

# Multi-label Node Classification Task On Graph-structured Data

Tianqi Zhao

09.11.2023

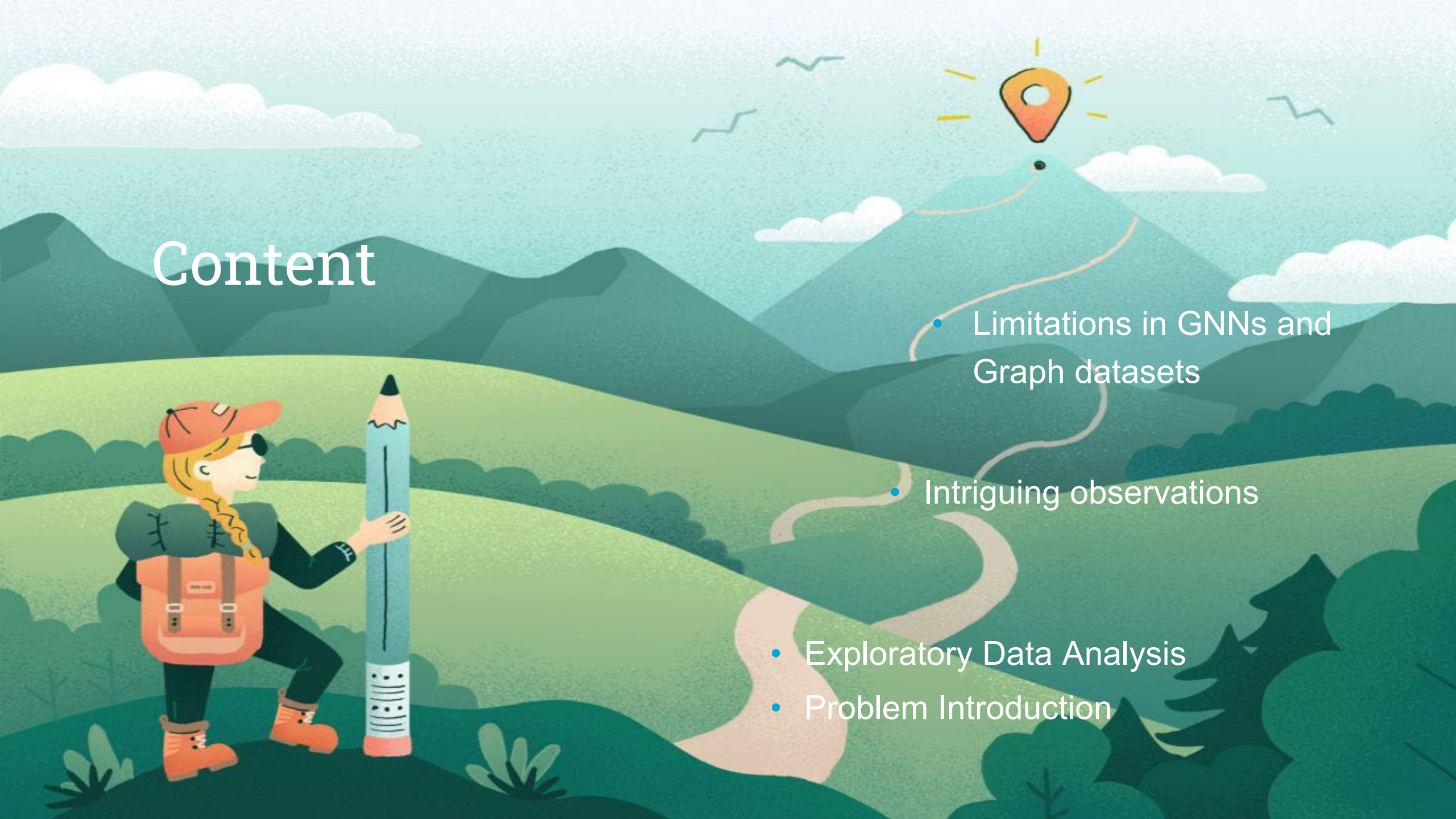
<https://openreview.net/forum?id=EZhkV2BjDP>





# Content

- Limitations in GNNs and Graph datasets
- Intriguing observations
- Exploratory Data Analysis
- Problem Introduction

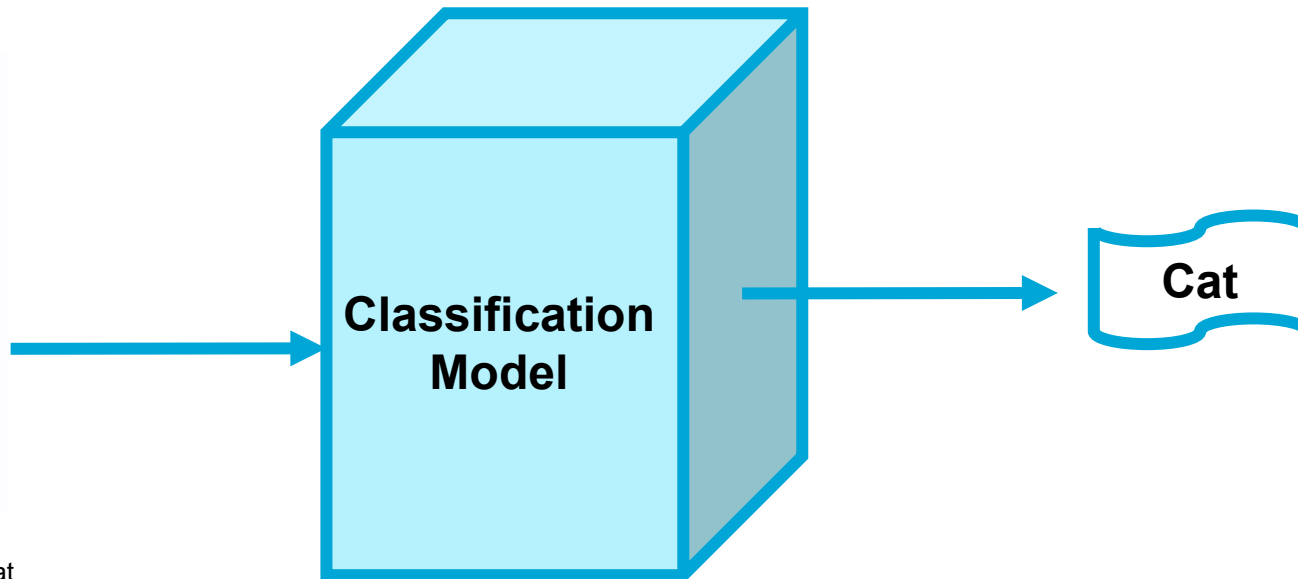


# Classification Task

- Input: description of the data
- Output: label



Image source: <https://stock.adobe.com/nl/search?k=cat>



# A More Realistic Scenario: Multi-label



Illustration of multi-label classification. One input Image contains multiple labels (cat and dog).

Image source: <https://www.pixelstalk.net/cute-dog-and-cat-wallpaper/=0>

# Multi-label Node Classification On Graphs

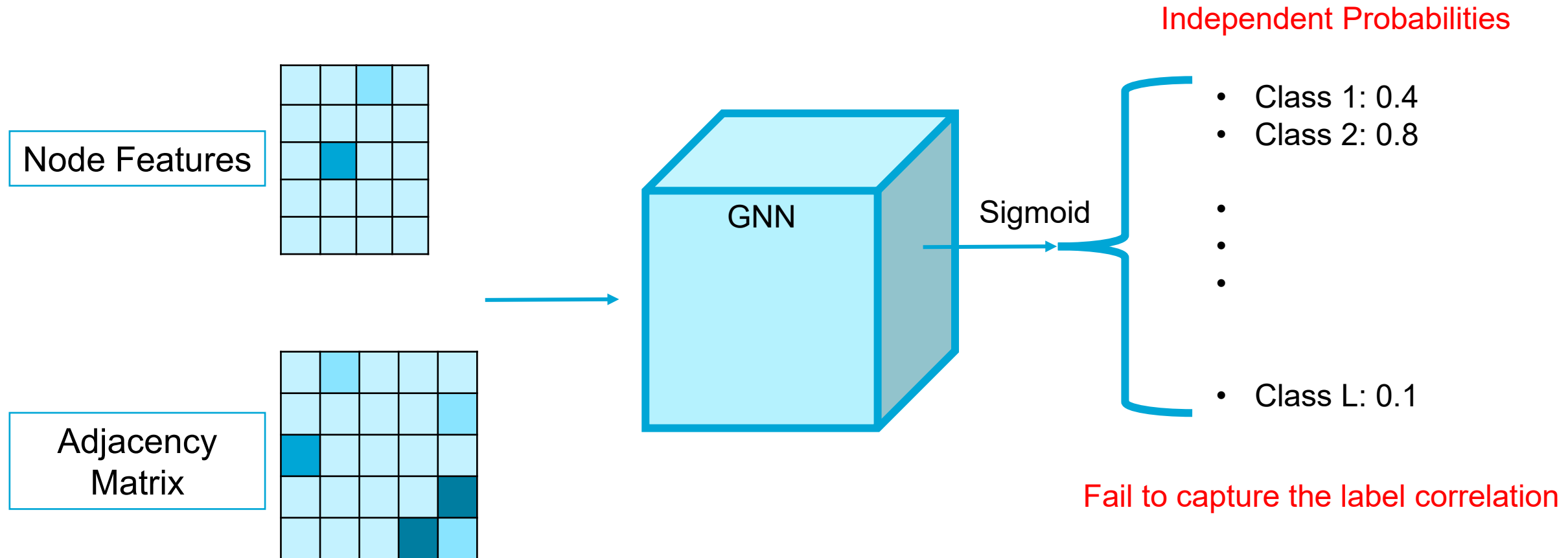
- **Graph:** Datapoints that are connected by the topological structure
- **Node classification task:**
  - **Given:** a graph, incomplete labelling
  - **Task:** predict the classes of the unlabeled nodes



Illustration of multi-label nodes in graphs

Image Source: <https://minutes.co/5-ways-to-improve-networking-skills-from-professionals-who-mastered-it/>

# GNNs Deployed For Multi-label Node Classification Task





# Performance of GNNs

- **Homophily in Single-label graphs:**
  - High homophily, GNNs show high performance
  - Most of the connected nodes share the same label
- **Homophily in Multi-label graphs:**
  - **Fraction of shared labels**
  - **Would GNNs work well on these datasets?**

Dataset	Label Homophily
Blogcatalog	0.10
Yelp	0.22
Ogbn-proteins	0.15

# Would GNNs designed for single-label low homophily graphs work on multi-label graphs?

- Low homophily: Opposites attract



# GNN Designed For Single Label Low Homophily Graphs: H2GCN

- Assumption: useful information in higher neighbourhoods

Method	BLOGCAT	YELP	OGB-PROTEINS	DBLP
MLP	0.043	0.096	0.026	0.350
DEEPWALK	<b>0.190</b>	0.096	0.044	0.585
LANC	<u>0.050</u>	OOM	<u>0.045</u>	0.836
GCN	0.037	0.131	<b>0.054</b>	<b>0.893</b>
GAT	0.041	0.150	0.021	0.829
GRAPHSAGE	0.045	<b>0.251</b>	0.027	<u>0.868</u>
H2GCN	<u>0.039</u>	<u>0.226</u>	<u>0.036</u>	<u>0.858</u>
GCN-LPA	0.043	0.116	0.023	0.801

# Different Semantics Of Homophily

- Low homophily: each neighbour only shares a small fraction of labels
- But all labels can be found among the neighbouring nodes

# Real-life Scenario



- Each friend shares part of one's interests: low homophily
- All the interests are in the direct neighbourhood

# Characteristics Of Ogbn-proteins

[Get Started](#)[Updates](#)[Large-Scale Challenge ▾](#)[Datasets ▾](#)[Leaderboards ▾](#)[Papers ▾](#)[Team](#)[Github](#)

## Leaderboard for [ogbn-proteins](#)

The ROC-AUC score on the test and validation sets. The higher, the better.

Package:  $\geq 1.1.1$

Rank	Method	Ext. data	Test ROC-AUC	Validation ROC-AUC	Contact	References	#Params	Hardware	Date
1	<b>LD+GAT</b>	Yes	0.8942 $\pm$ 0.0007	0.9527 $\pm$ 0.0007	<a href="#">Zhihao Shi (MIRA Lab, USTC &amp; CityBrain Lab, Alibaba Cloud)</a>	<a href="#">Paper</a> , <a href="#">Code</a>	664,233,700	GeForce RTX 3090 (24GB GPU)	Sep 27, 2023
2	<b>GIPA(Wide&amp;Deep)</b>	No	0.8917 $\pm$ 0.0007	0.9472 $\pm$ 0.0020	<a href="#">Houyi Li</a>	<a href="#">Paper</a> , <a href="#">Code</a>	17,438,716	Tesla V100-SXM2(32G)	Jan 19, 2023
3	<b>AGDN</b>	No	0.8865 $\pm$ 0.0013	0.9418 $\pm$ 0.0005	<a href="#">Chuxiong Sun</a>	<a href="#">Paper</a> , <a href="#">Code</a>	8,605,486	Tesla V100 (16GB GPU)	Sep 2, 2022
4	<b>RevGNN-Wide</b>	No	0.8824 $\pm$ 0.0015	0.9450 $\pm$ 0.0008	<a href="#">Guohao Li - DeepGCNs.org</a>	<a href="#">Paper</a> , <a href="#">Code</a>	68,471,608	NVIDIA RTX 6000 (48G)	Jun 16, 2021
5	<b>GAT+BOT+NGNN</b>	No	0.8809 $\pm$ 0.0016	0.9375 $\pm$ 0.0019	<a href="#">Xiang song (DGL team)</a>	<a href="#">Paper</a> , <a href="#">Code</a>	11,740,552	Tesla V100 (32GB)	Jan 23, 2022
6	<b>RevGNN-Deep</b>	No	0.8774 $\pm$ 0.0013	0.9326 $\pm$ 0.0006	<a href="#">Guohao Li - DeepGCNs.org</a>	<a href="#">Paper</a> , <a href="#">Code</a>	20,031,384	NVIDIA RTX 6000 (48G)	Jun 16, 2021
7	<b>GAT+BoT</b>	No	0.8765 $\pm$ 0.0008	0.9280 $\pm$ 0.0008	<a href="#">Yangkun Wang (DGL Team)</a>	<a href="#">Paper</a> , <a href="#">Code</a>	2,484,192	Tesla A100 (40GB GPU)	Jun 16,

Image source: [https://ogb.stanford.edu/docs/leader\\_nodeprop/](https://ogb.stanford.edu/docs/leader_nodeprop/)



# Characteristics Of Ogbn-proteins

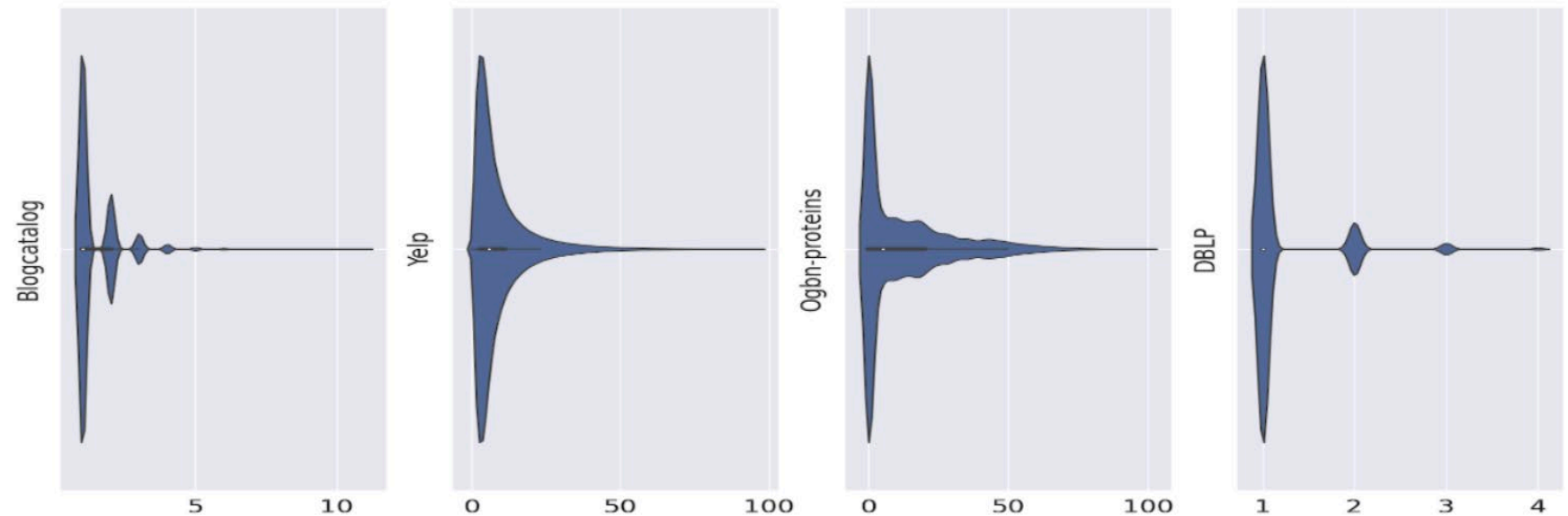
- Huge number of **unlabelled nodes**
  - in Ogbn-proteins, over **40% nodes without labels**
  - **89% test nodes without labels**
- Evaluation
  - Sparse label matrix, Long training epochs
  - encourages model to assign negative class and achieve high AUROC

	Model	MicroF1	MacroF1	AUROC	AP
OGB- PROTEINS	MLP	2.55	2.40	54.05	2.59
	DEEPWALK	<b>2.88</b>	<b>2.75</b>	<u>68.75</u>	4.41
	LANC	2.35	2.21	68.03	<u>4.48</u>
	GCN	<u>2.77</u>	<u>2.63</u>	<b>71.48</b>	<b>5.36</b>
	GAT	2.55	2.40	50.64	2.14
	GRAPHSAGE	2.59	2.43	55.83	2.68
	H2GCN	2.55	2.39	62.75	3.61
	GCN-LPA	2.56	2.41	53.22	2.33

# Characteristics Of Available Multi-label Graphs

- Label assignment
  - Most of the nodes have one label
    - in Blogcatalog, 72% nodes have one label
  - Nodes: bloggers
  - Labels: social groups
  - Edges: friendships

Is this true?



# Multi-label graph generator model

- Varying feature quality and homophily

Method	$r_{ori\_feat}$ <span style="color: red;">high</span> <span style="color: red;">→</span>					$r_{homo}$ <span style="color: red;">high</span> <span style="color: red;">→</span>				
	0.0	0.2	0.5	0.8	1.0	0.2	0.4	0.6	0.8	1.0
MLP	0.172	0.187	0.220	0.277	0.343	<b>0.343</b>	0.343	0.343	0.343	0.343
DEEPWALK	<b>0.487</b>	<b>0.487</b>	<b>0.487</b>	<b>0.487</b>	<b>0.487</b>	0.181	<b>0.522</b>	<b>0.813</b>	<b>0.869</b>	0.552
LANC	0.337	0.342	0.365	0.353	0.391	0.190	0.380	0.434	0.481	<u>0.629</u>
GCN	0.313	0.316	0.311	0.301	0.337	0.261	0.343	0.388	0.450	0.493
GAT	0.311	0.339	0.329	0.338	0.360	0.172	0.359	0.390	0.428	0.439
GRAPHSAGE	0.300	0.328	0.377	0.393	0.430	0.289	0.426	0.458	0.533	0.553
H2GCN	<u>0.376</u>	<u>0.401</u>	<u>0.427</u>	<u>0.442</u>	<u>0.467</u>	<u>0.297</u>	<u>0.484</u>	<u>0.512</u>	0.572	<b>0.652</b>
GCN-LPA	0.337	0.333	0.368	0.363	0.391	0.170	0.408	0.495	<u>0.604</u>	0.583

# Conclusions

- Multi-label Node Classification Task is still an open research field
- Check the datasets before the experiments
- Choose the right model on the datasets
- Use the right evaluation metric

Paper: Zhao et al. "Multi-label Node Classification On Graph-structured Data", TMLR 2023. <https://openreview.net/forum?id=EZhkV2BjDP>