# How Powerful are Graph Neural Networks?

*Keyulu Xu\*, Weihua Hu\*, Jure Leskovec, Stefanie Jegelka*

Reporter: 邵云帆 (19210240032)
yfshao19@fudan.edu.cn

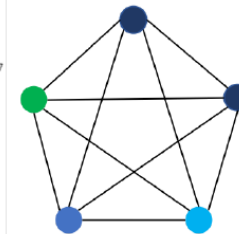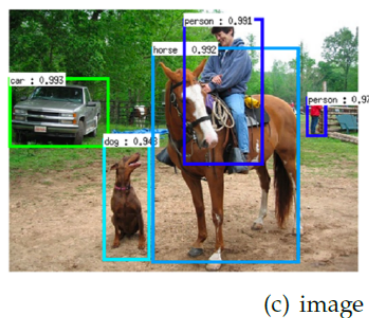https://arxiv.org/abs/1810.00826

# Introduction

- Success of deep learning on simple structured data
  - grids/images
  - sequences/text
  - more …

- Networks/graphs allow us to model complex domains
  - Social networks
  - Chemical structure
  - Biomedicine
  - more …



(a) physics

(b) molecule

(c) image

(d) text

(e) social network

(f) generation

# Applications

## Scene graph

- A structured formal graphical representation of an image.
- Encodes objects as nodes connected via pairwise relationships as edges.

# Applications

## Scene graph to image:

- **Image Generation from Scene Graphs** （CVPR 2018)
  - Generate scene layouts.
  - Use layouts generate images.

# Applications

## Image to scene graph

- **Scene Graph Generation by Iterative Message Passing** (CVPR 2018)
  - Use Region Proposal Network (RPN) generates bounding box proposals.
  - Predict node labels and edge labels.

# Applications

## VQA

- **Learning Conditioned Graph Structures for Interpretable Visual Question Answering**. (NeurIPS 2018)

- Input a question encoding, and a set of object bounding boxes.

- Obtain graph representations to learn image features.

# Graph Representation Learning

- **Input**: Graph with node features

- **Task**: Node/Graph classification

- **Goal**: Learning node/graph embeddings that capture structure information



Embedding space: $R^d$

# Graph Neural Networks (GNNs)

**Input**:

Graph $G = (V, E)$ with Node features $X_v$

**GNN**: Learning graph features using graph structure and node features.

- Aggregate

$$a_v^{(k)} = \text{AGGREGATE}^{(k)} \left( \left\{ h_u^{(k-1)} : u \in \mathcal{N}(v) \right\} \right)$$

- Combine

$$h_v^{(k)} = \text{COMBINE}^{(k)} \left( h_v^{(k-1)}, a_v^{(k)} \right)$$

- Readout

$$h_G = \text{READOUT}\left( \left\{ h_v^{(K)} \mid v \in G \right\} \right)$$

# Graph Neural Networks (GNNs)

Different variants of graph neural networks.

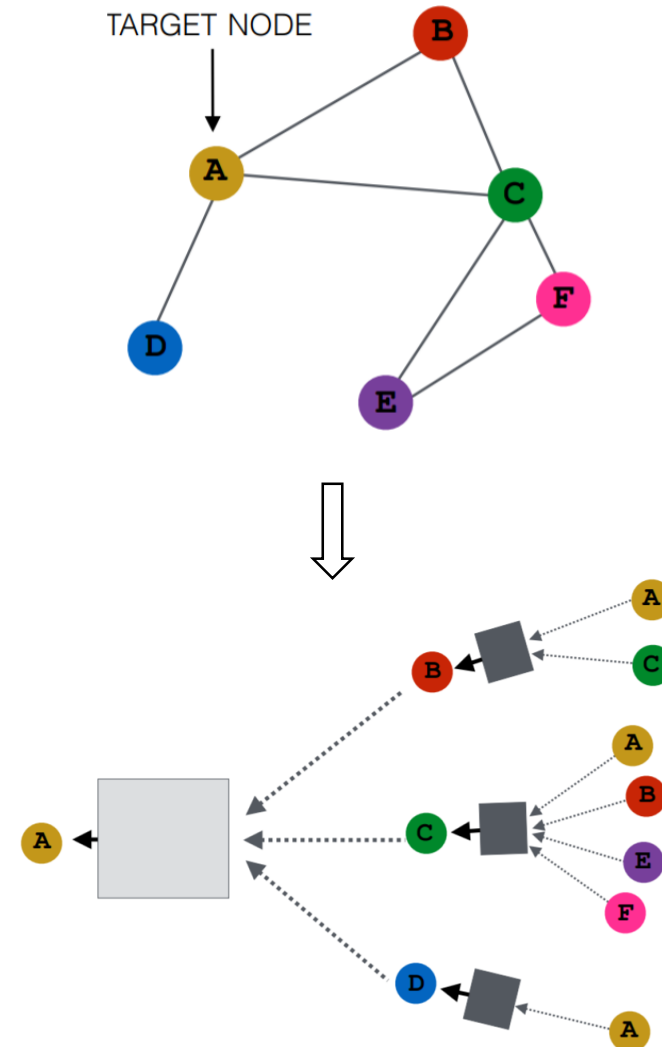| Name | Variant | Aggregator | Updater |
|---|---|---|---|
| Spectral Methods | ChebNet | $\mathbf{N}_k = \mathbf{T}_k(\tilde{\mathbf{L}})\mathbf{X}$ | $\mathbf{H} = \sum_{k=0}^{K} \mathbf{N}_k \mathbf{\Theta}_k$ |
| | $1^{st}$-order model | $\mathbf{N}_0 = \mathbf{X}$ <br> $\mathbf{N}_1 = \mathbf{D}^{-\frac{1}{2}}\mathbf{A}\mathbf{D}^{-\frac{1}{2}}\mathbf{X}$ | $\mathbf{H} = \mathbf{N}_0\mathbf{\Theta}_0 + \mathbf{N}_1\mathbf{\Theta}_1$ |
| | Single parameter | $\mathbf{N} = (\mathbf{I}_N + \mathbf{D}^{-\frac{1}{2}}\mathbf{A}\mathbf{D}^{-\frac{1}{2}})\mathbf{X}$ | $\mathbf{H} = \mathbf{N}\mathbf{\Theta}$ |
| | GCN | $\mathbf{N} = \tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}\mathbf{X}$ | $\mathbf{H} = \mathbf{N}\mathbf{\Theta}$ |
| Non-spectral Methods | Neural FPs | $\mathbf{h}_{\mathcal{N}_v}^t = \mathbf{h}_v^{t-1} + \sum_{k=1}^{\mathcal{N}_v} \mathbf{h}_k^{t-1}$ | $\mathbf{h}_v^t = \sigma(\mathbf{h}_{\mathcal{N}_v}^t \mathbf{W}_L^{\mathcal{N}_v})$ |
| | DCNN | Node classification: <br> $\mathbf{N} = \mathbf{P}^*\mathbf{X}$ <br><br> Graph classification: <br> $\mathbf{N} = 1_N^T\mathbf{P}^*\mathbf{X}/N$ | $\mathbf{H} = f\left(\mathbf{W}^c \odot \mathbf{N}\right)$ |
| | GraphSAGE | $\mathbf{h}_{\mathcal{N}_v}^t = \text{AGGREGATE}_t\left(\{\mathbf{h}_u^{t-1}, \forall u \in \mathcal{N}_v\}\right)$ | $\mathbf{h}_v^t = \sigma\left(\mathbf{W}^t \cdot [\mathbf{h}_v^{t-1}\|\mathbf{h}_{\mathcal{N}_v}^t]\right)$ |

Graph Neural Networks: A Review of Methods and Applications. J. Zhou, G. Cui, Z. Zhang, et al.

# Graph Neural Networks (GNNs)

| | | | |
|---|---|---|---|
| Graph Attention Networks | GAT | $\alpha_{vk} = \dfrac{\exp\left(\text{LeakyReLU}\left(\mathbf{a}^T[\mathbf{W}\mathbf{h}_v\|\mathbf{W}\mathbf{h}_k]\right)\right)}{\sum_{j\in\mathcal{N}_v}\exp(\text{LeakyReLU}(\mathbf{a}^T[\mathbf{W}\mathbf{h}_v\|\mathbf{W}\mathbf{h}_j]))}$ <br> $\mathbf{h}^t_{\mathcal{N}_v} = \sigma\left(\sum_{k\in\mathcal{N}_v}\alpha_{vk}\mathbf{W}\mathbf{h}_k\right)$ <br> Multi-head concatenation: <br> $\mathbf{h}^t_{\mathcal{N}_v} = \Big\|_{m=1}^{M}\sigma\left(\sum_{k\in\mathcal{N}_v}\alpha^m_{vk}\mathbf{W}^m\mathbf{h}_k\right)$ <br> Multi-head average: <br> $\mathbf{h}^t_{\mathcal{N}_v} = \sigma\left(\frac{1}{M}\sum_{m=1}^{M}\sum_{k\in\mathcal{N}_v}\alpha^m_{vk}\mathbf{W}^m\mathbf{h}_k\right)$ | $\mathbf{h}^t_v = \mathbf{h}^t_{\mathcal{N}_v}$ |
| Gated Graph Neural Networks | GGNN | $\mathbf{h}^t_{\mathcal{N}_v} = \sum_{k\in\mathcal{N}_v}\mathbf{h}^{t-1}_k + \mathbf{b}$ | $\mathbf{z}^t_v = \sigma(\mathbf{W}^z\mathbf{h}^t_{\mathcal{N}_v} + \mathbf{U}^z\mathbf{h}^{t-1}_v)$ <br> $\mathbf{r}^t_v = \sigma(\mathbf{W}^r\mathbf{h}^t_{\mathcal{N}_v} + \mathbf{U}^r\mathbf{h}^{t-1}_v)$ <br> $\widetilde{\mathbf{h}^t_v} = \tanh(\mathbf{W}\mathbf{h}^t_{\mathcal{N}_v} + \mathbf{U}(\mathbf{r}^t_v\odot\mathbf{h}^{t-1}_v))$ <br> $\mathbf{h}^t_v = (1-\mathbf{z}^t_v)\odot\mathbf{h}^{t-1}_v + \mathbf{z}^t_v\odot\widetilde{\mathbf{h}^t_v}$ |
| Graph LSTM | Tree LSTM (Child sum) | $\mathbf{h}^t_{\mathcal{N}_v} = \sum_{k\in\mathcal{N}_v}\mathbf{h}^{t-1}_k$ | $\mathbf{i}^t_v = \sigma(\mathbf{W}^i\mathbf{x}^t_v + \mathbf{U}^i\mathbf{h}^t_{\mathcal{N}_v} + \mathbf{b}^i)$ <br> $\mathbf{f}^t_{vk} = \sigma\left(\mathbf{W}^f\mathbf{x}^t_v + \mathbf{U}^f\mathbf{h}^{t-1}_k + \mathbf{b}^f\right)$ <br> $\mathbf{o}^t_v = \sigma(\mathbf{W}^o\mathbf{x}^t_v + \mathbf{U}^o\mathbf{h}^t_{\mathcal{N}_v} + \mathbf{b}^o)$ <br> $\mathbf{u}^t_v = \tanh(\mathbf{W}^u\mathbf{x}^t_v + \mathbf{U}^u\mathbf{h}^t_{\mathcal{N}_v} + \mathbf{b}^u)$ <br> $\mathbf{c}^t_v = \mathbf{i}^t_v\odot\mathbf{u}^t_v + \sum_{k\in\mathcal{N}_v}\mathbf{f}^t_{vk}\odot\mathbf{c}^{t-1}_k$ <br> $\mathbf{h}^t_v = \mathbf{o}^t_v\odot\tanh(\mathbf{c}^t_v)$ |
| | Tree LSTM (N-ary) | $\mathbf{h}^{ti}_{\mathcal{N}_v} = \sum_{l=1}^{K}\mathbf{U}^i_l\mathbf{h}^{t-1}_{vl}$ <br> $\mathbf{h}^{tf}_{\mathcal{N}_v k} = \sum_{l=1}^{K}\mathbf{U}^f_{kl}\mathbf{h}^{t-1}_{vl}$ <br> $\mathbf{h}^{to}_{\mathcal{N}_v} = \sum_{l=1}^{K}\mathbf{U}^o_l\mathbf{h}^{t-1}_{vl}$ <br> $\mathbf{h}^{tu}_{\mathcal{N}_v} = \sum_{l=1}^{K}\mathbf{U}^u_l\mathbf{h}^{t-1}_{vl}$ | $\mathbf{i}^t_v = \sigma(\mathbf{W}^i\mathbf{x}^t_v + \mathbf{h}^{ti}_{\mathcal{N}_v} + \mathbf{b}^i)$ <br> $\mathbf{f}^t_{vk} = \sigma(\mathbf{W}^f\mathbf{x}^t_v + \mathbf{h}^{tf}_{\mathcal{N}_v k} + \mathbf{b}^f)$ <br> $\mathbf{o}^t_v = \sigma(\mathbf{W}^o\mathbf{x}^t_v + \mathbf{h}^{to}_{\mathcal{N}_v} + \mathbf{b}^o)$ <br> $\mathbf{u}^t_v = \tanh(\mathbf{W}^u\mathbf{x}^t_v + \mathbf{h}^{tu}_{\mathcal{N}_v} + \mathbf{b}^u)$ <br> $\mathbf{c}^t_v = \mathbf{i}^t_v\odot\mathbf{u}^t_v + \sum_{l=1}^{K}\mathbf{f}^t_{vl}\odot\mathbf{c}^{t-1}_{vl}$ <br> $\mathbf{h}^t_v = \mathbf{o}^t_v\odot\tanh(\mathbf{c}^t_v)$ |

Graph Neural Networks: A Review of Methods and Applications. J. Zhou, G. Cui, Z. Zhang, et al.

# How powerful are GNNs

Motivation

- **Many GNN variants** are designed based on **empirical intuition**, heuristics, and experimental trial-and error.

- There is **little theoretical understanding** of the properties and limitations of GNNs.

- Formal analysis of GNNs' **representational capacity** is limited.

# How powerful are GNNs

Contribution

- Figure out the **discriminative power of GNNs**
  - Discriminative power: Map different graphs to different embedding

- Show the **maximum capacity** of GNNs theoretically

- Identify **structures that cannot be distinguished** by popular GNN variants

- Develop a **new GNN architecture** (GIN) and show its discriminative power

# WL Test

- Weisfeiler-Lehman Graph Isomorphism Test

Aggregation



Given labeled graphs G and G′

1st iteration
Result of steps 1 and 2: multiset-label determination and sorting

Combine

1st iteration
Result of step 3: label compression

1st iteration
Result of step 4: relabeling

Weisfeiler-Lehman Graph Kernels. N. Shervashidze, et al.

- Weisfeiler-Lehman Graph Isomorphism Test

**Aggregation**

Given labeled graphs G and G'

1st iteration
Result of steps 1 and 2: multiset-label determination and sorting



**Readout**

End of the 1st iteration
Feature vector representations of G and G'

$$\varphi_{WLsubtree}^{(1)}(G) = (2, 1, 1, 1, 1, 2, 0, 1, 0, 1, 1, 0, 1)$$

$$\varphi_{WLsubtree}^{(1)}(G') = (1, 2, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1)$$

Counts of original node labels     Counts of compressed node labels

$$k_{WLsubtree}^{(1)}(G,G') = \langle \varphi_{WLsubtree}^{(1)}(G), \varphi_{WLsubtree}^{(1)}(G') \rangle = 11.$$

e

**Con**

1st it
Result of step 3:

| 1,4 | → | 6 |
| 2,3 | → | 7 |
| 2,35 | → | 8 |
| 2,45 | → | 9 |

c

ration
4: relabeling

d

[Weisfeiler-Lehman Graph Kernels](). N. Shervashidze, et al.

# WL Test

- Weisfeiler-Lehman Graph Isomorphism Test

Aggregation



Readout

Given labeled graphs G and G'

1st iteration
Result of steps 1 and 2: multiset-label determination and sorting

End of the 1st iteration
Feature vector representations of G and G'

$$\varphi^{(1)}_{WLsubtree}(G) = (2, 1, 1, 1, 1, 2, 0, 1, 0, 1, 1, 0, 1)$$

$$\varphi^{(1)}_{WLsubtree}(G') = (1, 2, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1)$$

Counts of original node labels

Counts of compressed node labels

$$k^{(1)}_{WLsubtree}(G,G') = <\varphi^{(1)}_{WLsubtree}(G), \varphi^{(1)}_{WLsubtree}(G')> = 11.$$

WL Test and GNNs are similar!

1st it
Result of step 3:

| | | |
|---|---|---|
| 1,4 | ⟶ | 6 |
| 2,3 | ⟶ | 7 |
| 2,35 | ⟶ | 8 |
| 2,45 | ⟶ | 9 |

c

ration
4: relabeling

Weisfeiler-Lehman Graph Kernels. N. Shervashidze, et al.

# Theorem Framework

- Graph nodes as Multiset
- Both WL test & GNNs capture graph structures



**An overview of our theoretical framework.** Middle panel: rooted subtree structures (at the blue node) that the WL test uses to distinguish different graphs. Right panel: if a GNN's aggregation function captures the *full multiset* of node neighbors, the GNN can capture the rooted subtrees in a recursive manner and be as powerful as the WL test.

# Theorem Framework

**Theorem 1:**

GNNs can be at most as powerful as the Weisfeiler-Lehman graph isomorphism test (a.k.a. canonical labeling or color refinement)

**Theorem 2:**

A maximally powerful GNN would never map two different neighborhoods, i.e., multisets of feature vectors, to the same representation. This means its aggregation scheme must be <span style="color:red">injective</span>.

# Theorem Framework

**Theorem 3.** *Let $\mathcal{A} : \mathcal{G} \to \mathbb{R}^d$ be a GNN. With a sufficient number of GNN layers, $\mathcal{A}$ maps any graphs $G_1$ and $G_2$ that the Weisfeiler-Lehman test of isomorphism decides as non-isomorphic, to different embeddings if the following conditions hold:*

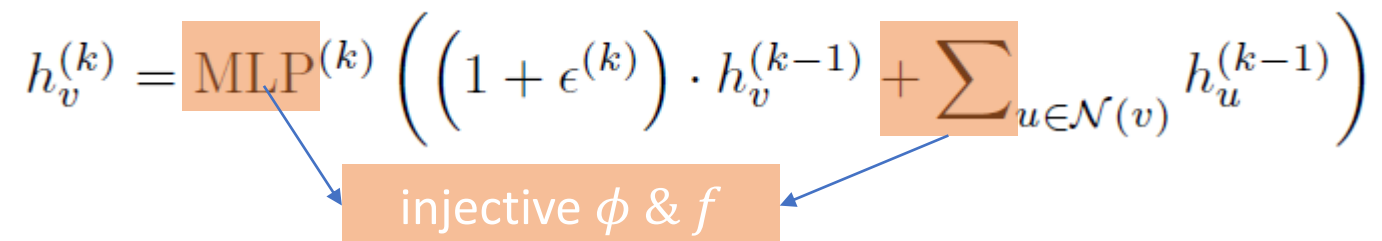a) *$\mathcal{A}$ aggregates and updates node features iteratively with*

$$h_v^{(k)} = \phi\left(h_v^{(k-1)}, f\left(\left\{h_u^{(k-1)} : u \in \mathcal{N}(v)\right\}\right)\right),$$

*where the functions $f$, which operates on multisets, and $\phi$ are* injective.

b) *$\mathcal{A}$'s graph-level readout, which operates on the multiset of node features $\left\{h_v^{(k)}\right\}$, is* injective.

Refer the paper for more theorems.

# GIN

In theory, $h_v^k = \phi\big(f\big(\{h_v^{k-1}, for\ v\ in\ V\}\big)\big)$

$$h_v^{(k)} = \mathrm{MLP}^{(k)} \left( \left( 1 + \epsilon^{(k)} \right) \cdot h_v^{(k-1)} + \sum_{u \in \mathcal{N}(v)} h_u^{(k-1)} \right)$$

injective $\phi$ & $f$

- Based on Theorems, we can:
  - Use sum pooling as **Aggregator & Readout**
  - Use MLP as **Combiner**

- Generally, there may exist many other powerful GNNs. GIN is the one being simple.
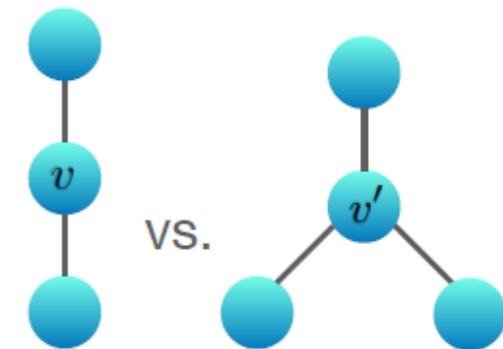
# Choice of Aggregator

- Choose between **sum, mean and max** pooling.
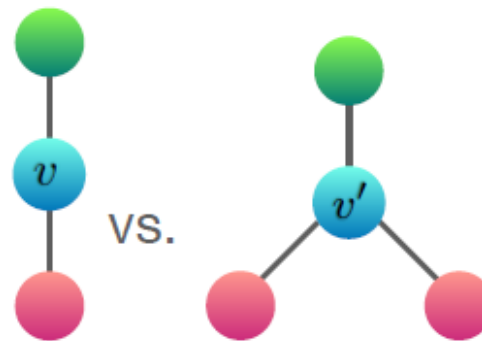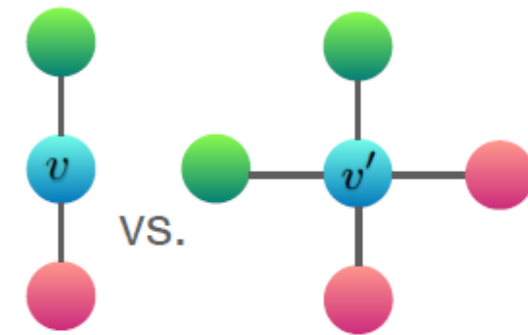
Blue as 20
Green as 10
Red as 5



(a) Mean and Max both fail        (b) Max fails        (c) Mean and Max both fail

**Examples of graph structures that mean and max aggregators fail to distinguish.** Between the two graphs, nodes $v$ and $v'$ get the same embedding even though their corresponding graph structures differ.

# Choice of Aggregator

- **Sum**: captures the full multiset

- **Mean**: captures the proportion/distribution of elements of a given type

- **Max**: reduces the multiset to a simple set



Input  sum - multiset > mean - distribution > max - set

# Choice of Combiner

- **1-layer perceptrons are not sufficient**

**Lemma** *There exist finite multisets $X_1 \neq X_2$ so that for any linear mapping $W$,* $\sum_{x \in X_1} \text{ReLU}(Wx) = \sum_{x \in X_2} \text{ReLU}(Wx).$

So the GNN layers degenerate into simply summing over neighborhood features.

- GIN **use 2-layer MLP** instead.

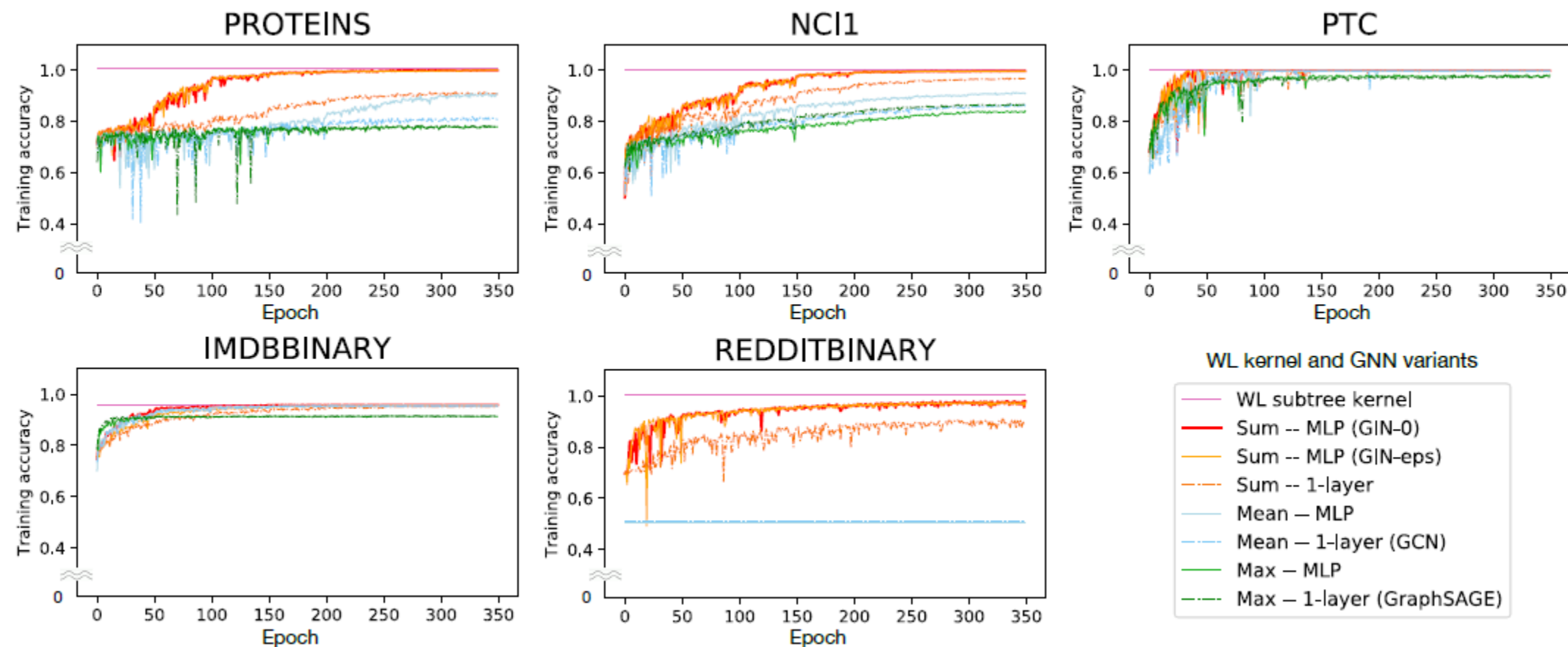# Dataset

- **9 graph classification benchmarks**
  - 4 bioinformatics datasets
    - MUTAG, PTC, NCI1, PROTEINS
  - 5 social network datasets
    - COLLAB, IMDB-BINARY, IMDB-MULTI, REDDITBINARY and REDDIT-MULTI5K

- Remove node features, force GNNs mainly learn from graph structures.

# Experiment

- Expressive power demonstrated by **training accuracy**

How Powerful are Graph Neural Networks?

# Test set performance

| Datasets | IMDB-B | IMDB-M | RDT-B | RDT-M5K | COLLAB | MUTAG | PROTEINS | PTC | NCI1 |
|---|---|---|---|---|---|---|---|---|---|
| # graphs | 1000 | 1500 | 2000 | 5000 | 5000 | 188 | 1113 | 344 | 4110 |
| # classes | 2 | 3 | 2 | 5 | 3 | 2 | 2 | 2 | 2 |
| Avg # nodes | 19.8 | 13.0 | 429.6 | 508.5 | 74.5 | 17.9 | 39.1 | 25.5 | 29.8 |
| WL subtree | $73.8 \pm 3.9$ | $50.9 \pm 3.8$ | $81.0 \pm 3.1$ | $52.5 \pm 2.1$ | $78.9 \pm 1.9$ | $90.4 \pm 5.7$ | $75.0 \pm 3.1$ | $59.9 \pm 4.3$ | $\mathbf{86.0 \pm 1.8}$ * |
| DCNN | 49.1 | 33.5 | – | – | 52.1 | 67.0 | 61.3 | 56.6 | 62.6 |
| PATCHYSAN | $71.0 \pm 2.2$ | $45.2 \pm 2.8$ | $86.3 \pm 1.6$ | $49.1 \pm 0.7$ | $72.6 \pm 2.2$ | $\mathbf{92.6 \pm 4.2}$ * | $75.9 \pm 2.8$ | $60.0 \pm 4.8$ | $78.6 \pm 1.9$ |
| DGCNN | 70.0 | 47.8 | – | – | 73.7 | 85.8 | 75.5 | 58.6 | 74.4 |
| AWL | $74.5 \pm 5.9$ | $51.5 \pm 3.6$ | $87.9 \pm 2.5$ | $54.7 \pm 2.9$ | $73.9 \pm 1.9$ | $87.9 \pm 9.8$ | – | – | – |
| SUM–MLP (GIN-0) | $\mathbf{75.1 \pm 5.1}$ | $\mathbf{52.3 \pm 2.8}$ | $\mathbf{92.4 \pm 2.5}$ | $\mathbf{57.5 \pm 1.5}$ | $\mathbf{80.2 \pm 1.9}$ | $89.4 \pm 5.6$ | $\mathbf{76.2 \pm 2.8}$ | $64.6 \pm 7.0$ | $\mathbf{82.7 \pm 1.7}$ |
| SUM–MLP (GIN-$\epsilon$) | $74.3 \pm 5.1$ | $52.1 \pm 3.6$ | $\mathbf{92.2 \pm 2.3}$ | $57.0 \pm 1.7$ | $\mathbf{80.1 \pm 1.9}$ | $89.0 \pm 6.0$ | $75.9 \pm 3.8$ | $63.7 \pm 8.2$ | $\mathbf{82.7 \pm 1.6}$ |
| SUM–1-LAYER | $74.1 \pm 5.0$ | $\mathbf{52.2 \pm 2.4}$ | $90.0 \pm 2.7$ | $55.1 \pm 1.6$ | $\mathbf{80.6 \pm 1.9}$ | $\mathbf{90.0 \pm 8.8}$ | $\mathbf{76.2 \pm 2.6}$ | $63.1 \pm 5.7$ | $82.0 \pm 1.5$ |
| MEAN–MLP | $73.7 \pm 3.7$ | $\mathbf{52.3 \pm 3.1}$ | $50.0 \pm 0.0$ | $20.0 \pm 0.0$ | $79.2 \pm 2.3$ | $83.5 \pm 6.3$ | $75.5 \pm 3.4$ | $\mathbf{66.6 \pm 6.9}$ | $80.9 \pm 1.8$ |
| MEAN–1-LAYER (GCN) | $74.0 \pm 3.4$ | $51.9 \pm 3.8$ | $50.0 \pm 0.0$ | $20.0 \pm 0.0$ | $79.0 \pm 1.8$ | $85.6 \pm 5.8$ | $76.0 \pm 3.2$ | $64.2 \pm 4.3$ | $80.2 \pm 2.0$ |
| MAX–MLP | $73.2 \pm 5.8$ | $51.1 \pm 3.6$ | – | – | – | $84.0 \pm 6.1$ | $76.0 \pm 3.2$ | $64.6 \pm 10.2$ | $77.8 \pm 1.3$ |
| MAX–1-LAYER (GraphSAGE) | $72.3 \pm 5.3$ | $50.9 \pm 2.2$ | – | – | – | $85.1 \pm 7.6$ | $75.9 \pm 3.2$ | $63.9 \pm 7.7$ | $77.7 \pm 1.5$ |

**Test set classification accuracies (%).**

# Summary

- The most powerful GNNs are **as powerful as** the WL test.

- **Powerful GNNs** have **injective** aggregation and graph readout.

- **GIN** is maximally powerful GNN. Key is to use **sum and MLP**.

# Thank You!

## Q & A