

HOW POWERFUL ARE GRAPH NEURAL NETWORKS ?

刘希阳

6-29

1 INTRODUCTION

- GNNs revolutionizing graph representation learning
- But, limited understanding of their representational properties and limitations
- Design of new GNNs is mostly based on:
 - empirical intuition
 - Heuristics
 - experimental trial-and-error

1 INTRODUCTION

- This paper
 - characterize how expressive different GNN variants are in learning to represent and distinguish between different graph structures
 - show that **GNNs** are **at most as powerful as the WL test** in distinguishing graph structures.
 - identify graph structures that cannot be distinguished by popular GNN variants, such as GCN (Kipf & Welling, 2017) and GraphSAGE (Hamilton et al., 2017a)
 - develop a simple neural architecture, Graph Isomorphism Network (GIN)

2 PRELIMINARIES

- Common Graph Neural Networks

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

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

- For Graphsage

$$a_v^{(k)} = \text{MAX} \left(\left\{ \text{ReLU} \left(W \cdot h_u^{(k-1)} \right), \forall u \in \mathcal{N}(v) \right\} \right)$$

- For GCN

$$h_v^{(k)} = \text{ReLU} \left(W \cdot \text{MEAN} \left\{ h_u^{(k-1)}, \forall u \in \mathcal{N}(v) \cup \{v\} \right\} \right)$$

2 PRELIMINARIES

- Common Prediction Task:
 - For node classification:
 - Directly using $h_v^{(k)}$
 - For graph classification:
 - Form a representation of whole graph with READOUT function
$$h_G = \text{READOUT}(\{h_v^{(K)} \mid v \in G\}).$$
 - READOUT function can be
 - Summation
 - Sophisticated graph-level pooling function
 - etc.

2 PRELIMINARIES

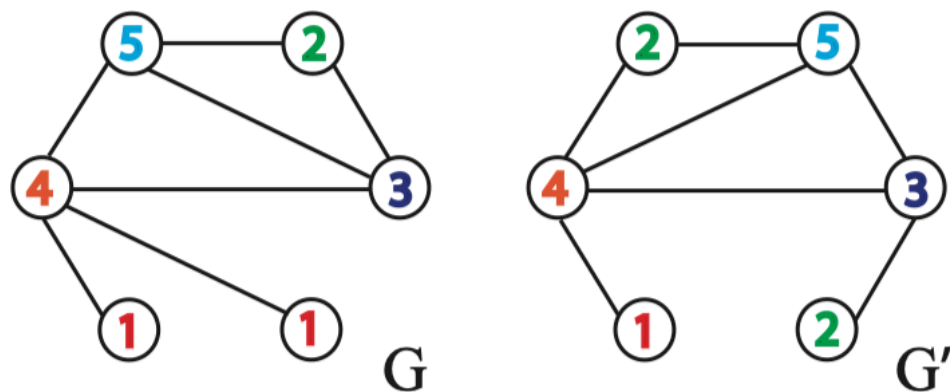
- Weisfeiler-Lehman test(1968)
 - For graph isomorphism problem(NPC problem)
 - effective and computationally efficient
- Two phases:
 - **Aggregates** the labels of nodes and their neighborhoods
 - **Hashes** the aggregated labels into unique new labels
- GCN framework is a similar to WL test

2 PRELIMINARIES

- Weisfeiler-Lehman Graph Kernels(2011)
- Based on the WL test, Shervashidze et al. proposed the WL subtree kernel that measures the similarity between graphs

WL subtree kernel

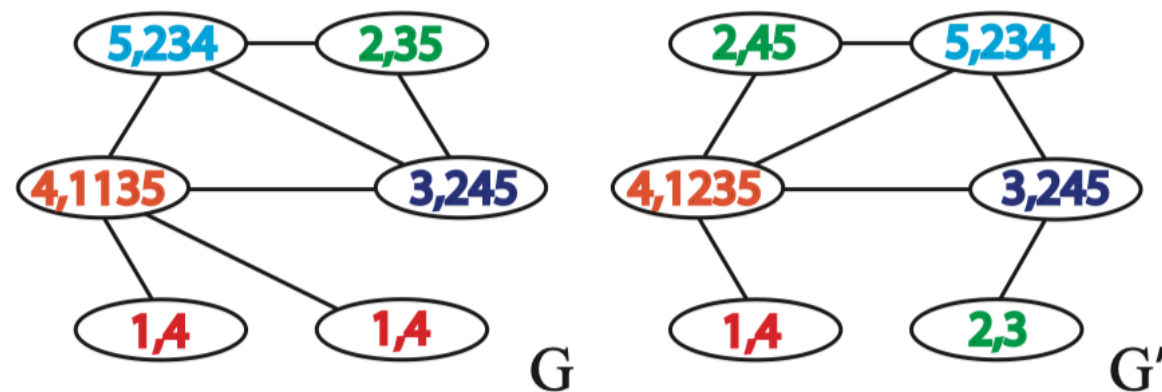
Given labeled graphs G and G'



a

1st iteration

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



b

1st iteration

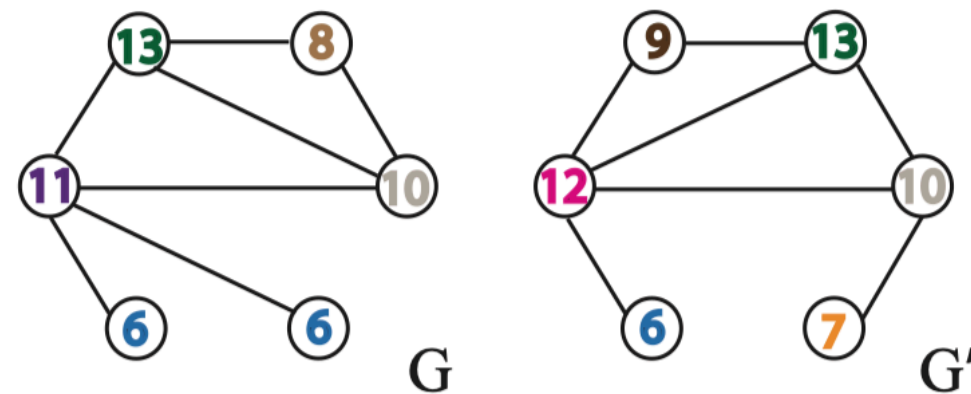
Result of step 3: label compression



c

1st iteration

Result of step 4: relabeling



d

2 PRELIMINARIES

- Weisfeiler-Lehman Graph Kernels(2011)

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

$$\varphi_{WLsubtree}^{(1)}(G) = (\textcolor{red}{2}, \textcolor{green}{1}, \textcolor{blue}{1}, \textcolor{red}{1}, \textcolor{blue}{1}, \textcolor{blue}{2}, \textcolor{orange}{0}, \textcolor{brown}{1}, \textcolor{brown}{0}, \textcolor{gray}{1}, \textcolor{purple}{1}, \textcolor{magenta}{0}, \textcolor{green}{1})$$
$$\varphi_{WLsubtree}^{(1)}(G') = (\textcolor{red}{1}, \textcolor{green}{2}, \textcolor{blue}{1}, \textcolor{red}{1}, \textcolor{blue}{1}, \textcolor{blue}{1}, \textcolor{orange}{1}, \textcolor{brown}{0}, \textcolor{brown}{1}, \textcolor{gray}{1}, \textcolor{purple}{0}, \textcolor{magenta}{1}, \textcolor{green}{1})$$

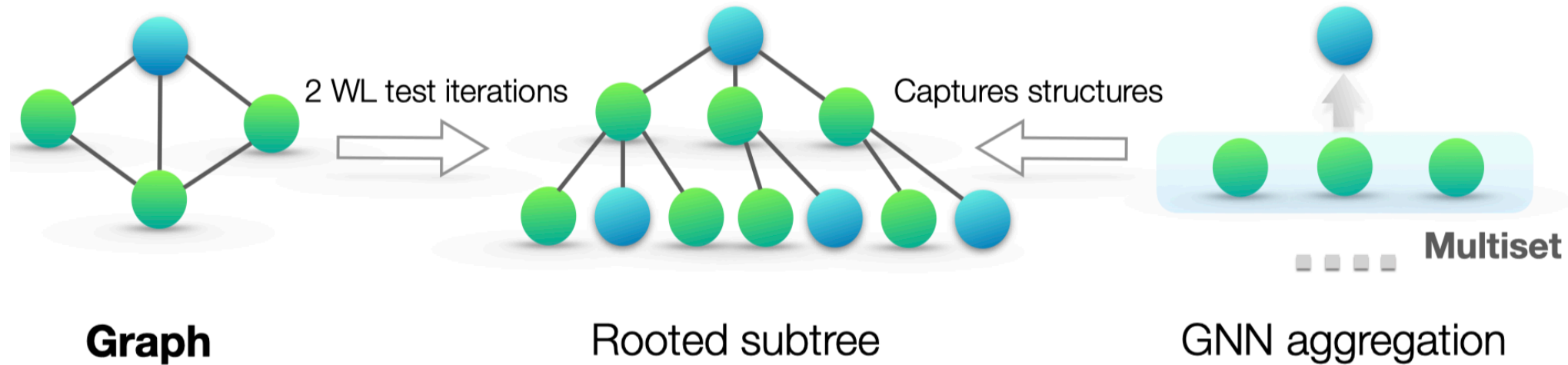
└──────────┘ └──────────────────────────┘

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

3 THEORETICAL FRAMEWORK



- To study the representational power of a GNN:
 - whether a GNN maps two neighborhoods (i.e., two multisets) to the same embedding or representation
 - A maximally powerful GNN would **never map** two **different neighborhoods** to the **same representation**
 - **This means its aggregation scheme must be injective**

4 BUILDING POWERFUL GRAPH NEURAL NETWORKS

Lemma 2. *Let G_1 and G_2 be any two non-isomorphic graphs. If a graph neural network $\mathcal{A} : \mathcal{G} \rightarrow \mathbb{R}^d$ maps G_1 and G_2 to different embeddings, the Weisfeiler-Lehman graph isomorphism test also decides G_1 and G_2 are not isomorphic.*

- Any aggregation-based GNN is at most as powerful as the WL test in distinguishing different graphs.
- GNN is as powerful as the WL test, if
 - neighbor aggregation is injective
 - graph-level readout function is injective
- an important benefit of GNNs beyond distinguishing different graphs
 - capturing **similarity** of graph structures
 - WL test are essentially one-hot encodings
 - GNN embed the subtrees to low-dimensional space

4 BUILDING POWERFUL GRAPH NEURAL NETWORKS

- GRAPH ISOMORPHISM NETWORK (GIN)

- generalizes the WL test and hence achieves maximum discriminative power among GNNs

- For node embedding

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

- For graph embedding

$$h_G = \text{CONCAT} \left(\text{READOUT} \left(\left\{ h_v^{(k)} \mid v \in G \right\} \right) \mid k = 0, 1, \dots, K \right).$$

5 LESS POWERFUL BUT STILL INTERESTING GNNS

- Study two aspects of the aggregator
 - 1-layer perceptrons instead of MLPs
 - mean or max-pooling instead of the sum
- MLP: Universal approximation theorem (Hornik et al., 1989; Hornik, 1991)
- Unlike models using MLPs, the 1-layer perceptron (even with the bias term) is not a universal approximator of multiset functions

Lemma 7. *There exist finite multisets $X_1 \neq X_2$ so that for any linear mapping W ,*
$$\sum_{x \in X_1} \text{ReLU}(Wx) \neq \sum_{x \in X_2} \text{ReLU}(Wx).$$

5 LESS POWERFUL BUT STILL INTERESTING GNNS

- **Mean** aggregators **captures** the proportion/**distribution** of elements of a given type
- **Max** aggregator ignores multiplicities, **learns** sets with **distinct elements**

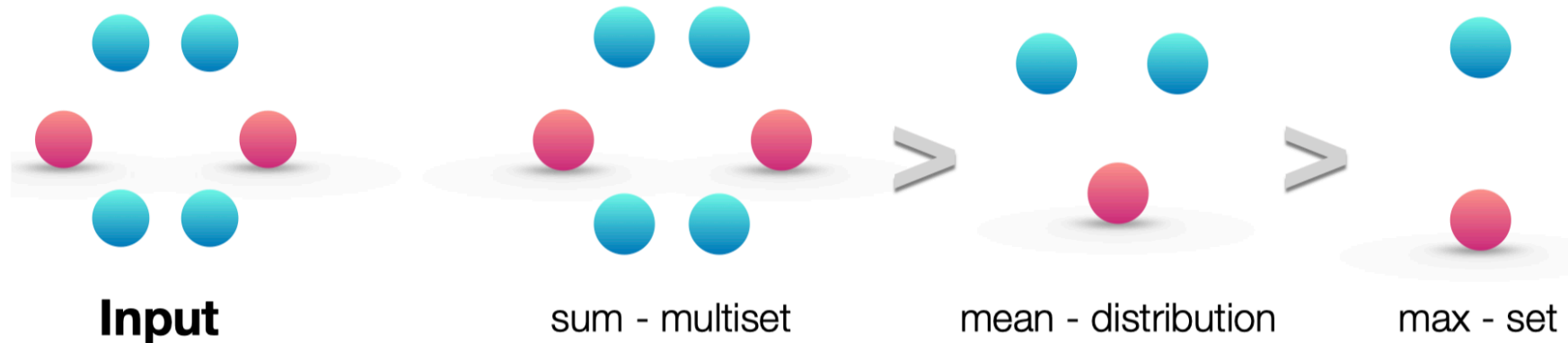


Figure 2: **Ranking by expressive power for sum, mean and max aggregators over a multiset.**

5 LESS POWERFUL BUT STILL INTERESTING GNNS

- **Mean** aggregator is as **powerful** as the sum aggregator when node features are diverse and **rarely repeat**
- **Max-pooling** may be **suitable** for tasks where it is important to identify representative elements or the “**skeleton**”, rather than to distinguish the exact structure or distribution

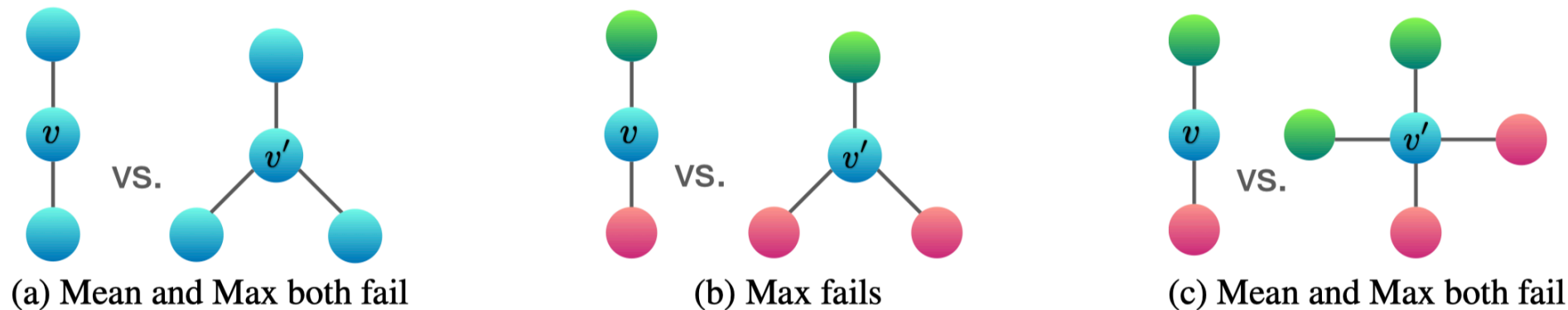


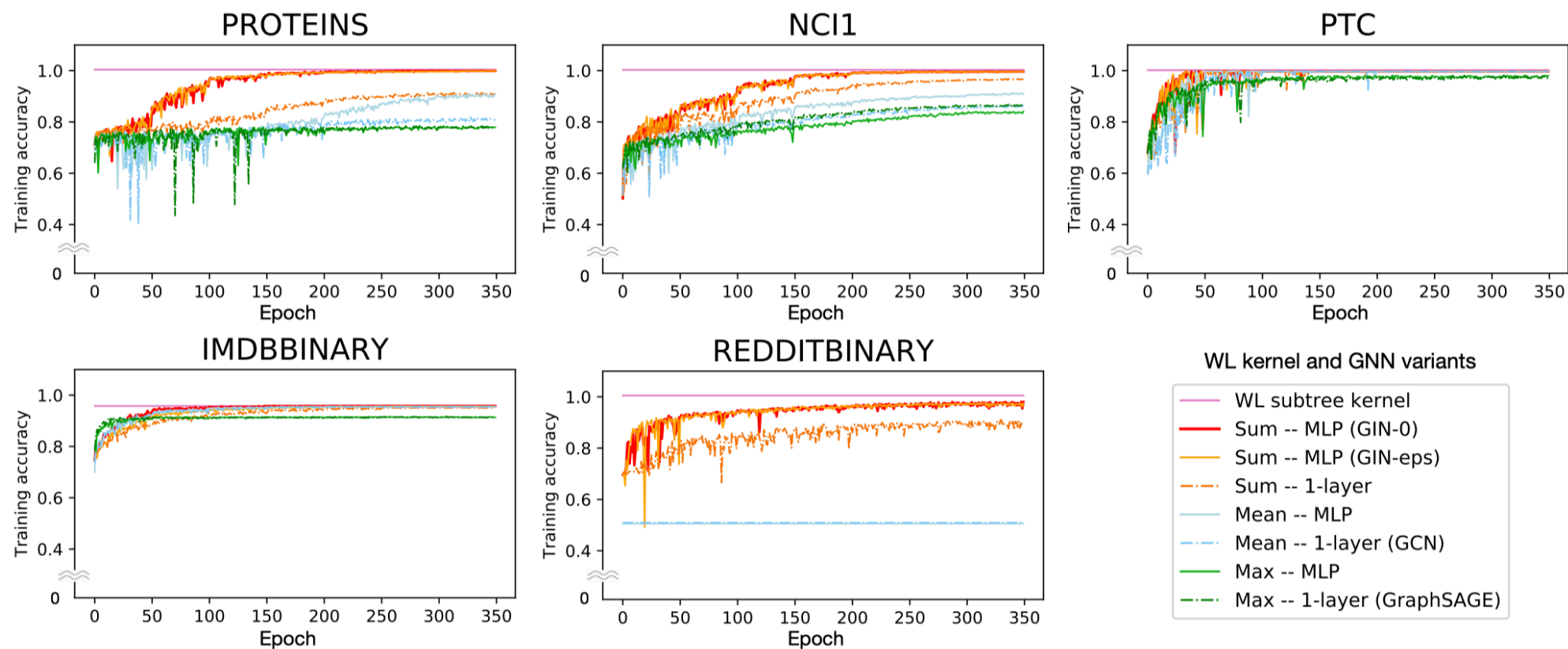
Figure 3: **Examples of graph structures that mean and max aggregators fail to distinguish.**

6 EXPERIMENTS

- 9 graph classification benchmarks
 - 4 bioinformatics datasets (MUTAG, PTC, NCI1, PROTEINS)
 - 5 social network datasets (COLLAB, IMDB-BINARY, IMDB-MULTI, REDDITBINARY and REDDIT-MULTI5K)
- Models and configurations:
 - 10-fold cross-validation
 - 5 GNN layers (including the input layer)
 - MLPs have 2 layers
 - learning rate 0.01

6 EXPERIMENTS

- Training set performance



6 EXPERIMENTS

- Test set performance

Datasets	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
Baselines	WL subtree	73.8 ± 3.9	50.9 ± 3.8	81.0 ± 3.1	52.5 ± 2.1	78.9 ± 1.9	90.4 ± 5.7	75.0 ± 3.1	59.9 ± 4.3	86.0 ± 1.8 *
	DCNN	49.1	33.5	–	–	52.1	67.0	61.3	56.6	62.6
	PATCHYSAN	71.0 ± 2.2	45.2 ± 2.8	86.3 ± 1.6	49.1 ± 0.7	72.6 ± 2.2	92.6 ± 4.2 *	75.9 ± 2.8	60.0 ± 4.8	78.6 ± 1.9
	DGCNN	70.0	47.8	–	–	73.7	85.8	75.5	58.6	74.4
	AWL	74.5 ± 5.9	51.5 ± 3.6	87.9 ± 2.5	54.7 ± 2.9	73.9 ± 1.9	87.9 ± 9.8	–	–	–
GNN variants	SUM-MLP (GIN-0)	75.1 ± 5.1	52.3 ± 2.8	92.4 ± 2.5	57.5 ± 1.5	80.2 ± 1.9	89.4 ± 5.6	76.2 ± 2.8	64.6 ± 7.0	82.7 ± 1.7
	SUM-MLP (GIN- ϵ)	74.3 ± 5.1	52.1 ± 3.6	92.2 ± 2.3	57.0 ± 1.7	80.1 ± 1.9	89.0 ± 6.0	75.9 ± 3.8	63.7 ± 8.2	82.7 ± 1.6
	SUM-1-LAYER	74.1 ± 5.0	52.2 ± 2.4	90.0 ± 2.7	55.1 ± 1.6	80.6 ± 1.9	90.0 ± 8.8	76.2 ± 2.6	63.1 ± 5.7	82.0 ± 1.5
	MEAN-MLP	73.7 ± 3.7	52.3 ± 3.1	50.0 ± 0.0	20.0 ± 0.0	79.2 ± 2.3	83.5 ± 6.3	75.5 ± 3.4	66.6 ± 6.9	80.9 ± 1.8
	MEAN-1-LAYER (GCN)	74.0 ± 3.4	51.9 ± 3.8	50.0 ± 0.0	20.0 ± 0.0	79.0 ± 1.8	85.6 ± 5.8	76.0 ± 3.2	64.2 ± 4.3	80.2 ± 2.0
	MAX-MLP	73.2 ± 5.8	51.1 ± 3.6	–	–	–	84.0 ± 6.1	76.0 ± 3.2	64.6 ± 10.2	77.8 ± 1.3
	MAX-1-LAYER (GraphSAGE)	72.3 ± 5.3	50.9 ± 2.2	–	–	–	85.1 ± 7.6	75.9 ± 3.2	63.9 ± 7.7	77.7 ± 1.5

7 CONCLUSION

- ❑ Weighted average via attention and LSTM pooling aggregation were not covered.
- ❑ Sum aggregator outperforms mean and max-pooling.
- ❑ MLP outperforms 1-layer perception which commonly used.

THANKS

6-29