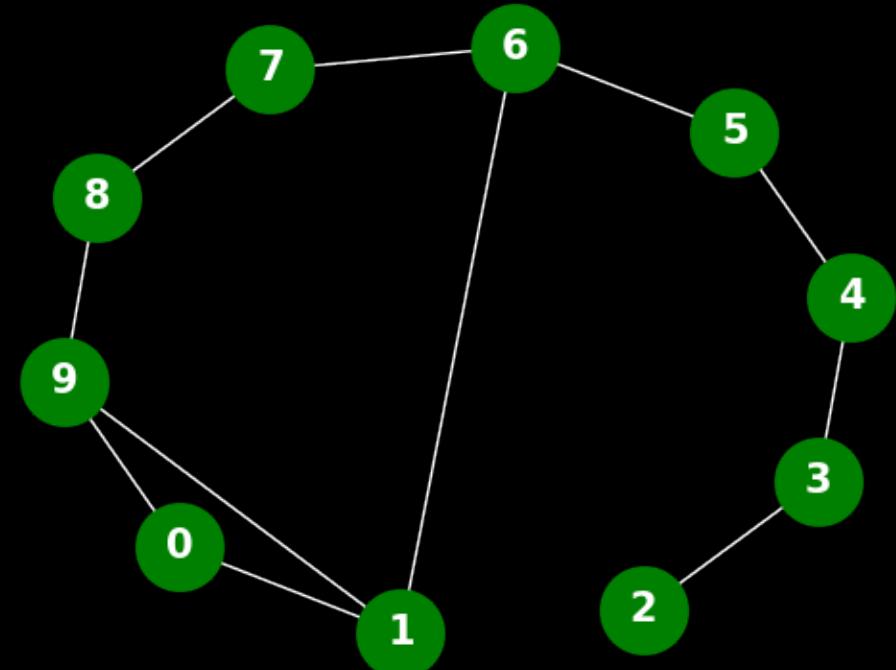
Building a network with a reduced sensitivity to these input perturbations, can be phrased as a constraint on the Lipschitz constant of the network, which should be “as small as possible” with respect to a suitable norm.

$$\|\mathcal{N}_\theta(X + \delta) - \mathcal{N}_\theta(X)\| \leq c\|\delta\|$$

# For graphs it is more complicated..



Attack to the  
connectivity  
of the Graph



e.g. A hacker adding or removing  
friendships on Facebook

# Adversarial attacks on Graphs

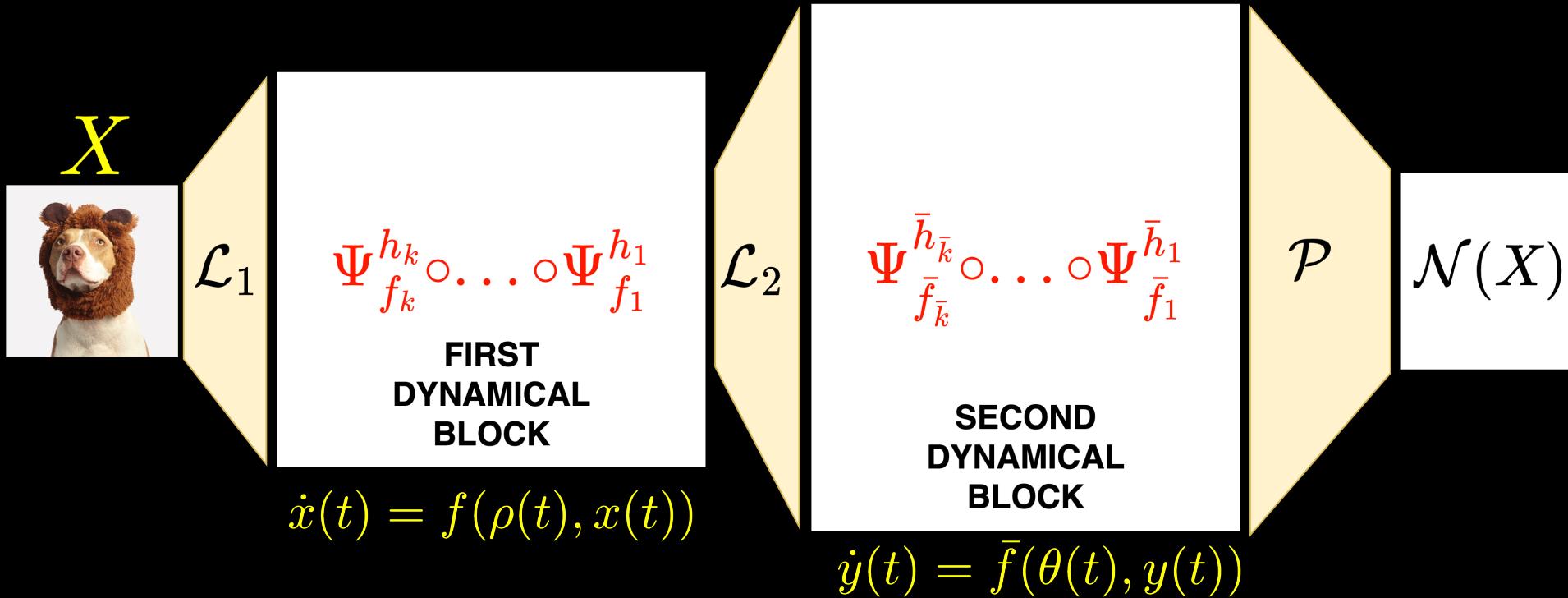
$$\begin{aligned} F_* &= F + \delta F, & \|\delta F\|_F &\leq \varepsilon_1 \\ A_* &= A + \delta A, & \sum_{i,j=1}^n |(\delta A)_{ij}| \\ && = \|\text{vec}(\delta A)\|_1 &\leq \varepsilon_2 \end{aligned}$$

Attacks do not break the properties of symmetry generally

Goal:

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

# Dynamical systems-based neural networks



# Our proposed architecture

Coupled Systems Graph Neural Network (CSGNN)

$$\begin{cases} \dot{F}(t) = -G(A(t))^T \sigma(G(A(t))F(t)W)W^T \\ \dot{A}(t) = \sigma(M(A(t))) \\ F(0) = F^{(0)}, \ A(0) = A^{(0)} \end{cases}$$

And the solution is approximated with Explicit Euler steps of “small-enough” step size to obtain the neural network.

# Some details on the functions

## Graph Gradient Operator

$$\begin{aligned}(G(A)F)_{ijk} &= A_{ij} (F_{ik} - F_{jk}), & i, j &= 1, \dots, n \\ && k &= 1, \dots, c \\ (G(A)^T O)_{ik} &= \sum_{j=1}^n (A_{ij} O_{ijk} - A_{ji} O_{jik}), & i &= 1, \dots, n, \\ && k &= 1, \dots, c\end{aligned}$$

## Linear Equivariant Map

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

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

# Non-expansivity of the system

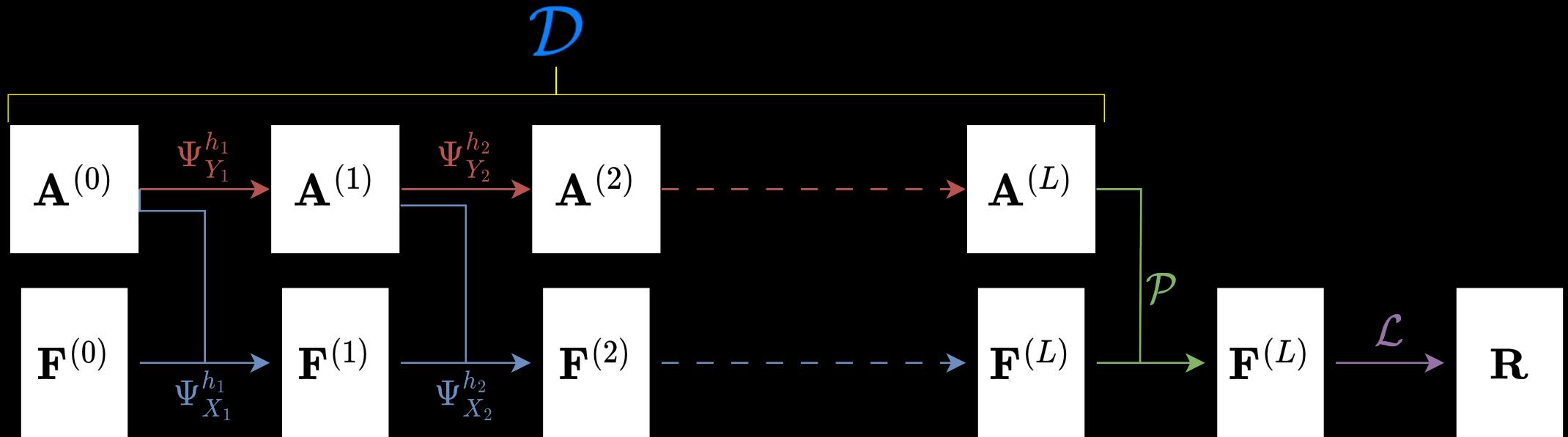
If  $\sigma(x) = \max\{ax, x\}$ ,  $a \in (0, 1)$ , then for a suitable choice of the coefficient  $k_1$ , the two individual systems are contractive, i.e.

$$\begin{aligned}\|F(t) - F_*(t)\|_F &\leq e^{-\nu_1 t} \left\| F^{(0)} - F_*^{(0)} \right\|_F, \nu_1 > 0, t \geq 0 \\ \|\text{vec}(A(t)) - \text{vec}(A_*(t))\|_1 &\leq e^{-\nu_2 t} \left\| \text{vec}(A^{(0)}) - \text{vec}(A_*^{(0)}) \right\|_1, \nu_2 > 0, t \geq 0\end{aligned}$$

and there is a pair of constants  $m_1, m_2 > 0$  such that the coupled system satisfies

$$\begin{aligned}&m_1 \|F(t) - F_*(t)\|_F + m_2 \|\text{vec}(A(t)) - \text{vec}(A_*(t))\|_1 \\ &\leq m_1 \left\| F^{(0)} - F_*^{(0)} \right\|_F + m_2 \left\| \text{vec}(A^{(0)}) - \text{vec}(A_*^{(0)}) \right\|_1\end{aligned}$$

# CSGNN



$$\left(F^{(0)}, A^{(0)}\right) := (\mathcal{K}(F_*), A_*)$$

$$\Psi_{X_i}^{h_i}(F, A) = F - h_i G(A)^T \sigma(G(A) F W_i) W_i^T$$

$$\Psi_{Y_i}^{h_i}(A) = A + h_i \sigma(M_i(A))$$

# Focus on the feature updates

If  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a non-decreasing 1-Lipschitz function, then the explicit Euler update is non-expansive in the Frobenius norm for a small enough step-size, i.e.

$$\left\| \Psi_{X_i}^{h_i}(F + \delta F, A) - \Psi_{X_i}^{h_i}(F, A) \right\|_F \leq \|\delta F\|_F,$$
$$\delta F \in \mathbb{R}^{n \times c}$$

# Focus on the adjacency updates

If  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a non-decreasing 1-Lipschitz function, then the explicit Euler update is non-expansive in the vectorized 1-norm for a small enough step-size when

$$k_1 = \left( \alpha - \sum_{i=2}^9 |k_i| \right), \quad \alpha \leq 0.$$

This means that:

$$\left\| \text{vec}(\Psi_{Y_i}^{h_i}(A + \delta A)) - \text{vec}(\Psi_{Y_i}^{h_i}(A)) \right\|_1 \leq \|\text{vec}(\delta A)\|_1,$$
$$\delta A \in \mathbb{R}^{n \times n}$$

# Robustness of the network

If the assumptions of the two previous theorems hold, and

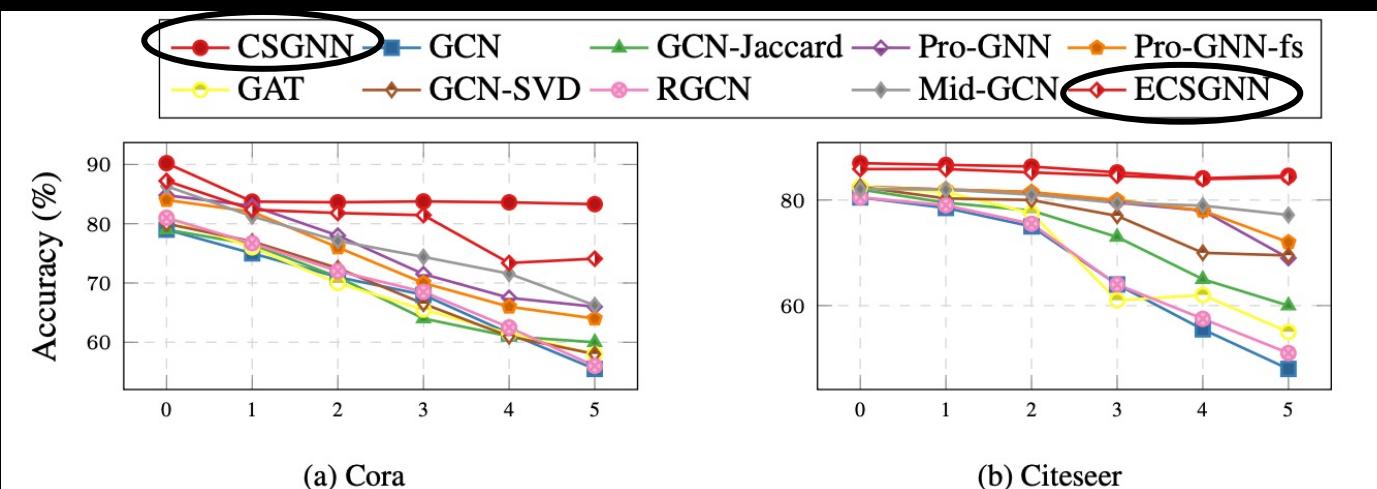
$$\begin{aligned} F_* &= F + \delta F, & \|\delta F\|_F &\leq \varepsilon_1 \\ A_* &= A + \delta A, & \|\text{vec}(\delta A)\|_1 &\leq \varepsilon_2 \end{aligned}$$

it follows

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

# Some experimental results

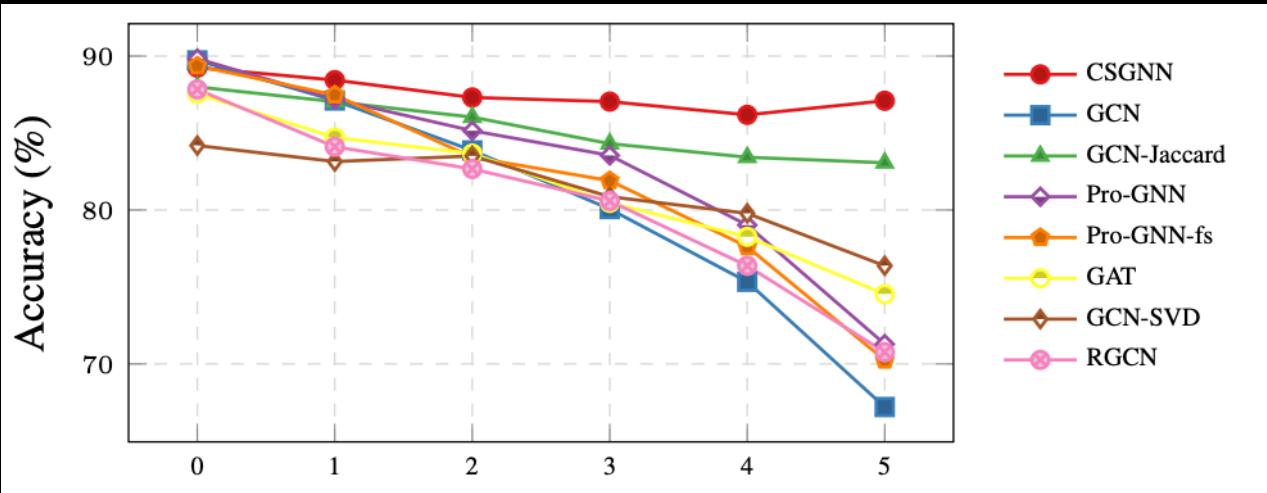
Method	Cora			Citeseer		
	nettack	metattack	random	nettack	metattack	random
CSGNN <sub>noAdj</sub>	81.90	70.25	77.19	82.20	70.17	71.28
CSGNN	83.29	74.46	78.38	84.60	72.94	72.70



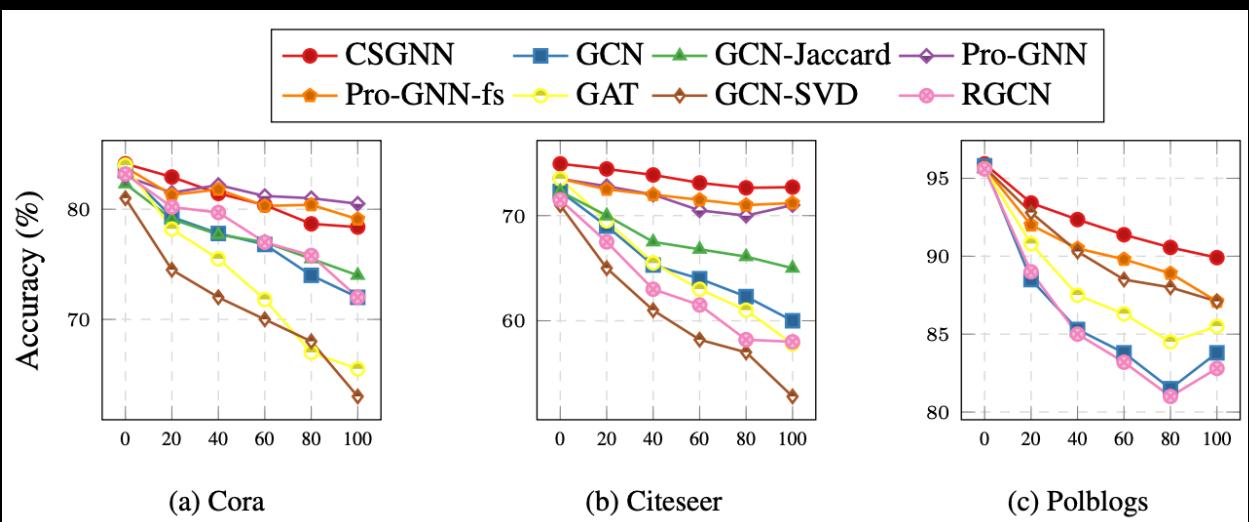
Node classification accuracy (%) of ECSGNN and other baselines, under a targeted attack generated by nettack. The horizontal axis describes the number of perturbations per node.

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

# Some experimental results



Classification accuracy for The Pubmed dataset using Nettack as attack method.



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

**Thank you for the  
attention**