

# FAST AND DEEP GRAPH NEURAL NETWORK<sup>[4]</sup>



# PROBLEM

**WHY GRAPHS:** graphs are relevant data structures that provide a useful abstraction for many kinds of real data (*molecules, social networks, transportation systems...*)

**Graph Neural Networks (GNNs)** have emerged as powerful tools for managing it [1]

from flat to  
structured domains

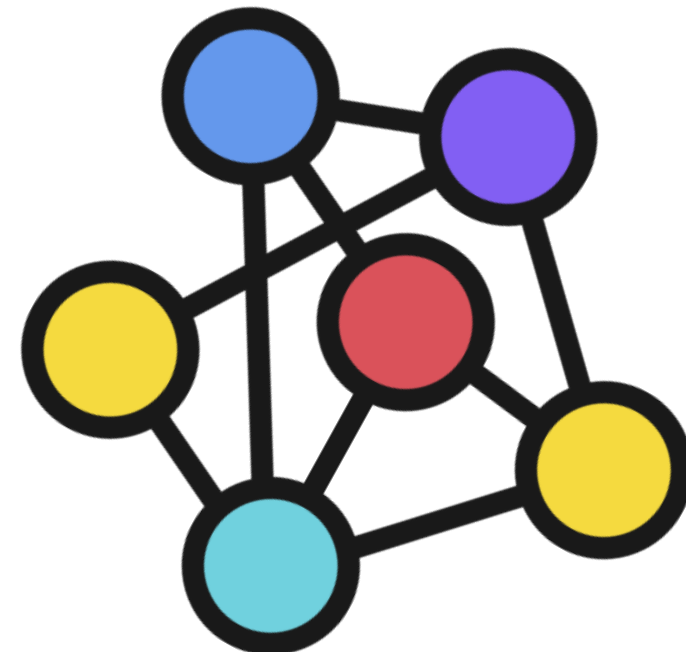
+

from shallow to  
deep architectures



richer representations but  
**high computational cost** [2]

**Solution:** Reservoir Computing [3] + Recursive Processing of Graphs [4]

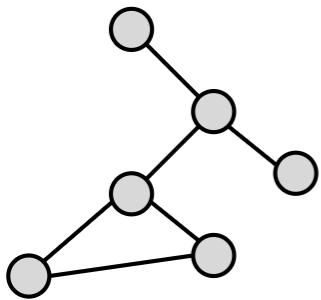


# FAST AND DEEP GNN

It combines:

1. The capability of stable dynamic systems for the **graph embedding**
2. The potentiality of **deep organization** of the GNN architecture
3. The extreme efficiency of a randomized neural network

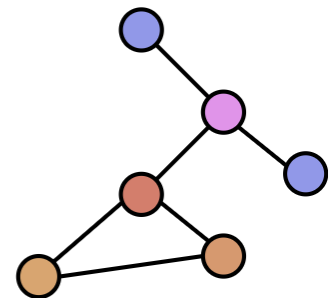
**core idea** is to exploit the fixed point of the recursive/dynamical system to embed the input graphs



input graph



**Efficiency** comes both from  
**sparsity** and **weight randomness**



fixed point embedding

iterative state computation

# PROPOSED METHOD

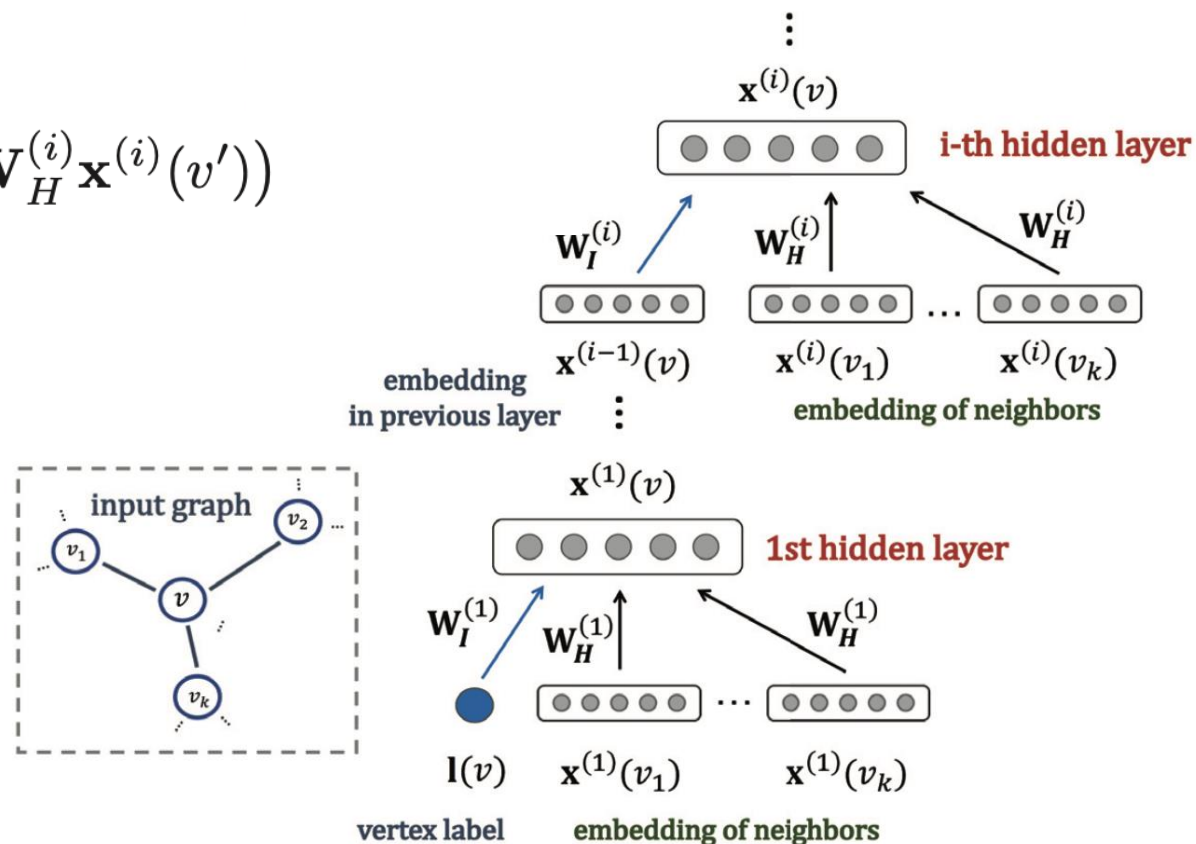
**Problem:** Given a **set graphs**  $G = (V_G, E_G)$  and its adjacency matrix  $A_G \in R^{N_G \times N_G}$  with values in  $[0,1]$  we want to **classify** them

**FDGNN [5]:** 
$$\mathbf{x}^{(i)}(v) = \tanh(\mathbf{W}_I^{(i)} \mathbf{u}^{(i)}(v) + \sum_{v' \in \mathcal{N}(v)} \mathbf{W}_H^{(i)} \mathbf{x}^{(i)}(v'))$$

Where:

- $i$  refers to the  $i$ -th layer
- $\mathbf{u}^{(i)}(v)$  is the current input:
  - $\mathbf{u}^{(1)}(v) = \mathbf{l}(v)$
  - $\mathbf{u}^{(i)}(v) = \mathbf{x}^{(i-1)}(v)$

**Initialization** to 0 value for node embedding



# PROPOSED METHOD

Using  $\mathbf{U}^{(i)}$  and  $\mathbf{X}^{(i)}$  as **column-wise** collection of  $\mathbf{u}^{(i)}(v)$  and  $\mathbf{x}^{(i)}(v)$ :

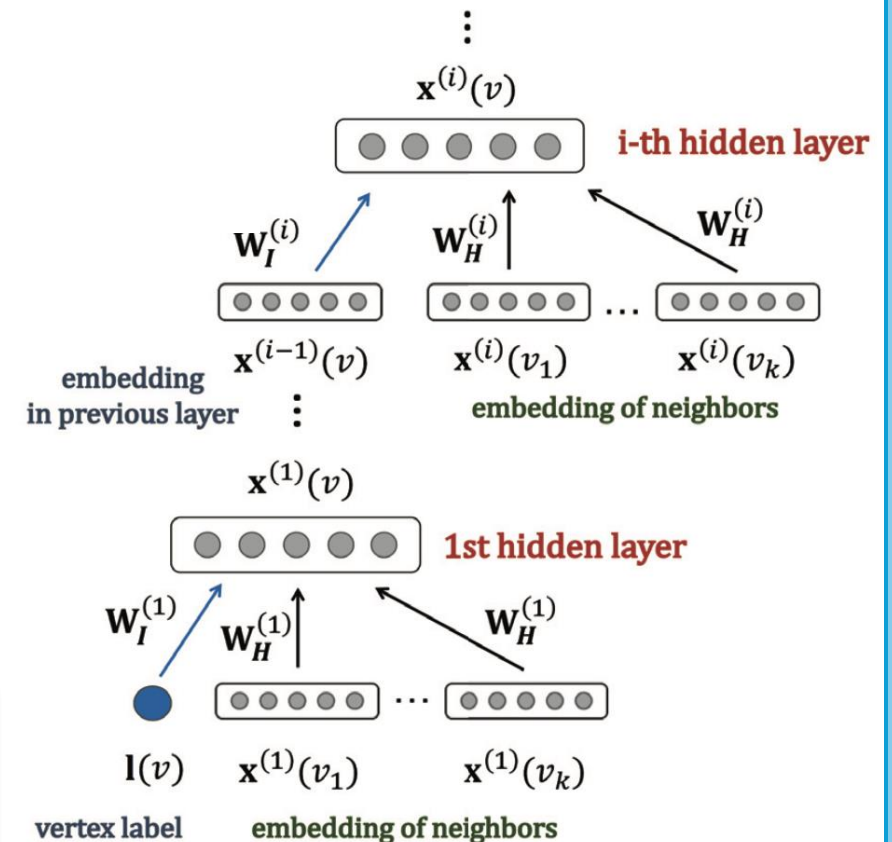
**FDGNN:**

$$\mathbf{X}^{(i)} = F^{(i)}(\mathbf{U}^{(i)}, \mathbf{X}^{(i)}) = \tanh(\mathbf{W}_I^{(i)} \mathbf{U}^{(i)} + \mathbf{W}_H^{(i)} \mathbf{X}^{(i)} \mathbf{A})$$

With  $F^{(i)}: R^{U^{(i)} \times N} \times R^{H^{(i)} \times N} \rightarrow R^{H^{(i)} \times N}$

**nota:** in case of mutual dependencies among vertexes (cycles...) this equation might not admit a unique solution

To ensure **uniqueness of neural representation**, asymptotical stability is required. Both the input  $\mathbf{U}^{(i)}$  and the adjacency matrix  $\mathbf{A}$  plays a fundamental role



# EMBEDDING PROCEDURE

---

**Algorithm 1** Layer-wise State Convergence in GESN

---

**Require:** Graph  $G = (V, E)$ , vertex labels, number of layers  $L$ , threshold  $\varepsilon$ ,  
max iterations  $\nu$

```
1:  $X^{(0)} \leftarrow$  vertex labels
2: for each layer  $i = 1$  to  $L$  do
3:    $X_0^{(i)} \leftarrow 0$ 
4:    $t \leftarrow 0$ 
5:   repeat
6:      $X_{t+1}^{(i)} \leftarrow F^{(i)}(U^{(i)}, X_t^{(i)})$ 
7:      $t \leftarrow t + 1$ 
8:   until  $\|X_t^{(i)} - X_{t-1}^{(i)}\| < \varepsilon$  or  $t \geq \nu$ 
9: end for
10: return  $X^{(L)}$ 
```

---

THE ONLY FREE  
PARAMETERS!

POOLING SUM  
FUNCTION FOR  
GRAPH-LEVEL TASK

**Output computation:**  $\mathbf{y} = \mathbf{W}_Y \tanh \left( \mathbf{W}_\phi \sum_{v \in V} \mathbf{x}^{(L)}(v) \right)$

Elements in  $\mathbf{W}_\phi$  are randomly initialized and rescaled to have unitary L2-norm

# STABILITY CONDITION

Let  $F_t$  and  $X_t$  denote the function  $F$  and the state  $X$  at the  $t$  –  $th$  iteration:

$$\mathbf{X}_t = F(\mathbf{U}, \mathbf{X}_{t-1}) = F(\mathbf{U}, F(\mathbf{U}, \mathbf{X}_{t-2})) = \dots = F(\mathbf{U}, F(\mathbf{U}, F(\dots (F(\mathbf{U}, \mathbf{X}_0)) \dots))).$$

Assuming that input and state space are **compact sets**

*note:* the focus of the analysis is on a generic layer  $i$  of the architecture so the index is dropped for the ease of notation

## Graph Embedding Stability (GES)

**Def:** For every input  $\mathbf{U}$  to the current layer, and for every  $X_0, Z_0$  initial states for the neural embeddings in the current layer, it results that:

$$\|F_t(\mathbf{U}, \mathbf{X}_0) - F_t(\mathbf{U}, \mathbf{Z}_0)\| \rightarrow 0 \quad as \quad t \rightarrow \infty.$$

# ORIGINAL CONDITION

## Sufficient condition for GES:

For every input  $\mathbf{U}$  to the current layer, if  $\|W_H\| k < 1$  then  $F$  has dynamics that satisfy the GES property.

*(too restrictive in practical applications)*

## Necessary condition for GES:

Assume that a  $k$ -regular graph with null vertices labels is an admissible input for the system. If  $F$  has dynamics that satisfy the GES property, then  $\rho(W_H) k < 1$ .

This conditions rely on **local information** ( $k$  value) but a good bound will consider **global connectivity information**. For **undirected** graphs ( $\mathbf{A}$  is symmetric) we consider:

$$\alpha^* = \rho(\mathbf{A}) = \|\mathbf{A}\|, \quad \text{because } \mathbf{A} \text{ symmetric}$$

The following conditions allows to develop a very effective dynamic, close to the edge of chaos [6]

With  $k$  the maximum among the sizes of the neighborhoods of the vertices



# SUFFICIENT CONDITION

For every input  $\mathbf{U}$  to the current layer, if  $\|\mathbf{W}_H\| \alpha^* < 1$ , where  $\alpha^*$  is the input graph spectral norm, then  $F$  has dynamics that satisfy the GES property [9].

**Proof:**

Considering: 
$$\mathbf{X}_t = \begin{cases} \tanh((\mathbf{I} \otimes \mathbf{W}_I)\mathbf{U} + (\mathbf{A} \otimes \mathbf{W}_H)\mathbf{X}_{t-1}) & t > 0 \\ \mathbf{X}_0 & t = 0 \end{cases}$$
 Where: 
$$A \otimes B = \begin{bmatrix} a_{11}B & \cdots & a_{1n}B \\ \vdots & \ddots & \vdots \\ a_{m1}B & \cdots & a_{mn}B \end{bmatrix}$$

Considering  $\mathbf{I} \otimes \mathbf{W}_I$  as **input to reservoir weights** and  $\mathbf{A} \otimes \mathbf{W}_H$  as **reservoir recurrent units**, given the **Lipschitz continuity** of  $F$ , consider the evolution of the state difference at **time**  $T$  when starting from  $\mathbf{x}_0$  and  $\mathbf{x}'_0$ :

$$\begin{aligned} \|F_T(\mathbf{u}, \mathbf{x}_0) - F_T(\mathbf{u}, \mathbf{x}'_0)\| &= \|F(\mathbf{u}, \mathbf{x}_{T-1}) - F(\mathbf{u}, \mathbf{x}'_{T-1})\| \\ &\leq \|\mathbf{A} \otimes \mathbf{W}_H\| \|\mathbf{x}_{T-1} - \mathbf{x}'_{T-1}\| \\ &\dots \\ &\leq \|\mathbf{A} \otimes \mathbf{W}_H\|^T \|\mathbf{x}_0 - \mathbf{x}'_0\| \\ &= (\alpha^* \|\mathbf{W}_H\|)^T \|\mathbf{x}_0 - \mathbf{x}'_0\|. \end{aligned}$$

**A perturbation on the initial state thus propagates through iterations as  $(\alpha^* \|\mathbf{W}_H\|)^T$**

# NECESSARY CONDITION

If  $F$  has dynamics that satisfy the GES property under null input  $\mathbf{u} = \mathbf{0}$ , then  $\rho(\mathbf{W}_H) \alpha^* < 1$ , where  $\alpha^*$  is the input graph spectral radius[9]

**Proof:**

Consider the linearised version of 
$$\mathbf{X}_t = \begin{cases} \tanh((\mathbf{I} \otimes \mathbf{W}_I)\mathbf{U} + (\mathbf{A} \otimes \mathbf{W}_H)\mathbf{X}_{t-1}) & t > 0 \\ \mathbf{X}_0 & t = 0 \end{cases},$$

then around the zero state for null input,  $\tilde{\mathbf{x}} = (\mathbf{A} \otimes \mathbf{W}_H)\tilde{\mathbf{x}}$ . If the condition  $\rho(\mathbf{A} \otimes \mathbf{W}_H) = \alpha^* \rho(\mathbf{W}_H) < 1$  is violated, then the system is unstable around the zero state, and therefore the GES property is not satisfied

**Remark 1:** In both sufficient and necessary condition  $\|\mathbf{A}\| = \rho(\mathbf{A}) = \alpha^*$ , assuming that the input graph is undirected. Both can be extended to **directed graphs** as  $\|\mathbf{W}_H\| < 1/\|\mathbf{A}\|$  and  $\rho(\mathbf{W}_H) < 1/\rho(\mathbf{A})$ , respectively

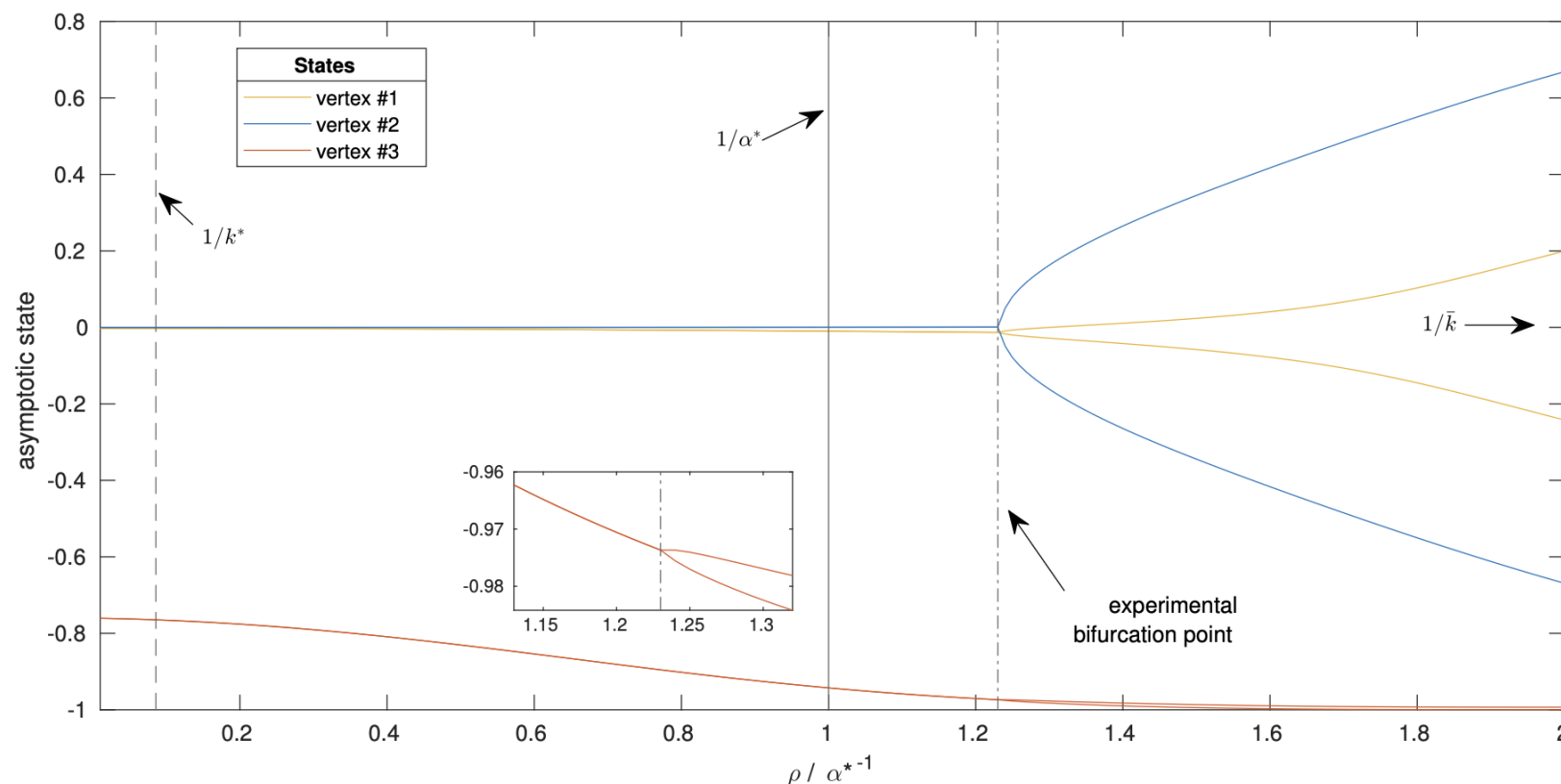
**Remark 2:** to deal with **weighted graph**, its enough to include the graph laplatian  $(\mathbf{D} - \mathbf{A})$  in the equation, keeping valid GES conditions

With  $k$  the maximum among the sizes of the neighborhoods of the vertices

# BIFURCATION DIAGRAM<sup>[9]</sup>

## Setup:

- 1000 state transition function iterations
- single-unit reservoir
- different values of  $\rho$
- mono-dimensional random vertex input feature



*asymptotic state values for three vertices of the Cora graph, which represents the citation network of 2708 scientific publications [7]*

**$1/k^*$  and  $1/avg(k)$  provides good bounds for hyperparameter selection on desired spectral radius**

# WEIGHT INITIALIZATION

**Hidden-to-hidden weight initialization:** randomly initialize from u.d. between  $[-1, 1]$  and rescaled to have the desired spectral radius

Since each graph in the dataset has its own  $\alpha^*$  and vertex degrees  $k^*$ , the standard practice [5], [8] consists in using the average values across the graphs in the dataset to constrain **recurrent weights initialization**

**Input to hidden weight initialization:** randomly sample from u.d. between  $[-\omega^{(i)}, \omega^{(i)}]$ .

# COMPUTATIONAL COST

**FDGNN:**  $\mathbf{X}^{(i)} = F^{(i)}(\mathbf{U}^{(i)}, \mathbf{X}^{(i)}) = \tanh(\mathbf{W}_I^{(i)} \mathbf{U}^{(i)} + \mathbf{W}_H^{(i)} \mathbf{X}^{(i)} \mathbf{A})$

Assuming no more than  $C$  connections among neurons in hidden layers (sparsity), the cost of computing  $\mathbf{X}^{(i)}$  for each layer is:  $\mathcal{O}((C + k) H N)$

The entire process of graph embedding cost is:  $\mathcal{O}(L \nu (C + k) H N)$

*With  $\nu$  the max number of iterations*

**cost of the encoding is the same for both training and test**

Readout trained by direct method for efficiency reasons.



# EXPERIMENTS

**FDGNN** was evaluated and compared with s.o.a. model on **9 public benchmark datasets** for **graph classification**

Datasets comes from **cheminformatics** (*dataset of proteins that are classified as enzymes or non-enzymes*) and **Social Networks** (*e.g. movie collaboration datasets containing actor/actress, target is the movie genre*)

## Experimental setting:

- all the hidden layers of the architecture in the graph embedding component shared the same values of the hyper-parameters
- $H$  fixed to 50 and 500 according to the dataset
- $C = 1$  each hidden neuron has one feedforward and one recurrent connection
- $k$  is selected as the average on the graphs in the dataset (here we use  $k$  for condition stability, not  $\alpha^*$ )

# RESULTS

## Test accuracy averaged over the 10 folds of the cross-validation

	MUTAG	PTC	COX2	PROTEINS	NCI1	IMDB-b	IMDB-m	REDDIT	COLLAB
FDGNN	<b>88.51</b> $\pm 3.77$	<b>63.43</b> $\pm 5.35$	<b>83.39</b> $\pm 2.88$	<b>76.77</b> $\pm 2.86$	77.81 $\pm 1.62$	<b>72.36</b> $\pm 3.63$	<b>50.03</b> $\pm 1.25$	<b>89.48</b> $\pm 1.00$	74.44 $\pm 2.02$
FDGNN <sub>(L=1)</sub>	87.38 $\pm 6.55$	<b>63.43</b> $\pm 5.35$	82.41 $\pm 2.67$	<b>76.77</b> $\pm 2.86$	77.11 $\pm 1.52$	71.79 $\pm 3.37$	49.34 $\pm 1.70$	87.74 $\pm 1.61$	73.82 $\pm 2.32$
GNN (Uwents et al. 2011)	80.49 $\pm 0.81$	-	-	-	-	-	-	-	-
RelNN (Uwents et al. 2011)	87.77 $\pm 2.48$	-	-	-	-	-	-	-	-
DGCNN (Zhang et al. 2018)	85.83 $\pm 1.66$	58.59 $\pm 2.47$	-	75.54 $\pm 0.94$	74.44 $\pm 0.47$	70.03 $\pm 0.86$	47.83 $\pm 0.85$	-	73.76 $\pm 0.49$
PGC-DGCNN (Tran, Navarin, and Sperduti 2018)	87.22 $\pm 1.43$	61.06 $\pm 1.83$	-	76.45 $\pm 1.02$	76.13 $\pm 0.73$	71.62 $\pm 1.22$	47.25 $\pm 1.44$	-	<b>75.00</b> $\pm 0.58$
DCNN (Tran, Navarin, and Sperduti 2018)	-	-	-	61.29 $\pm 1.60$	56.61 $\pm 1.04$	-	-	-	-
PSCN (Tran, Navarin, and Sperduti 2018)	-	-	-	75.00 $\pm 2.51$	76.34 $\pm 1.68$	71.00 $\pm 2.29$	45.23 $\pm 2.84$	-	72.60 $\pm 2.15$
GK (Zhang et al. 2018)	81.39 $\pm 1.74$	55.65 $\pm 0.46$	-	71.39 $\pm 0.31$	62.49 $\pm 0.27$	65.87 $\pm 0.98$	43.89 $\pm 0.38$	77.34 $\pm 0.18$	72.84 $\pm 0.56$
DGK (Yanardag and Vishwanathan 2015)	82.66 $\pm 1.45$	57.32 $\pm 1.13$	-	71.68 $\pm 0.50$	62.48 $\pm 0.25$	66.96 $\pm 0.56$	44.55 $\pm 0.52$	78.04 $\pm 0.39$	73.09 $\pm 0.25$
RW (Zhang et al. 2018)	79.17 $\pm 2.07$	55.91 $\pm 0.32$	-	59.57 $\pm 0.09$	-	-	-	-	-
PK (Zhang et al. 2018)	76.00 $\pm 2.69$	59.50 $\pm 2.44$	81.00 $\pm 0.20$	73.68 $\pm 0.68$	82.54 $\pm 0.47$	-	-	-	-
WL (Zhang et al. 2018)	84.11 $\pm 1.91$	57.97 $\pm 2.49$	83.20 $\pm 0.20$	74.68 $\pm 0.49$	<b>84.46</b> $\pm 0.45$	-	-	-	-
KCNN (Nikolentzos et al. 2018)	-	62.94 $\pm 1.69$	-	75.76 $\pm 0.28$	77.21 $\pm 0.22$	71.45 $\pm 0.15$	47.46 $\pm 0.21$	81.85 $\pm 0.12$	74.93 $\pm 0.14$
CGMM (Bacciu, Errica, and Micheli 2018)	85.30	-	-	-	-	-	-	-	-

Typically from 3 to 5 layers are enough effective

# RESULTS

## Execution time of FDGNN on single core, without GPU

Task	Training	Test
MUTAG	0.56'' $\pm 0.33$	0.06'' $\pm 0.04$
PTC	0.16'' $\pm 0.03$	0.02'' $\pm 0.00$
COX2	1.36'' $\pm 0.42$	0.15'' $\pm 0.05$
PROTEINS	2.16'' $\pm 0.47$	0.24'' $\pm 0.04$
NCI1	2.00' $\pm 0.45$	13.36'' $\pm 3.02$
IMDB-b	7.46'' $\pm 3.14$	0.83'' $\pm 0.35$
IMDB-m	8.68'' $\pm 1.73$	0.98'' $\pm 0.22$
REDDIT	2.47' $\pm 0.01$	16.49'' $\pm 0.28$
COLLAB	22.86' $\pm 4.70$	2.54' $\pm 0.52$

## FDGNN execution time compared with s.o.a GNN models

FDGNN	GNN	GIN	WL
0.56'' $\pm 0.33$	202.28'' $\pm 166.87$	499.24'' $\pm 2.25$	1.16'' $\pm 0.03$
	361x	872x	2x



# RESULTS UNDER NEW CONDITIONS

TABLE II

HOLD-OUT ACCURACY PEAKS FOR DIFFERENT RESERVOIR RADII

Task	$1/k_*$	$1/\alpha^*$	$1/\bar{k}$
NCI1	$75.6 \pm 0.7$	$77.2 \pm 0.7$	$77.8 \pm 0.6$
IMDB-Binary	$71.3 \pm 0.7$	$69.3 \pm 1.0$	$68.7 \pm 1.2$
Reddit-Binary	$78.0 \pm 0.6$	$89.0 \pm 0.5$	$82.6 \pm 1.2$
Reddit-Multi-5K	$53.4 \pm 0.5$	$57.2 \pm 0.7$	$52.7 \pm 0.9$
Reddit-Multi-12K	$33.4 \pm 1.5$	$42.9 \pm 0.8$	$35.3 \pm 1.2$

- distribution of **NCI1** is tightly concentrated around the mean
- **Reddit datasets** exhibit a long tail in  $\alpha^*$  distribution, with their means skewed towards larger values: instability produced by larger spectral radius graphs is compensated by the high number of small spectral radius that are moving toward EOS.

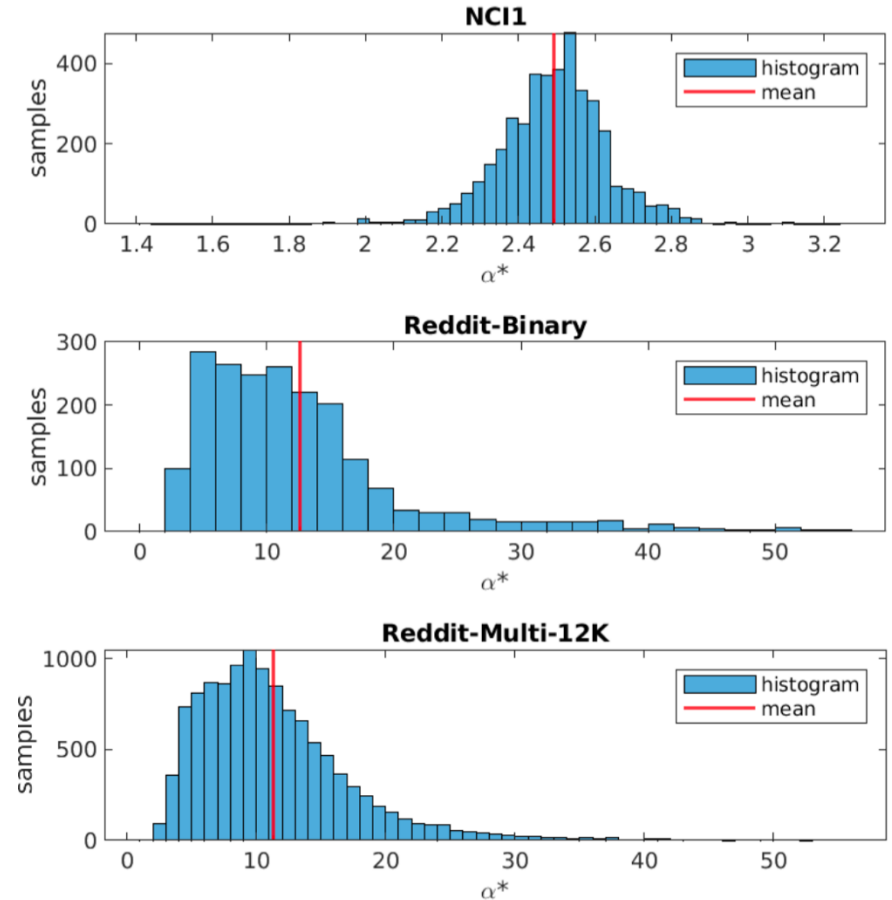
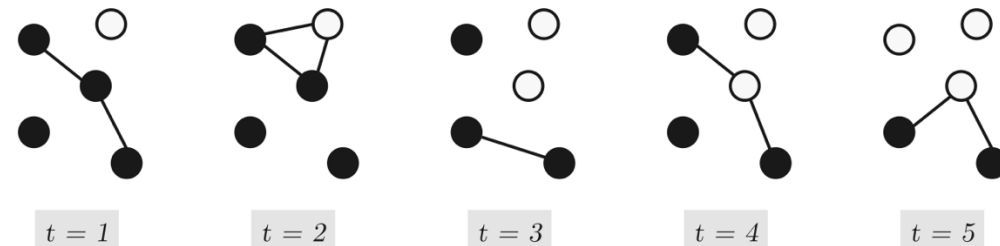


Fig. 4. Distribution of  $\alpha^*$  on three datasets.

# A FARTHER STEP... DynGESN



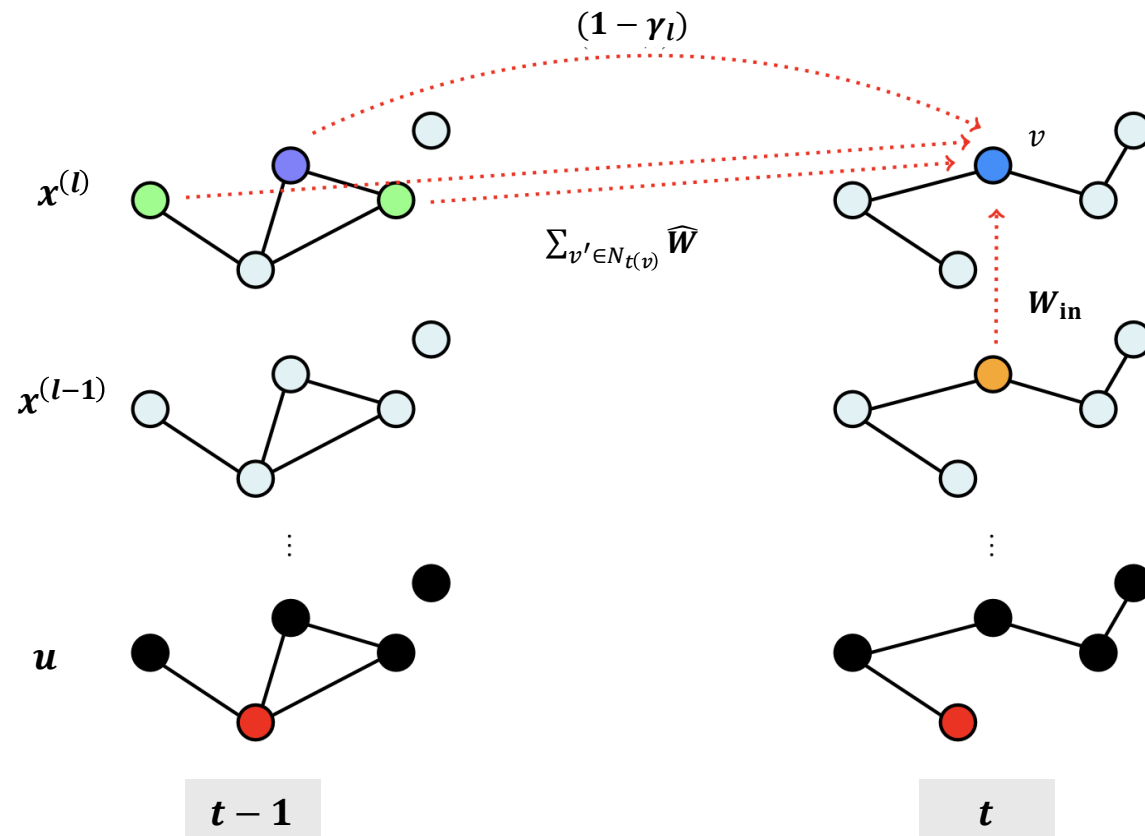
A **dynamic graph**  $G$  as a pair  $(V, E)$  with  $E = \{(v_1; v_2; t) | v_1; v_2 \in V; t \in 1 \dots T\}$  the set of edges  $v_1 \rightarrow v_2$  between a pair of vertices at a time-step  $t$ . The graph is **static** if graph topology is preserved across time, **dynamic** otherwise. A discrete-time dynamic graph can also be viewed as a **sequence of static graphs**.

**DynGESN [10] vertex-wise definition:**

$$\mathbf{x}_t^{(l)}(\nu) = \gamma_l \tanh \left( \mathbf{W}_{\text{in}}^{(l)} \mathbf{x}_t^{(\ell-1)}(\nu) + \sum_{\nu' \in N_t(\nu)} \hat{\mathbf{W}}^{(l)} \mathbf{x}_{t-1}^{(l)}(\nu') \right) + (1 - \gamma_l) \mathbf{x}_{t-1}^{(l)}(\nu)$$

Graph embedding via **pooling operation** using final final state for each layer:

$$\mathbf{X}_{\mathcal{G}} = \begin{bmatrix} \sum_{v \in \mathcal{V}} \mathbf{x}_T^{(1)}(v) \\ \vdots \\ \sum_{v \in \mathcal{V}} \mathbf{x}_T^{(L)}(v) \end{bmatrix} \in \mathcal{X}^L \subset \mathbb{R}^{HL}$$



# CONCLUSION

---

## Novelties:

- A **fast** and deep GNN model: combining RC with deep GNN architectures
- A **better bound on GES property**

## Weakness:

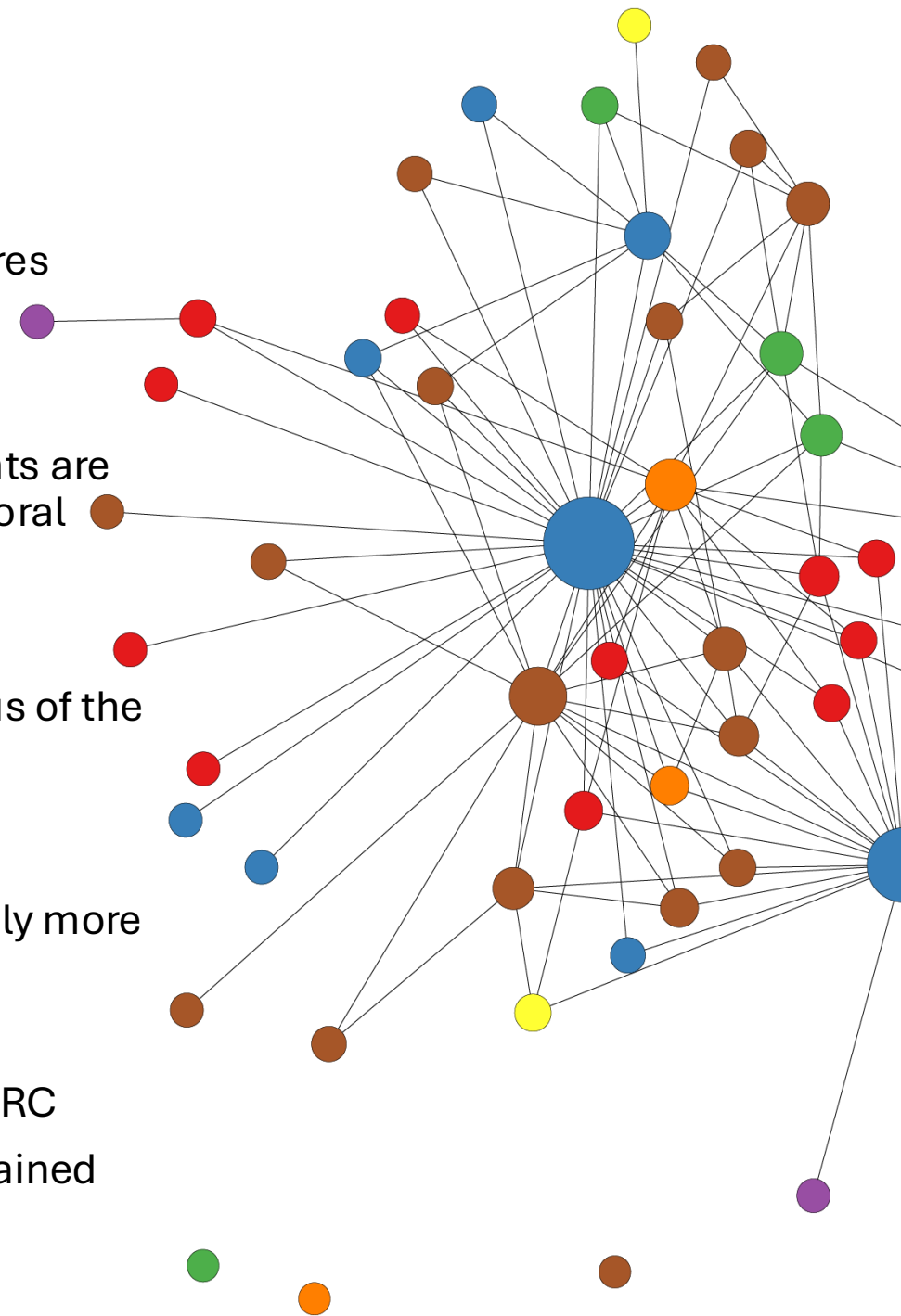
- **Limited Expressiveness** Compared to e2e GNN: since reservoir weights are fixed, Graph ESNs can underperform on tasks requiring adaptive temporal modeling or long-term dependencies

## Strengths:

- Better bounds allows to better target the selection of the spectral radius of the reservoir towards the stability limit
- **fast** and **stable** training under GES property
- through the **deep architecture** the model is able to build a progressively more effective representation of the input graphs

## Remarks:

- Model selection is the most **crucial** and **time consuming** operation in RC
- This bound can be applied in the **initialization of end-to-end** trained message passing models



# BIBLIOGRAPHY

---

- [1] D. Bacciu, F. Errica, A. Micheli, and M. Podda, “A gentle introduction to deep learning for graphs,” *Neural Networks*, vol. 129, pp. 203–221, 2020
- [2] R. Schwartz, J. Dodge, N. A. Smith, and O. Etzioni, “Green AI,” *Communications of the ACM*, vol. 63, no. 12, pp. 54–63, nov 2020
- [3] M. Lukosevicius and H. Jaeger, “Reservoir computing approaches to recurrent neural network training,” *Computer Science Review*, vol. 3, no. 3, pp. 127–149, 2009
- [4] C. Gallicchio and A. Micheli, “Graph echo state networks,” in *The 2010 International Joint Conference on Neural Networks (IJCNN)*. IEEE, July 2010, pp. 3967–3974
- [5] **C. Gallicchio and A. Micheli, “Fast and deep graph neural networks,” in *The Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI-20)*, November 2020**
- [6] J. Boedecker, O. Obst, J. T. Lizier, N. M. Mayer, and M. Asada, “Information processing in echo state networks at the edge of chaos,” *Theory in Biosciences*, vol. 131, no. 3, pp. 205–213, 2012
- [7] A. K. McCallum, K. Nigam, J. Rennie, and K. Seymore, “Automating the Construction of Internet Portals with Machine Learning,” *Information Retrieval*, vol. 3, no. 2, pp. 127–163, 2000
- [8] C. Gallicchio and A. Micheli, “Ring reservoir neural networks for graphs,” in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020
- [9] **D. Tortorella, C. Gallicchio, A. Micheli, Spectral bounds for graph echo state network stability, in: *The 2022 International Joint Conference on Neural Networks*, 2022.**
- [10] A. Micheli, D. Tortorella, Discrete-time dynamic graph echo state networks, *Neurocomputing* 496 (2022) 85–95.