# GraphLand: Evaluating Graph Machine Learning Models on Diverse Industrial Data

Extended Abstract Track Submissions

#### **Anonymous Author(s)**

Anonymous Affiliation
Anonymous Email

## **Abstract**

Although data that can be naturally represented as graphs is widespread in realworld applications across diverse industries, popular graph ML benchmarks for node property prediction only cover a surprisingly narrow set of data domains, and graph neural networks (GNNs) are often evaluated on just a few academic citation networks. This issue is particularly pressing in light of the recent growing interest in designing graph foundation models. These models are supposed to be able to transfer to diverse graph datasets from different domains, and yet the proposed graph foundation models are often evaluated on a very limited set of datasets from narrow applications. To alleviate this issue, we introduce Graph-Land: a benchmark of 14 diverse graph datasets for node property prediction from a range of different industrial applications. GraphLand allows evaluating graph ML models on a wide range of graphs with diverse sizes, structural characteristics, and feature sets, all in a unified setting. Further, GraphLand allows investigating such previously underexplored research questions as how realistic temporal distributional shifts under transductive and inductive settings influence graph ML model performance. We evaluate a range of GNNs on GraphLand datasets and show that they significantly outperform graph-agnostic models in realistic settings. We also evaluate currently available general-purpose graph foundation models and find that they fail to produce competitive results on our proposed datasets.

## 1 Introduction

2

5

6

8

9

10

11

12

13

14 15

16

18

19

20

21

23

28

31

34

35

36

37

38

Recently, there has been a significant push for data-centric approaches in machine learning. In particular, high-quality, realistic, reliable, and diverse benchmarks are paramount for proper evaluation of the performance of machine learning methods. In the field of graph machine learning (GML), there has recently been a lot of criticism of existing popular benchmark datasets concerning such aspects as lacking practical relevance [1], low structural diversity that leaves most of the possible graph structure space not represented [2, 3], low application domain diversity [1], graph structure not being beneficial for the considered tasks [1, 4–6], potential bugs in the data collection processes leading to incorrect labels [7] and duplicated graph nodes [8]. While there have recently been efforts to create more realistic graph benchmarks, they focus on more specific domains (e.g., 3D molecular data) that require specialized models. At the same time, benchmarks for the standard and most widespread GML setting of node property prediction in a single large graph have received considerably less attention, and evaluation of the performance of classic graph neural networks and recent graph foundation models is still often limited to a few academic citation networks despite this setting and models developed for it having vast real-world applications in diverse industries.

We believe the historical focus of GNN evaluation on academic citation networks, which represent only a single (and a rather narrow) application domain, is primarily a consequence of the availability of open data of this type, rather than its relevance to real-world applications. At the same time, some of the most classical and simultaneously practically important examples of real-world graphs —

social networks, web graphs, and road networks — are surprisingly rarely used for GNN evaluation, perhaps due to the lack of easily accessible high-quality datasets.

Recently, there has been a lot of interest in developing graph foundation models (GFMs) — models that after large-scale pretraining can be applied to diverse graph datasets without or with minimal fine-tuning [9, 10]. Proper evaluation of such models thus requires the use of a diverse set of realistic graph datasets. However, currently the proposed GFMs are frequently evaluated only on text-attributed graphs (and mostly citation networks), thus overlooking the problem of transferring to graphs with different node feature sets, which is required for truly general GFMs as graphs in real-world applications from different domains often come with completely different node feature sets.

It has been argued that due to the unavailability of diverse realistic industrial datasets for researchers it is worth shifting the evaluation of GML models to synthetic datasets [2, 11]. However, we believe that it is important to evaluate models on real-world data as much as possible, both to obtain unbiased estimates of model performance in realistic scenarios and to showcase the potential of GML in industrial applications. Thus, it is desirable to have open and easily accessible diverse and realistic graph datasets.

55

56

57

62

63

64

65

82

83

87

89

91

92

In our work, we aim to alleviate the issue of a lack of realistic GNN benchmarks for node property prediction by introducing GraphLand: a collection of graphs and associated machine learning tasks collected from a variety of industrial applications that represent real-world GML usage. GraphLand significantly extends the set of available datasets for GML model evaluation, providing in a unified format 14 graph datasets, many of which represent applications or structural properties that have not been covered by standard GML benchmarks before. The datasets in GraphLand have been collected both from open data that has been underutilized or not utilized at all in the field of GML, and from newly released data from services of a large technological company for which the use of GML has internally proven its usefulness. A key feature of GraphLand is its diversity, with graphs spanning a wide range of domains, sizes, and structural properties, and having rich node features with different types, meanings, and distributions.

For datasets in GraphLand, we provide several data splits, including a realistic temporal one, which allows for investigating practically important questions previously underexplored in GML literature: how temporal distributional shifts affect the performance of GML models in both transductive and inductive settings.

We run extensive experiments on GraphLand datasets with a range of GNNs and several openly available GFMs. We find that GNNs can achieve great results in industrial applications with attention-based GNNs often performing better than more classic ones. However, their performance can be strongly affected by temporal distributional shifts and dynamically evolving graph structure, which highlights the importance of developing models more resilient to such changes. Further, we find that currently available GFMs perform poorly on our datasets and fail to achieve results competitive with more classic methods.

We hope that GraphLand will allow more diverse and realistic evaluation of GML models, as well as encourage research into currently underexplored directions such as designing GML models that are more resilient to temporal distributional shifts and dynamically evolving graphs, and designing GFMs that are truly generalizable to graph data from different domains with different node feature sets.

## 2 Limitations of Popular Graph Machine Learning Benchmarks

By far the most popular datasets used in modern GML literature are the three academic citation networks cora, citeseer, and pubmed [12–16]. These datasets became so widespread perhaps because they were used by the foundational work on modern GNNs by Kipf and Welling [17]. However, these datasets only cover a single and rather narrow application of paper subject prediction in citation networks. Another popular set of datasets for GML was introduced by Shchur et al. [18] and includes academic coauthorship networks coauthor-cs and coauthor-physics, and e-commerce co-purchasing networks amazon-computers and amazon-photo. However, all the aforementioned datasets together only cover three applications, while GML methods can be used in a much wider variety of settings. Later, larger-scale graph datasets were introduced in the Open Graph Benchmark (OGB) [19]. However, out of the five node property prediction datasets, three are academic citation networks (ogbn-arxiv, ogbn-mag, ogbn-papers100M) and one more is an e-commerce co-purchasing network (ogbn-products). Thus, OGB does not significantly expand the

96

97

98

102

103

104

113

123

125

126

127

130

131

135

137 138

139

141

143

144

145

146

range of real-world applications available for evaluating GML models. Further, all the aforementioned datasets only provide textual descriptions as node features. However, co-purchasing networks, which represent a very practically important application of GML, in realistic settings come with rich product metadata that can be represented as numerical and categorical features. There is currently a lack of datasets with such metadata available among popular GML benchmarks. Further, all the aforementioned datasets are homophilous, i.e., edges in them typically connect nodes of the same class. It has been shown by Huang et al. [20] that even very simple models can provide strong results in homophilous networks. Thus, it is important for standard GML benchmarks to also include a wide selection of *non-homophilous* graphs, i.e., graphs in which edges do not have the tendency to connect nodes of the same class. For a long time, the only popular source of non-homophilous graph datasets was the benchmark from Pei et al. [21]. However, it was recently shown by Platonov et al. [8] that these datasets have numerous problems including duplicated nodes, small size (leading to noisy evaluation metric estimates), and insufficient class representation (such as the texas dataset having a class that consists of only a single node). While Platonov et al. [8] introduced several new non-homophilous graph datasets, they were meant to be used to reliably reevaluate the performance of different models in absence of homophily, rather than represent realistic GML applications. Thus, some of these datasets are synthetic (minesweeper), semi-synthetic (roman-empire), or have limited node features despite the original data source potentially providing more information about the nodes (tolokers, questions).

Overall, it can be seen that the currently popular graph datasets for node property prediction do not allow evaluating GNNs and other GML models on a wide range of practically impactful industrial applications.

**Text-attributed graphs and generalization of graph foundation models.** Most of the datasets frequently used for node property prediction only have textual descriptions as node attributes. However, graphs representing real-world networks often have rich and diverse node attributes that go beyond just texts and encompass a variety of numerical and categorical features with different meanings and distributions. There has recently been a lot of interest in developing general-purpose GFMs that are expected to generalize to graphs from different domains [9, 10]. A key challenge for such GFMs is being able to adapt to graphs with different node feature sets which is required for a model truly generalizable to different domains. Yet, the GFMs that have been proposed in the current literature typically overlook this challenge and are often only evaluated on text-attributed graphs [22–24]. Textual attributes can be easily projected to a common latent feature space by applying pretrained text encoders based on Large Language Models, thus allowing a single GFM to work with different text-attributed graphs. The prevalence of such text-attributed graphs in GML benchmarks has led to most of current GFM research overlooking the problem of generalization to different node feature sets, since it is not required to solve tasks from standard benchmarks. However, this problem is very important for real-world industrial applications of GML in which graphs often come attributed with a mixture of numerical and categorical node features. GFMs must therefore be able to work with such features to effectively solve practical tasks. The problem of devising a single foundation model that can work with arbitrary numerical and categorical features has received significant attention from the ML for tabular data community, where such features are standard, and first successful attempts to develop such a model have recently emerged [25-31]. However, these ideas have not yet spread to the GML community (likely because current standard GML benchmarks do not require working with non-textual node features) despite the significant benefits of designing a successful GFM that can handle arbitrary node feature sets for practical applications.

# 3 GraphLand: A Collection of Diverse Industrial Graph Datasets

GraphLand is a collection of 14 graph datasets with node property prediction tasks (either classification or regression). Some of these datasets are newly released for this benchmark, while others are collected from open data sources that are underutilized or not utilized at all in current GML benchmarking. In selecting datasets for GraphLand, we aim to fulfill the following desiderata: datasets should come from diverse fields representing impactful industrial applications of GML, graphs should exhibit a range of different sizes and structural characteristics, nodes should have rich features, graph structure should be beneficial for the considered tasks.

Here we briefly describe our datasets, while detailed information is provided in Appendix A. First, we describe our newly released datasets. web-fraud, web-topics, and web-traffic represent

**Table 1:** Experimental results under the RL (random low) data split in the transductive setting. The best result and those statistically indistinguishable from it are highlighted in orange. TLE stands for time limit exceeded (24 hours); RTE stands for runtime error in the official code of GFMs.

	mu	ticlass classification	1	binary classification				
	hm-categories	pokec-regions	web-topics	tolokers-2	city-reviews	artnet-exp	web-fraud	
best const. pred. ResMLP	$\begin{array}{c} 19.46 \pm 0.00 \\ 37.72 \pm 0.18 \end{array}$	$3.77 \pm 0.00$ $4.88 \pm 0.01$	$\begin{array}{c} 28.36 \pm 0.00 \\ 42.41 \pm 0.02 \end{array}$	$\begin{array}{c} 21.82 \pm 0.00 \\ 41.16 \pm 1.13 \end{array}$	$12.09 \pm 0.00 \\ 71.32 \pm 0.11$	$\begin{array}{c} 10.00 \pm 0.00 \\ 35.07 \pm 2.34 \end{array}$	$0.66 \pm 0.00 \\ 8.77 \pm 0.18$	
GCN GraphSAGE GAT GT	$61.70 \pm 0.35$ $56.75 \pm 0.53$ $67.96 \pm 0.33$ $69.23 \pm 0.50$	$34.96 \pm 0.38$ $37.88 \pm 0.41$ $46.17 \pm 0.32$ $46.47 \pm 0.16$	$46.45 \pm 0.10 47.41 \pm 0.13 48.25 \pm 0.05 48.00 \pm 0.05$	$\begin{array}{c} 51.32 \pm 0.96 \\ 53.73 \pm 0.53 \\ 53.78 \pm 1.34 \\ 54.50 \pm 1.20 \end{array}$	$77.15 \pm 0.28$ $77.82 \pm 0.13$ $77.67 \pm 0.13$ $76.97 \pm 0.21$	$43.09 \pm 0.38  42.65 \pm 0.59  46.62 \pm 0.32  45.16 \pm 0.46$	$10.02 \pm 0.18$ $12.11 \pm 0.23$ $13.32 \pm 0.29$ $12.74 \pm 0.42$	
OpenGraph (ICL) AnyGraph (ICL) GCOPE (FT)	$\begin{array}{c} 9.49 \pm 0.93 \\ 15.47 \pm 2.36 \\ 19.51 \pm 0.07 \end{array}$	$\begin{array}{c} 1.73 \pm 0.31 \\ 24.65 \pm 1.51 \\ \text{TLE} \end{array}$	$\begin{array}{c} \text{RTE} \\ 6.67 \pm 3.88 \\ \text{TLE} \end{array}$	$\begin{array}{c} 40.49 \pm 0.31 \\ 31.33 \pm 2.89 \\ 28.67 \pm 1.42 \end{array}$	$58.44 \pm 1.08$ $64.37 \pm 1.29$ $67.38 \pm 1.23$	$\begin{array}{c} 15.65 \pm 1.23 \\ 13.14 \pm 1.15 \\ 16.10 \pm 2.79 \end{array}$	$\begin{array}{c} \text{RTE} \\ 0.68 \pm 0.03 \\ \text{TLE} \end{array}$	

**(b)** Results for regression datasets.  $R^2$  is reported for all datasets.

	hm-prices	avazu-ctr	city-roads-M	city-roads-L	twitch-views	artnet-views	web-traffic
best const. pred. ResMLP	$\begin{array}{c} 0.00 \pm 0.00 \\ 62.66 \pm 0.37 \end{array}$	$0.00 \pm 0.00$ $24.54 \pm 0.36$	$\begin{array}{c} 0.00 \pm 0.00 \\ 54.77 \pm 0.15 \end{array}$	$\begin{array}{c} 0.00 \pm 0.00 \\ 46.47 \pm 0.29 \end{array}$	$0.00 \pm 0.00$ $13.35 \pm 0.02$	$\begin{array}{c} 0.00 \pm 0.00 \\ 29.71 \pm 0.60 \end{array}$	$\begin{array}{c} 0.00 \pm 0.00 \\ 72.42 \pm 0.05 \end{array}$
GCN GraphSAGE GAT GT	$70.54 \pm 0.21$	$30.47 \pm 0.27$ $31.84 \pm 0.24$ $33.20 \pm 0.20$ $30.87 \pm 0.47$	$59.05 \pm 0.16$ $57.51 \pm 0.53$ $59.11 \pm 0.20$ $58.05 \pm 0.58$	$53.26 \pm 0.14$ $52.43 \pm 0.25$ $53.43 \pm 0.20$ $53.38 \pm 0.12$	$75.55 \pm 0.05$ $66.87 \pm 0.11$ $72.93 \pm 0.17$ $72.19 \pm 0.14$	$55.99 \pm 0.26$ $49.79 \pm 0.51$ $53.36 \pm 0.78$ $54.23 \pm 0.22$	$82.07 \pm 0.14$ $83.50 \pm 0.11$ $84.68 \pm 0.06$ $84.49 \pm 0.07$

a part of the Internet (web-graph). artnet-views and artnet-exp represent a social network of art creators. city-roads-M and city-roads-L represent road networks of two major cities. city-reviews represents a network of users of a review of places and organizations service. Further, we describe datasets obtained from open sources that were previously underutilized. hm-categories and hm-prices represent co-purchasing networks. avazu-ctr represents a network of devices from which advertisements can be viewed. pokec-regions represents a large social networks. twitch-views represents a network of content streamers. tolokers-2 represents a network of crowdsourcing platform workers.

For each dataset, we prepare several different train/val/test splits. Random-low (RL) and random-high RH are random splits with low and high labeled rates. Temporal-high TH is a temporal split. Additionally, temporal-high/inductive (THI) setting allows evaluation models under an inductive settings by presenting different snapshots of dynamically evolving networks for train, val, and test. More details on our data splits are provided in Appendix C

## 4 Experiments

In Table 1 we present our experimental results for the RL split. It can be seen that classic GNNs (GCN [17] and GraphSAGE[32]) demonstrate strong results significantly outperforming graphagnostic model ResMLP and demonstrating the potential of GNNs in industrial applications. Attention augmented GNNs GAT [33] and neighborhood-attention Graph Transformer (GT) [34] also produce strong results often outperforming classic GNNs. Further, we evaluate currently openly available GFMs. Despite a lot of works proposing different GFMs, we find only three open models that can perform node classification in graphs with arbitrary node features (OpenGraph [35], AnyGraph[36], GCOPE [37]) and none that can perform node regression. In our experiments, these models produce poor results and cannot compete with classic GNNs.

We present results for other splits in Appendix D where we show that GNN performance significantly decreases under temporal distributional shifts and in the inductive setting, suggesting adaptation of models to these realistic settings is an important direction for future research.

We hope GraphLand will encourage the evaluation of GML methods under more realistic and diverse settings, the development of GML methods that are more resilient to temporal distributional shifts and dynamically changing graph structure, and the development of more consistently performing GFMs that can handle different node feature sets that go beyond just textual descriptions.

#### References

180

186

187

188

195

196

197

198

199

200

- [1] Maya Bechler-Speicher, Ben Finkelshtein, Fabrizio Frasca, Luis Müller, Jan Tönshoff, Antoine Siraudin, Viktor Zaverkin, Michael M. Bronstein, Mathias Niepert, Bryan Perozzi, Mikhail Galkin, and Christopher Morris. Position: Graph Learning Will Lose Relevance Due To Poor Benchmarks. In *International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=nDFpl2lhoH. 1
  - [2] John Palowitch, Anton Tsitsulin, Brandon Mayer, and Bryan Perozzi. GraphWorld: Fake Graphs Bring Real Insights for GNNs. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 3691–3701, 2022. 1, 2
- [3] Seiji Maekawa, Koki Noda, Yuya Sasaki, and Makoto Onizuka. Beyond real-world benchmark
   datasets: An empirical study of node classification with GNNs. Advances in Neural Information
   Processing Systems, 35:5562–5574, 2022. 1
- [4] Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. A fair comparison of graph neural networks for graph classification. *International Conference on Learning Representations* (ICLR), 2020. 1
  - [5] Zhengdao Li, Yong Cao, Kefan Shuai, Yiming Miao, and Kai Hwang. Rethinking the effectiveness of graph classification datasets in benchmarks for assessing gnns. *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, 2024.
  - [6] Corinna Coupette, Jeremy Wayland, Emily Simons, and Bastian Rieck. No metric to rule them all: Toward principled evaluations of graph-learning datasets. In *International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=XbmBNwrfG5. 1
- Yuwen Li, Miao Xiong, and Bryan Hooi. GraphCleaner: Detecting Mislabelled Samples in
   Popular Graph Learning Benchmarks. In *International Conference on Machine Learning*, pages
   203 20195–20209. PMLR, 2023. 1
- [8] Oleg Platonov, Denis Kuznedelev, Michael Diskin, Artem Babenko, and Liudmila Prokhorenkova. A critical look at the evaluation of GNNs under heterophily: Are we really making progress? *International Conference on Learning Representations (ICLR)*, 2023. 1, 3, 10, 12, 13
- Yuxiang Wang, Wenqi Fan, Suhang Wang, and Yao Ma. Towards Graph Foundation Models: A
   Transferability Perspective. arXiv preprint arXiv:2503.09363, 2025. 2, 3
- 210 [10] Haitao Mao, Zhikai Chen, Wenzhuo Tang, Jianan Zhao, Yao Ma, Tong Zhao, Neil Shah, Mikhail Galkin, and Jiliang Tang. Position: Graph foundation models are already here. In International Conference on Machine Learning, 2024. URL https://openreview.net/forum?id=Edz0QXKKAo. 2, 3
- [11] Minji Yoon, Yue Wu, John Palowitch, Bryan Perozzi, and Ruslan Salakhutdinov. Graph
   Generative Model for Benchmarking Graph Neural Networks. In *International Conference on Machine Learning*. PMLR, 2023.
- [12] C Lee Giles, Kurt D Bollacker, and Steve Lawrence. CiteSeer: An automatic citation indexing system. In *Proceedings of the third ACM conference on Digital libraries*, pages 89–98, 1998. 2
- 219 [13] Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating the construction of internet portals with machine learning. *Information Retrieval*, 3:127–163, 2000.
- Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- 224 [15] Galileo Namata, Ben London, Lise Getoor, Bert Huang, and U Edu. Query-driven active surveying for collective classification. In *10th international workshop on mining and learning with graphs*, volume 8, page 1, 2012.
- Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning
   with graph embeddings. In *International Conference on Machine Learning*, pages 40–48.
   PMLR, 2016. 2
- [17] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *International Conference on Learning Representations (ICLR)*, 2017. 2, 4, 13

- [18] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018. 2
- [19] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
   Catasta, and Jure Leskovec. Open Graph Benchmark: Datasets for Machine Learning on Graphs.
   Advances in Neural Information Processing Systems, 33:22118–22133, 2020. 2, 12
- [20] Qian Huang, Horace He, Abhay Singh, Ser-Nam Lim, and Austin R Benson. Combining
   Label Propagation and Simple Models Out-performs Graph Neural Networks. *International Conference on Learning Representations (ICLR)*, 2021. 3
- [21] Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-GCN: Geometric Graph Convolutional Networks. *International Conference on Learning Representations* (ICLR), 2020. 3, 12
- [22] Hao Liu, Jiarui Feng, Lecheng Kong, Ningyue Liang, Dacheng Tao, Yixin Chen, and Muhan
   Zhang. One for All: Towards Training One Graph Model for All Classification Tasks. *International Conference on Learning Representations (ICLR)*, 2024. 3
- Yuhan Li, Peisong Wang, Zhixun Li, Jeffrey Xu Yu, and Jia Li. ZeroG: Investigating Cross-dataset Zero-shot Transferability in Graphs. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1725–1735, 2024.
- Yufei He, Yuan Sui, Xiaoxin He, and Bryan Hooi. UniGraph: Learning a Unified Cross-Domain Foundation Model for Text-Attributed Graphs. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 1*, pages 448–459, 2025. 3
- [25] Noah Hollmann, Samuel Müller, Katharina Eggensperger, and Frank Hutter. TabPFN: A
   Transformer That Solves Small Tabular Classification Problems in a Second. *International Conference on Learning Representations (ICLR)*, 2023. 3
- Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin Hoo, Robin Tibor Schirrmeister, and Frank Hutter. Accurate predictions on small data with a tabular foundation model. *Nature*, 637(8045):319–326, 2025.
- 258 [27] Han-Jia Ye, Qile Zhou, Huai-Hong Yin, De-Chuan Zhan, and Wei-Lun Chao. Rethinking 259 Pre-Training in Tabular Data: A Neighborhood Embedding Perspective. *arXiv preprint* 260 *arXiv:2311.00055*, 2023.
- [28] Han-Jia Ye, Si-Yang Liu, and Wei-Lun Chao. A Closer Look at TabPFN v2: Understanding Its Strengths and Extending Its Capabilities. *arXiv preprint arXiv:2502.17361*, 2025.
- Andreas C Mueller, Carlo A Curino, and Raghu Ramakrishnan. MotherNet: Fast Training and Inference via Hyper-Network Transformers. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=6H4jRWKFc3.
- [30] Junwei Ma, Valentin Thomas, Rasa Hosseinzadeh, Hamidreza Kamkari, Alex Labach, Jesse C
   Cresswell, Keyvan Golestan, Guangwei Yu, Maksims Volkovs, and Anthony L Caterini. Tab DPT: Scaling Tabular Foundation Models on Real Data. arXiv preprint arXiv:2410.18164,
   2024.
- [31] Jingang Qu, David Holzmüller, Gaël Varoquaux, and Marine Le Morvan. TabICL: A Tabular
   Foundation Model for In-Context Learning on Large Data. In *International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=0VvD1PmNzM. 3
- [32] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in Neural Information Processing Systems*, 30, 2017. 4, 12, 13
- 275 [33] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua 276 Bengio. Graph attention networks. *International Conference on Learning Representations* 277 (*ICLR*), 2018. 4, 13
- Yunsheng Shi, Zhengjie Huang, Shikun Feng, Hui Zhong, Wenjin Wang, and Yu Sun. Masked
   label prediction: Unified message passing model for semi-supervised classification. *Proceedings* of the Thirtieth International Joint Conference on Artificial Intelligence, 2021. 4, 13
- [35] Lianghao Xia, Ben Kao, and Chao Huang. OpenGraph: Towards Open Graph Foundation
   Models. Proceedings of the 2024 Conference on Empirical Methods in Natural Language
   Processing, 2024. 4, 13

- 284 [36] Lianghao Xia and Chao Huang. Anygraph: Graph foundation model in the wild. *arXiv preprint* arXiv:2408.10700, 2024. 4, 13
- [37] Haihong Zhao, Aochuan Chen, Xiangguo Sun, Hong Cheng, and Jia Li. All in One and One for
   All: A Simple yet Effective Method towards Cross-domain Graph Pretraining. In *Proceedings* of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pages
   4443–4454, 2024. 4, 13
- 290 [38] Steve Wang and Will Cukierski. Click-through rate prediction, 2014. URL https://kaggle.com/competitions/avazu-ctr-prediction. 9
- 292 [39] Sergei Ivanov and Liudmila Prokhorenkova. Boost then convolve: Gradient boosting meets 293 graph neural networks. In *International Conference on Learning Representations (ICLR)*, 2021. 294 URL https://openreview.net/forum?id=ebS5NUfoMKL. 9
- [40] Carlos García Ling et al. H&M Personalized Fashion Recommendations, 2022. URL https://kaggle.com/competitions/h-and-m-personalized-fashion-recommendations. 9
- Lubos Takac and Michal Zabovsky. Data analysis in public social networks. In *International scientific conference and international workshop present day trends of innovations*, volume 1, 2012. 10
- [42] Derek Lim, Felix Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Bhalerao, and
   Ser Nam Lim. Large scale learning on non-homophilous graphs: New benchmarks and strong
   simple methods. Advances in Neural Information Processing Systems, 34:20887–20902, 2021.
   10, 11
- Benedek Rozemberczki and Rik Sarkar. Twitch gamers: a dataset for evaluating proximity preserving and structural role-based node embeddings, 2021. 10
- Daniil Likhobaba, Nikita Pavlichenko, and Dmitry Ustalov. Toloker Graph: Interaction of Crowd Annotators, 2023. URL https://github.com/Toloka/TolokerGraph. 10
- Stefano Boccaletti, Ginestra Bianconi, Regino Criado, Charo I Del Genio, Jesús Gómez-Gardenes, Miguel Romance, Irene Sendina-Nadal, Zhen Wang, and Massimiliano Zanin. The structure and dynamics of multilayer networks. *Physics reports*, 544(1):1–122, 2014. 10
- 311 [46] Oleg Platonov, Denis Kuznedelev, Artem Babenko, and Liudmila Prokhorenkova. Characteriz-312 ing Graph Datasets for Node Classification: Homophily-Heterophily Dichotomy and Beyond. 313 Advances in Neural Information Processing Systems, 36:523–548, 2023. 11
- [47] Mikhail Mironov and Liudmila Prokhorenkova. Revisiting Graph Homophily Measures. *Learning on Graphs Conference*, 2024. 11, 12
- Shurui Gui, Xiner Li, Limei Wang, and Shuiwang Ji. GOOD: A Graph Out-of-Distribution Benchmark. *Advances in Neural Information Processing Systems*, 35:2059–2073, 2022. 12
- Gleb Bazhenov, Denis Kuznedelev, Andrey Malinin, Artem Babenko, and Liudmila Prokhorenkova. Evaluating Robustness and Uncertainty of Graph Models Under Structural Distributional Shifts. *Advances in Neural Information Processing Systems*, 36:75567–75594, 2023. 12
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 13
- [51] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint* arXiv:1607.06450, 2016. 13
- Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Revisiting Deep
   Learning Models for Tabular Data. Advances in Neural Information Processing Systems, 34:
   18932–18943, 2021. 13
- [53] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
   Message Passing for Quantum Chemistry. In *International Conference on Machine Learning*,
   pages 1263–1272. PMLR, 2017. 13
- [54] Dan Hendrycks and Kevin Gimpel. Gaussian Error Linear Units (GELUs). arXiv preprint
   arXiv:1606.08415, 2016. 13

- [55] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen.
   Graph Contrastive Learning with Augmentations. Advances in Neural Information Processing
   Systems, 33:5812–5823, 2020. 13
- [56] Jun Xia, Lirong Wu, Jintao Chen, Bozhen Hu, and Stan Z Li. SimGRACE: A Simple Framework
   for Graph Contrastive Learning without Data Augmentation. In *Proceedings of the ACM web* conference 2022, pages 1070–1079, 2022. 13
- [57] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014. URL http://jmlr.org/papers/v15/srivastava14a.html. 14
- [58] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 14
- [59] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
   Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An Imperative
   Style, High-Performance Deep Learning Library. Advances in Neural Information Processing
   Systems, 32, 2019. 14
- [60] Minjie Wang, Da Zheng, Zihao Ye, Quan Gan, Mufei Li, Xiang Song, Jinjing Zhou, Chao Ma,
   Lingfan Yu, Yu Gai, et al. Deep Graph Library: A Graph-Centric, Highly-Performant Package
   for Graph Neural Networks. arXiv preprint arXiv:1909.01315, 2019. 14

# A GraphLand Benchmark Details

#### A.1 Dataset descriptions

In this section, we provide a more detailed description of the GraphLand datasets. Note that none of the proposed datasets contain any personal information. For review purposes, our datasets can be found at this anonymous Kaggle storage, while data format description and our code can be found at this anonymous repository.

web-fraud, web-topics, and web-traffic These three datasets are web-graphs — they represent a segment of the Internet. The nodes are websites, and a directed edge connects two nodes if at least one user followed a link from one website to the other in a selected period of time. We prepared three datasets with the same graph but different tasks: in web-fraud, the task is to predict which websites are fraudulent (strongly imbalanced binary classification); in web-topics, the task is to predict the topic that a website belongs to (multiclass classification); and in web-traffic, the task is to predict how many users visited a website in a specific period of time (regression). With almost 3 million nodes, this is one of the largest publicly available attributed graphs that is not a citation network. Nodes in this graph have more than two hundred features, examples of which include the number of videos on the website (numerical feature), the website's zone and whether the website is on a free hosting (categorical features).

artnet-views and artnet-exp These two datasets represent a social network of art creators. The nodes are users, and an edge connects two nodes if the users are friends. We prepared two datasets with the same graph but different tasks: in artnet-views, the task is to predict how many views a user receives in a specific period of time (regression); and in artnet-exp, the task is to predict which users create explicit art content (binary classification). The examples of node features in this graph include user interests (categorical features).

city-roads-M and city-roads-L. These datasets are obtained from the logs of a navigation service and represent the road networks of two major cities, with the second one being several times larger than the first. The nodes are segments of roads, and a directed edge connects two nodes if the segments are incident to each other and moving from one segment to the other is permitted by traffic rules. The task is to predict the average travel speed on a road segment at a specific timestamp (regression). The features include various information about the road segment such as binary indicators of whether there is a bike dismount sign, whether the road segment ends with a crosswalk or a toll post, whether the road segment is in poor condition, whether it is restricted for trucks, and whether it has a mass transit lane (categorical features). The examples of numerical features are the length of the road and the geographic coordinates of the road endpoints.

**city-reviews** This dataset is obtained from the logs of a review service in which users can leave reviews and ratings for places and organizations in two major cities. The nodes are users, and an edge connects two nodes if the users often leave reviews for the same organizations. The graph is undirected. The task is fraud detection — to predict which users leave fraudulent reviews (binary classification). The node features are based on user interactions with the service and their examples include the share of negative reviews among all reviews the user has left (numerical feature) and the browser that is used to access the service by the user (categorical feature).

avazu-ctr This dataset is based on open data that has been introduced at the Kaggle competition organized by Avazu [38]. The data contains information about interactions between devices used to access the Internet, websites, and advertisements. In our graph, the nodes are devices, and an edge connects two nodes if the devices often visit the same websites. The graph is undirected. A smaller version of a similar dataset has been used by Ivanov and Prokhorenkova [39]; however, it contained only a small subset of devices, while for our dataset we collected data for all the available devices which makes our graph more than 50 times larger. The task is to predict the advertisement click-through rate (CTR) observed on devices (regression). Nodes in this graph have more than two hundred numerical features; however, most of them were anonymized in the original data source.

hm-categories and hm-prices These datasets are based on open data that has been introduced at the Kaggle competition organized by H&M [40]. The graph represents a co-purchasing network. The nodes are products, and an edge connects two nodes if the products are often bought by the same customers. The graph is undirected. We prepared two datasets with the same graph but different tasks: in hm-categories, the task is to predict the product category (multiclass classification), and in hm-prices, the task is to predict the product price (regression). The node features in this dataset

include product metadata such as product color (categorical feature), as well as information obtained from product purchasing statistics such as what proportion of product purchases occurs on different weekdays (numerical features).

**pokec-regions** This dataset is based on the data from Takac and Zabovsky [41]. It represents the online social network Pokec. The nodes are users, and a directed edge connects two nodes if one user has marked the other one as a friend. While this graph is quite popular in network analysis as an example of a classic social network, it is relatively rarely used for machine learning, with the exception of Lim et al. [42] who use the same graph as us but with different task, node features, and data split. In our dataset, the task is to predict which region a user is from (extreme multiclass classification with 183 classes). The node features in our dataset are based on user profile information, examples of them include the profile completion proportion (numerical feature) and binary indicators of whether different profile fields are filled (categorical features).

**twitch-views** This dataset is based on the data from Rozemberczki and Sarkar [43]. It represents the live-streaming network Twitch. The nodes are users, and an edge connects two nodes if both users follow each other. The task is to predict how many views a user gets in a specific period of time (regression). The node features are based on user profile information and examples of them include user language and affiliate status (categorical features).

tolokers-2 This is a new version of the dataset tolokers from Platonov et al. [8], Likhobaba et al. [44] with a significantly extended set of node features. It is based on the data from the Toloka crowdsourcing platform and the graph represents a network of tolokers (workers). The nodes are tolokers, and an edge connects two nodes if these tolokers have worked on the same task. The graph is undirected. The task is fraud detection — to predict which tolokers have been banned in one of the projects (binary classification). The new node features include various performance statistics of workers, such as the number of approved assignments and the number of skipped assignments (numerical features), as well as worker's profile information, such as their education level (categorical feature).

For all datasets, we provide random stratified RL and RH data splits. Further, we provide temporal TH data split (with the possibility of using the inductive learning setting THI) for all datasets with the exception of city-roads-M and city-roads-L datasets (since well-established road network graphs typically do not evolve over time significantly), as well as city-reviews and web-traffic datasets (since for them some of the necessary temporal information was not available).

## A.2 Dataset properties

410

411

412

413

415

416

417

418

419

423

424

425

426

427

430

431

432

438

439

440

442

443

444

445

446

450

451

452

453

454

455

458

459

461

A key characteristic of our benchmark is its diversity. As described above, our graphs come from different domains and have different prediction tasks. Their edges are also constructed in different ways (based on user interactions, activity similarity, physical connections, etc.). However, the proposed datasets also differ in many other ways. Some properties of our graphs are presented in Table 2 (see below for the details on how the provided characteristics are defined). First, note that the sizes of our datasets range from 11K to 3M nodes. The smaller graphs can be suitable for computeintensive models, while the larger graphs can provide a moderate scaling challenge. The average and median degrees of our graphs also vary significantly and our benchmark has both sparse and relatively dense graphs, including graphs with the average degree in the order of hundreds which is larger than the average degrees of most datasets used in current GML research (such graphs may highlight the importance of attention-based GNNs with their soft edge selection mechanisms). The average distance between two nodes in our graphs varies from 2.45 for hm-categories and hm-prices to 194 for city-roads-L; and graph diameter (maximum distance) varies from 8 for twitch-views to 553 for city-roads-L. Further, we report the values of clustering coefficients which show how typical closed node triplets are for the graph. In the literature, there are two definitions of clustering coefficients [45]: the global clustering coefficient and the average local clustering coefficient. We have both graphs where the clustering coefficients are high and graphs where they are almost zero, as well as graphs where global and local clustering coefficients significantly disagree (which is possible for graphs with imbalanced degree distributions). The degree assortativity coefficient is defined as the Pearson correlation coefficient of degrees among pairs of linked nodes. For most of our graphs, the degree assortativity is either negative or close to zero, which means that nodes do not tend to connect to other nodes with similar degrees, while city-roads-M and city-roads-L datasets are the exceptions — for them the degree assortativity is positive and large.

**Table 2:** Characteristics of the proposed GraphLand datasets.

		node classification				node regression								
	hm-categories	pokec-regions	web-topics	tolokers-2	city-reviews	artnet-exp	web-frand	hm-prices	avazu-ctr	city-roads-M	city-roads-L	twitch-views	artnet-views	web-traffic
# nodes	46.5K	1.6M	2.9M	11.8K	148.8K	50.4K	2.9M	46.5K	76.3K	57.1K	142.3K	168.1K	50.4K	2.9M
# edges	10.7M	22.3M	12.4M	519.0K	1.2M	280.3K	12.4M	10.7M	11.0M	107.1K	231.6K	6.8M	280.3K	12.4M
avg degree	460.92	27.32	8.56	88.28	15.66	11.12	8.56	460.92	288.04	3.75	3.26	80.87	11.12	8.56
median degree	45	13	2	30	4	2	2	45	71	4	3	32	2	2
avg distance	2.45	4.68	3.08	2.79	4.91	4.42	3.08	2.45	3.55	126.75	194.05	2.88	4.42	3.08
diameter	13	14	36	11	19	13	36	13	14	383	553	8	13	36
global clustering	0.27	0.05	0.00	0.23	0.26	0.03	0.00	0.27	0.24	0.00	0.00	0.02	0.03	0.00
avg local clustering	0.70	0.11	0.33	0.53	0.41	0.08	0.33	0.70	0.85	0.00	0.00	0.16	0.08	0.33
degree assortativity	-0.35	0.00	-0.14	-0.08	0.01	0.03	-0.14	-0.35	-0.30	0.70	0.74	-0.09	0.03	-0.14
# classes	21	183	28	2	2	2	2	N/A	N/A	N/A	N/A	N/A	N/A	N/A
unbiased homophily	0.38	0.98	0.55	0.10	0.69	0.28	0.32	N/A	N/A	N/A	N/A	N/A	N/A	N/A
target assortativity	N/A	N/A	N/A	N/A	N/A	N/A	N/A	0.12	0.18	0.74	0.72	-0.41	0.19	-0.21
# node features	35	56	263	16	37	75	266	41	260	26	26	4	50	267

464

465

466

467

469

470

471

472

473

475

478

479

480

481 482

483

485

486

487

488

489

491 492

493

494

495

496

497

498

500

Further, let us discuss the graph-label relationships in our datasets. To measure the similarity of labels of connected nodes for regression datasets, we use target assortativity — the Pearson correlation coefficient of target values between pairs of connected nodes. For instance, for the city-roads-M and city-roads-L datasets, the target assortativity is positive and quite large, which shows that nodes tend to connect to other nodes with similar target values (which is expected for the task of speed prediction in road networks), while for the twitch-views and web-traffic datasets, the target assortativity is negative. For classification datasets, the similarity of neighbors' labels is usually called *homophily*: in homophilous datasets, nodes tend to connect to nodes of the same class. How to properly measure homophily has recently attracted some research. It has been noted by Lim et al. [42] and Platonov et al. [46] that homophily measures typically used in the literature – such as the proportion of edges connecting nodes of the same class — are not appropriate for comparing homophily levels between graphs with different numbers of classes and their size balance. Platonov et al. [46] proposed a set of properties that a homophily measure appropriate for use in such comparisons should satisfy and Mironov and Prokhorenkova [47] constructed the first known homophily measure that satisfies all these properties — *unbiased homophily*. Thus, in our work, we use unbiased homophily to measure the homophily levels of our datasets. Unbiased homophily (with  $\alpha = 0$ , see Mironov and Prokhorenkova [47] for more details) takes values in [-1, 1] with 1 indicating perfect homophily, -1 indicating perfect heterophily, and 0 indicating no preference between homophilous and heterophilous edges (such graphs are typically referred to as heterophilous in the literature, although a more appropriate term would be non-homophilous). Note that the values of unbiased homophily should not be compared to values of other homophily measures used in the literature; the unbiased homophily levels for some popular graph node classification datasets are provided in Mironov and Prokhorenkova [47]. Unbiased homophily indicates that among our datasets pokec-regions and city-reviews are homophilous, while the other ones are non-homophilous. Thus, our benchmark significantly expands the set of available non-homophilous graph datasets.

Finally, our datasets have diverse sets of node features consisting of numerical and categorical features with different meanings and distributions. All our datasets except twitch-views and tolokers-2 have at least several dozen node features, while some have several hundred node features.

Overall, our datasets are diverse in domain, scale, structural properties, graph-label relations, and node attributes. Coming from real-world GML applications, they may serve as a valuable tool for the research and development of GML methods for the industry.

**Computing dataset characteristics.** Further, we describe the characteristics that are used in Table 2. Note that, while some graphs in our benchmark are directed, we transformed all the graphs to be undirected before computing all the considered graph characteristics, since some of the characteristics are not defined for directed graphs.

Average degree and median degree are the average and median numbers of neighbors a node has, respectively. Since all our graphs are connected (when treated as undirected graphs), for any two nodes there is a path between them. Average distance is the average length of the shortest paths between all pairs of nodes, while diameter is the maximum length of the shortest paths between all pairs of nodes.

For our largest graphs — the ones used for the pokec-regions, web-traffic, web-fraud, and web-topics datasets — we approximate average distance with an average over distances for  $100 \, \mathrm{K}$  randomly sampled node pairs. Global clustering coefficient is computed as the tripled number of triangles divided by the number of pairs of adjacent edges (i.e., it is the fraction of closed triplets of nodes among all connected triplets). Average local clustering coefficient first computes the local clustering of each node, which is the fraction of connected pairs of its neighbors, and then averages the obtained values among all nodes. Degree assortativity is the Pearson correlation coefficient between the degrees of connected nodes. Further, target assortativity for regression datasets is the Pearson correlation coefficient between target values of connected nodes. For computing unbiased homophily, we follow Mironov and Prokhorenkova [47] and use the simplest version of this measure with the  $\alpha$  parameter set to 0.

# **B** Data Splits and Experimental Settings

503

504

505

506

507

508

510

511

512

514

515

516

518

519

521

523

524

530

531

532

533

534

535

538

539

540

545

547

548

549

550

553

In GML literature, there are two most popular settings for node property prediction regarding the relative sizes of train, validation, and test sets: one with a high label rate and one with a low label rate. In the high label rate setting, the train set encompasses 50% of all graph nodes or more. This setting is common in heterophilous benchmarks, e.g., it is used by datasets from Pei et al. [21] and Platonov et al. [8], and it is also used by most datasets from OGB [19]. In the low label rate setting, much smaller train set sizes are used (typically, no more than 10% of all graph nodes). This setting is commonly used with the classic cora, citeseer, pubmed citation networks, and it is also used by the ogbn-products dataset from OGB. Both of these settings can appear in real-world GML usage scenarios, depending on the resources available for data labeling. It is important to provide predetermined splits to ensure experiments in different works are run in the same setting and their results are comparable, but it is also important to accommodate different needs of different research projects. Thus, for datasets in our benchmark we provide fixed splits for both settings. We refer to these splits as the RL (random low) and RH (random high) splits. Specifically, the RL split randomly divides nodes into train/validation/test sets with 10%/10%/80% proportions, while the RH split randomly divides nodes into train/validation/test sets with 50%/25%/25% proportions.

The RL and RH data splits are random, as is common in current GML benchmarks. However, in real-world applications, data splits are often temporal, i.e., the labeled objects are the ones that appeared in the network earlier, while the ones that appeared later are not labeled and belong to the test set. Despite their prevalence in applications, temporal splits are very rarely used in current GML node property prediction benchmarks. To the best of our knowledge, the only such datasets with temporal splits available are citation networks from OGB, which represent only a single application. At the same time, temporal data splits may significantly affect the prediction problem and the model performance, as they often result in distributional shifts between train, validation, and test data. While some types of distributional shifts have been previously explored in the GML literature, e.g., shifts in node features [48] and shifts in graph structure [49], realistic temporal distributional shifts often combine shifts in several aspects of data simultaneously (e.g., shifts in the distributions of node features, labels, and graph structural characteristics), and the effect of such realistic shifts on GML model performance is currently under-explored. To close this gap, we provide a temporal split for most datasets in our benchmark. We refer to this split as the TH (temporal high) split; it divides nodes into train/validation/test sets with 50%/25%/25% proportions, i.e., exactly the same proportions as in the RH split, which allows comparing model results between the RH and TH splits to see how the complexity of the task changes when temporal distributional shifts are introduced.

Further, many real-world networks are not static, but evolve over time. Thus, in many applications, not only are there no labels available for nodes that appear in the network later, but the nodes themselves (with their attributes and incident edges) are not available at training time. This setting is known in GML as the *inductive* setting. In contrast to the *transductive* setting in which the whole graph is available at training time (including the nodes for which predictions should be made), in the inductive setting validation and test nodes are not available at training time. Despite temporally evolving graphs being common in practical applications, most standard node property prediction datasets only provide the transductive setting. While it is well-known that GNNs, in contrast to some other GML methods like shallow node embeddings, can work not only in the transductive but also in the inductive setting [32], it is not well-explored how the lack of complete graph information at training time in the inductive setting affects GNN performance. Moreover, when the inductive setting is used for GNN evaluation in the current literature, the validation and train nodes are typically

chosen randomly, which is not realistic, since in real-world applications the inductive setting is almost always induced by the temporal evolution of the graph. To fill this gap, for all datasets in our benchmark for which temporal information is available, we additionally provide the inductive experimental setting. We refer to this setting as THI (temporal high / inductive); it has the exact same data split as the transductive TH setting, but provides three snapshots of the graph: one for training, one for validation, and one for testing. This allows investigating how model performance changes between the transductive and the inductive setting. To the best of our knowledge, our work is the first to compare GNN performance between random and temporal data splits in the transductive setting under the same split ratios, and between transductive and inductive settings under the same temporal data split. This comparison allows us to investigate how much these differences can affect GNN performance.

# C Experimental Setup Details

#### C.1 Models

560

561

562

563

566

567

568

569

570

571

580

581

582

583

584

585

588

589

590

598

599

600

601

602

603

605

606

607

A graph-agnostic baseline. As a simple baseline, we use ResMLP — an MLP with skip-connections [50] and layer normalization [51]. This model does not have any information about the graph structure and operates on nodes as independent samples — we call such models *graph-agnostic*. It has been shown that such MLP-like models with skip-connections can serve as very strong baselines for industrial data with mixed numerical and categorical features [52].

**Graph neural networks.** We consider several representative GNN architectures. First, we use GCN [17] and GraphSAGE [32] as simple classical GNN models. For GraphSAGE, we use the version with the mean aggregation function, and we do not use the neighbor sampling technique proposed in the original paper, instead training the model on the full graph, like all other GNNs in our experiments. Further, we use two GNNs with attention-based neighborhood aggregation functions: GAT [33] and Graph Transformer (GT) [34]. Note that GT is a *local* graph transformer, i.e., each node only attends to its neighbors in the graph (in contrast to global graph transformers, in which each node attends to all other nodes in the graph, and which are thus not instances of the standard message-passing neural networks (MPNNs) framework of Gilmer et al. [53]). Following Platonov et al. [8], we equip all the considered GNNs with skip-connections and layer normalization, which we found important for their strong performance on our datasets. We also add a two-layer MLP with the GELU activation function [54] after every neighborhood aggregation block in GNNs. Our graph models are implemented in the same codebase as our ResMLP — we simply swap each residual block of ResMLP with a residual neighborhood aggregation block of the selected GNN architecture. Therefore, comparing the performance of ResMLP and GNNs allows us to see if graph information is helpful for the task. Indeed, in our experiments, GNNs significantly outperform graph-agnostic ResMLP on all our datasets, confirming the usefulness of the provided graph structure for the considered tasks.

**Graph foundation models.** Most currently available GFMs do not support node property prediction tasks in graphs with arbitrary node features. Of those that do, we were able to find only two models with open weights: OpenGraph [35] and AnyGraph [36]. Both OpenGraph and AnyGraph exploit the Transformer architecture and are pretrained with a link prediction objective on a mixture of different graph datasets. These methods differ in what data they can operate on. Specifically, OpenGraph only uses relational information and constructs node representations based on SVD factors of the adjacency matrix, while AnyGraph also uses the available node feature information and combines SVD factors for both the feature matrix and the adjacency matrix. Both these models were designed to be adapted to new node classification datasets without fine-tuning, using an in-context learning (ICL) setting instead. Specifically, they can perform link prediction in arbitrary graphs, and they cast any node classification task as a link prediction task where links to virtual nodes representing classes are predicted for unlabeled nodes. Further, we were able to reproduce the pretraining for one more GFM — GCOPE [37] (weights for this model are not publicly available, but training code is). GCOPE also applies a projection to node features (e.g., based on SVD or attention mechanism), but exploits additional virtual nodes as graph coordinators to simultaneously process different graph datasets at the pretraining stage. The authors use GraphCL [55] or SimGRACE [56] as the pretraining objective. In contrast to OpenGraph and AnyGraph, GCOPE uses fine-tuning for adaptation to new

node classification datasets. However, neither of the three models supports node regression. Further, we found that these GFMs cannot scale to large datasets.

### C.2 Experimental setup and hyperparameter selection details

614

626

627

628

636

637

638

641

642

648

649

650

651

652

655

658

Some of the graphs in our benchmark are directed. For our experiments, we converted directed graphs to undirected ones (by replacing each directed edge with an undirected one and then removing duplicated edges). We leave investigation of different ways to consider edge directions to further research.

We train all models 10 times with different random seeds to compute the mean and standard deviation of model performance, except for our largest datasets pokec-regions, web-traffic, web-fraud, web-topics, for which we train all models 5 times.

We train all our GNNs in a full-batch setting, i.e., we do not use any subgraph sampling techniques and train the models on the full graph. Our ResMLP baseline is implemented in the same codebase as our GNNs and thus is also trained in the full-batch setting.

Hyperparameter choice is extremely important for the performance of GNNs. Thus, we conducted hyperparameter search on the validation set for all models. Specifically, we found that the learning rate and dropout probability [57] are the most important hyperparameters for our GNN implementations on our datasets. Thus, we ran grid search selecting the learning rate from  $\{3 \times 10^{-5}, 1 \times 10^{-4}, 3 \times 10^{-5}, 1 \times 10^{-5}, 1 \times 10^{-4}, 3 \times 10^{-5}, 1 \times 1$  $10^{-4}$ ,  $1 \times 10^{-3}$ ,  $3 \times 10^{-3}$ } and dropout probability from  $\{0, 0.1, 0.2\}$  (note that the highest learning rate of  $3 \times 10^{-3}$  often resulted in NaN issues, however, we still included it in our hyperparameter search, as in our preliminary experiments we found it to be beneficial for some of our dataset/model combinations). In our preliminary experiments we found that the performance of our GNNs is quite stable for a wide variety of reasonable architecture hyperparameter values (we found the use of skip-connections and layer normalization to be important for this stability). Hence, for our final experiments, we kept these hyperparameters fixed. We set these values as follows: the number of graph neighborhood aggregation blocks to 3 and the hidden dimension to 512. The only exceptions to this hidden dimension size were made for our largest datasets: to avoid GPU out-of-memory issues, we decreased the hidden dimension to 400 for pokec-regions and to 200 for web-traffic, web-fraud, and web-topics. For GNNs with attention-based graph neighborhood aggregation (GAT and GT), the number of attention heads was set to 4. We used the Adam optimizer [58] in all our GNN experiments. We trained each model for 1000 steps and then selected the best step based on the performance on the validation set.

When applying deep learning models to data with numerical features, the preprocessing of these features is critically important. In our experiments, we considered two possible numerical feature transformation techniques: standard scaling and quantile transformation to standard normal distribution. We included them in the hyperparameter search for ResMLP and GNNs. In contrast, GBDT models do not need specialized preprocessing for numerical features and are not affected by their monotonic transformations. For categorical features, we used one-hot encoding for all models except for LightGBM and CatBoost, which support the use of categorical features directly and have their specialized strategies for working with them (XGBoost also offers such a feature, but it is currently marked as experimental, and we were not able to make it work). For regression datasets, neural models might perform better if the target variable is transformed. Therefore, in our experiments on regression datasets with ResMLP and GNNs, we considered the options of using the original targets or preprocessing targets with standard scaling, including these two options in the hyperparameter search.

Our GNNs are implemented using PyTorch [59] and DGL [60].

# **D** Additional Experimental Results

In the main text, we report our experimental results for the RL setting. In this section, we present results for other settings. In Table 3 we report combined results for the RH, TH, and THI settings for those datasets that have temporal data split. Additionally, in Tables 4, 5, 6 we separately report results for the RH, TH, and THI settings, respectively. In all tables with results, for each dataset, we highlight with color the best result as well as those results for which the mean differs from the best one by no more than the sum of the two results' standard deviations.

**Table 3:** Experimental results under the RH (random high), TH (temporal high), and THI (temporal high / inductive) settings. The best result and those statistically indistinguishable from it are highlighted in red for RH, violet for TH, and blue for THI.

		mul	ticlass classification	binary classification			
		hm-categories	pokec-regions	web-topics	tolokers-2	artnet-exp	web-fraud
best const. pred.	RH TH THI	$\begin{array}{c} 19.46 \pm 0.00 \\ 19.57 \pm 0.00 \\ 19.57 \pm 0.00 \end{array}$	$3.77 \pm 0.00$ $2.98 \pm 0.00$ $2.98 \pm 0.00$	$28.36 \pm 0.00$ $23.28 \pm 0.00$ $23.28 \pm 0.00$	$21.82 \pm 0.00$ $8.61 \pm 0.00$ $8.61 \pm 0.00$	$10.00 \pm 0.00 7.84 \pm 0.00 7.84 \pm 0.00$	$0.66 \pm 0.00$ $0.15 \pm 0.00$ $0.15 \pm 0.00$
ResMLP	RH TH THI	$\begin{array}{c} 43.12 \pm 0.25 \\ 32.44 \pm 0.54 \\ 32.44 \pm 0.54 \end{array}$	$\begin{array}{c} 5.09 \pm 0.01 \\ 4.18 \pm 0.01 \\ 4.18 \pm 0.01 \end{array}$	$\begin{array}{c} 44.55 \pm 0.08 \\ 35.49 \pm 0.03 \\ 35.49 \pm 0.03 \end{array}$	$\begin{array}{c} 45.96 \pm 0.46 \\ 21.72 \pm 6.69 \\ 21.72 \pm 6.69 \end{array}$	$\begin{array}{c} 43.55 \pm 0.23 \\ 37.48 \pm 0.51 \\ 37.48 \pm 0.51 \end{array}$	$13.52 \pm 0.21 2.83 \pm 0.26 2.83 \pm 0.26$
GCN	RH TH THI	$73.38 \pm 0.42 56.91 \pm 0.55 47.05 \pm 1.69$	$\begin{array}{c} 35.08 \pm 0.62 \\ 11.88 \pm 0.31 \\ 6.88 \pm 0.51 \end{array}$	$\begin{array}{c} 48.88 \pm 0.09 \\ 38.20 \pm 0.19 \\ 37.76 \pm 0.06 \end{array}$	$60.49 \pm 0.86$ $46.72 \pm 1.19$ $32.43 \pm 8.03$	$49.80 \pm 0.35 40.64 \pm 0.40 41.28 \pm 0.28$	$15.58 \pm 0.20  4.85 \pm 1.42  3.44 \pm 0.33$
GraphSAGE	RH TH THI	$73.34 \pm 0.68  59.62 \pm 0.51  48.11 \pm 2.08$	$\begin{array}{c} 40.76 \pm 0.21 \\ 16.60 \pm 0.28 \\ 8.04 \pm 0.26 \end{array}$	$\begin{array}{c} 50.05 \pm 0.03 \\ 39.00 \pm 0.09 \\ 38.04 \pm 0.18 \end{array}$	$58.42 \pm 0.92$ $17.05 \pm 7.65$ $30.86 \pm 9.48$	$\begin{array}{c} 48.49 \pm 0.37 \\ 40.50 \pm 0.84 \\ 40.53 \pm 0.40 \end{array}$	$20.47 \pm 0.17$ $16.01 \pm 2.22$ $13.88 \pm 1.32$
GAT	RH TH THI	$79.19 \pm 0.21$ $61.28 \pm 0.97$ $59.34 \pm 1.09$	$46.72 \pm 0.69 20.43 \pm 0.55 13.38 \pm 0.35$	$50.54 \pm 0.04$ $39.24 \pm 0.23$ $38.77 \pm 0.35$	$63.76 \pm 1.30$ $38.59 \pm 6.19$ $24.53 \pm 9.55$	$50.62 \pm 0.35$ $41.85 \pm 0.63$ $41.64 \pm 0.32$	$20.43 \pm 0.21$ $16.50 \pm 1.14$ $11.98 \pm 1.54$
GT	RH TH THI	$79.28 \pm 0.31$ $63.31 \pm 0.45$ $59.54 \pm 1.59$	$50.06 \pm 0.53$ $25.09 \pm 0.58$ $17.22 \pm 0.42$	$50.58 \pm 0.04$ $39.19 \pm 0.15$ $38.78 \pm 0.08$	$60.32 \pm 1.21 \\ 34.15 \pm 4.81 \\ 22.89 \pm 10.4$	$49.32 \pm 1.00$ $40.10 \pm 0.60$ $40.26 \pm 0.82$	$19.73 \pm 0.34$ $11.97 \pm 1.13$ $7.84 \pm 2.35$
OpenGraph (ICL)	RH TH THI	$11.69 \pm 0.84$ $5.76 \pm 1.03$ $5.76 \pm 1.03$	$2.56 \pm 0.42$ $0.80 \pm 0.45$ $0.80 \pm 0.45$	RTE RTE RTE	$44.62 \pm 1.35 9.12 \pm 1.74 9.12 \pm 1.74$	$23.72 \pm 1.86$ $16.19 \pm 1.36$ $16.19 \pm 1.36$	RTE RTE RTE
AnyGraph (ICL)	RH TH THI	$\begin{array}{c} 15.65 \pm 2.82 \\ 9.47 \pm 1.13 \\ 9.47 \pm 1.13 \end{array}$	$27.67 \pm 2.48$ $9.20 \pm 0.67$ $9.20 \pm 0.67$	$\begin{array}{c} 6.30 \pm 2.82 \\ 11.14 \pm 5.16 \\ 11.14 \pm 5.16 \end{array}$	$30.21 \pm 3.32$ $13.52 \pm 4.74$ $13.52 \pm 4.74$	$\begin{array}{c} 15.80 \pm 1.90 \\ 11.80 \pm 1.01 \\ 11.80 \pm 1.01 \end{array}$	$0.67 \pm 0.02$ $0.16 \pm 0.01$ $0.16 \pm 0.01$
GCOPE (FT)	RH TH THI	$\begin{array}{c} 19.99 \pm 0.12 \\ 19.14 \pm 0.58 \\ 15.69 \pm 3.44 \end{array}$	TLE TLE TLE	TLE TLE TLE	$31.79 \pm 1.95$ $8.46 \pm 1.15$ $10.73 \pm 1.48$	$\begin{array}{c} 23.86 \pm 2.33 \\ 20.83 \pm 1.18 \\ 19.10 \pm 0.96 \end{array}$	TLE TLE TLE

(b) Results for regression datasets.  $R^2$  is reported for all datasets.

		hm-prices	avazu-ctr	twitch-views	artnet-views
best const. pred.	RH TH THI	$0.00 \pm 0.00$ $-2.85 \pm 0.00$ $-2.85 \pm 0.00$	$0.00 \pm 0.00$ $0.00 \pm 0.00$ $0.00 \pm 0.00$	$\begin{array}{c} 0.00 \pm 0.00 \\ -22.31 \pm 0.00 \\ -22.31 \pm 0.00 \end{array}$	$0.00 \pm 0.00$ $-9.32 \pm 0.00$ $-9.32 \pm 0.00$
ResMLP	RH TH THI	$\begin{array}{c} 70.11 \pm 0.48 \\ 61.64 \pm 0.79 \\ 61.64 \pm 0.79 \end{array}$	$28.03 \pm 0.22$ $20.35 \pm 1.50$ $20.35 \pm 1.50$	$13.36 \pm 0.01$ $11.91 \pm 8.00$ $11.91 \pm 8.00$	$36.10 \pm 0.17$ $42.15 \pm 0.52$ $42.15 \pm 0.52$
GCN	RH TH THI	$79.76 \pm 0.76$ $65.20 \pm 0.84$ $64.31 \pm 0.82$	$34.96 \pm 0.11$ $37.49 \pm 0.26$ $34.78 \pm 0.48$	$77.12 \pm 0.11$ $68.17 \pm 0.24$ $63.58 \pm 0.54$	$61.02 \pm 0.13$ $54.44 \pm 0.43$ $53.73 \pm 0.47$
GraphSAGE	RH TH THI	$79.89 \pm 0.46 67.93 \pm 1.24 65.80 \pm 0.56$	$35.20 \pm 0.20$ $38.38 \pm 0.39$ $36.79 \pm 0.55$	$72.02 \pm 0.16 \\ 61.46 \pm 0.68 \\ 56.60 \pm 0.71$	$56.65 \pm 0.50$ $51.87 \pm 0.48$ $53.37 \pm 0.30$
GAT	RH TH THI	$81.68 \pm 0.41$ $70.83 \pm 0.99$ $69.74 \pm 1.50$	$35.74 \pm 0.19$ $39.21 \pm 0.17$ $37.18 \pm 1.02$	$76.06 \pm 0.30$ $66.32 \pm 0.59$ $61.41 \pm 1.30$	$\begin{array}{c} 59.01 \pm 0.52 \\ 52.30 \pm 0.34 \\ 52.44 \pm 0.59 \end{array}$
GT	RH TH THI	$80.90 \pm 0.42$ $69.70 \pm 0.84$ $67.33 \pm 2.05$	$34.38 \pm 0.31$ $38.27 \pm 0.27$ $36.49 \pm 0.83$	$75.57 \pm 0.15  65.83 \pm 0.24  60.67 \pm 1.02$	$58.97 \pm 0.25$ $51.67 \pm 0.54$ $52.26 \pm 0.48$

First, we notice that the observations from the results for the low label rate RL setting about the usefulness of graph structure, strong performance of classic and attention-based GNNs, and weak performance of GFMs also apply to high label rate settings. Next, we observe that temporal data splits are significantly more challenging for all models than random ones (with the exception of the avazu-ctr dataset). This is important as in real-world applications temporal distributional shifts are common, and not considering them can provide overly optimistic performance estimates. Further, the considered models perform significantly worse in the inductive setting than in the transductive one. These observations highlight the importance of developing GML methods that are more resilient to temporal distributional shifts and dynamic changes in the graph structure for industrial applications. In the absence of such methods, frequent retraining of GNNs on new data is recommended to achieve the best results. Note that GFMs that only utilize in-context learning to adapt to new graphs do not suffer from the transductive/inductive mismatch and thus represent a promising direction, but their performance is currently very weak compared with GNNs on all datasets.

**Table 4:** Experimental results under the RH (random high) data split in the transductive setting. The best result and those statistically indistinguishable from it are highlighted in red. TLE stands for time limit exceeded (24 hours); RTE stands for runtime error in the official code of GFMs.

	mu	lticlass classification	1	binary classification				
	hm-categories	pokec-regions	web-topics	tolokers-2	city-reviews	artnet-exp	web-fraud	
best const. pred. ResMLP	$19.46 \pm 0.00 \\ 43.12 \pm 0.25$	$3.77 \pm 0.00$ $5.09 \pm 0.01$	$28.36 \pm 0.00$ $44.55 \pm 0.08$	$\begin{array}{c} 21.82 \pm 0.00 \\ 45.96 \pm 0.46 \end{array}$	$12.09 \pm 0.00 75.21 \pm 0.08$	$10.00 \pm 0.00 43.55 \pm 0.23$	$0.66 \pm 0.00$ $13.52 \pm 0.21$	
GCN GraphSAGE GAT GT	$73.38 \pm 0.42$ $73.34 \pm 0.68$ $79.19 \pm 0.21$ $79.28 \pm 0.31$	$35.08 \pm 0.62$ $40.76 \pm 0.21$ $46.72 \pm 0.69$ $50.06 \pm 0.53$	$48.88 \pm 0.09 50.05 \pm 0.03 50.54 \pm 0.04 50.58 \pm 0.04$	$60.49 \pm 0.86 58.42 \pm 0.92 63.76 \pm 1.30 60.32 \pm 1.21$	$81.05 \pm 0.10$ $80.75 \pm 0.06$ $81.10 \pm 0.11$ $80.50 \pm 0.14$	$49.80 \pm 0.35$ $48.49 \pm 0.37$ $50.62 \pm 0.35$ $49.32 \pm 1.00$	$15.58 \pm 0.20  20.47 \pm 0.17  20.43 \pm 0.21  19.73 \pm 0.34$	
OpenGraph (ICL) AnyGraph (ICL) GCOPE (FT)	$11.69 \pm 0.84 15.65 \pm 2.82 19.99 \pm 0.12$	$\begin{array}{c} 2.56 \pm 0.42 \\ 27.67 \pm 2.48 \\ \text{TLE} \end{array}$	$\begin{array}{c} \text{RTE} \\ 6.30 \pm 2.82 \\ \text{TLE} \end{array}$	$\begin{array}{c} 44.62 \pm 1.35 \\ 30.21 \pm 3.32 \\ 31.79 \pm 1.95 \end{array}$	$62.96 \pm 0.84 \\ 65.04 \pm 1.41 \\ 69.74 \pm 0.36$	$\begin{array}{c} 23.72 \pm 1.86 \\ 15.80 \pm 1.90 \\ 23.86 \pm 2.33 \end{array}$	$\begin{array}{c} \text{RTE} \\ 0.67 \pm 0.02 \\ \text{TLE} \end{array}$	

## (b) Results for regression datasets. $R^2$ is reported for all datasets.

	hm-prices	avazu-ctr	city-roads-M	city-roads-L	twitch-views	artnet-views	web-traffic
best const. pred. ResMLP	$\begin{array}{c} 0.00 \pm 0.00 \\ 70.11 \pm 0.48 \end{array}$	$0.00 \pm 0.00$ $28.03 \pm 0.22$	$0.00 \pm 0.00$ $62.43 \pm 0.32$	$0.00 \pm 0.00$ $53.09 \pm 0.17$	$0.00 \pm 0.00$ $13.36 \pm 0.01$	$\begin{array}{c} 0.00 \pm 0.00 \\ 36.10 \pm 0.17 \end{array}$	$0.00 \pm 0.00$ $73.88 \pm 0.05$
GCN GraphSAGE GAT GT	$79.89 \pm 0.46$ $81.68 \pm 0.41$	$34.96 \pm 0.11 35.20 \pm 0.20 35.74 \pm 0.19 34.38 \pm 0.31$	$69.95 \pm 0.11$ $70.20 \pm 0.59$ $70.53 \pm 0.40$ $67.45 \pm 0.82$	$64.65 \pm 0.27$ $65.77 \pm 0.43$ $66.03 \pm 0.24$ $64.02 \pm 0.59$	$77.12 \pm 0.11$ $72.02 \pm 0.16$ $76.06 \pm 0.30$ $75.57 \pm 0.15$	$61.02 \pm 0.13$ $56.65 \pm 0.50$ $59.01 \pm 0.52$ $58.97 \pm 0.25$	$83.49 \pm 0.14$ $85.19 \pm 0.11$ $85.70 \pm 0.08$ $85.54 \pm 0.23$

# **E** Limitations and Broader Impact

The aim of our benchmark is to introduce a diverse set of graph datasets for node property prediction that covers a wide range of domains and graph structural properties, including those not encountered in commonly used datasets for GML model evaluation. However, data that can be naturally represented as graphs is so widespread across different domains that no benchmark can cover them all. Thus, our collection of 14 datasets still only covers a small part of the wide range of situations where modeling data as a graph can be useful. But we hope that it will encourage the GML research community to use more diverse sets of datasets and focus on practically relevant applications where graph-structured data appears.

Our benchmark includes datasets with realistic tasks such as fraud detection and user engagement prediction. Poorly performing machine learning models used for these tasks in real-world services can negatively affect the users of these services. For example, type I errors of fraud detection systems, i.e., wrongly predicting that an innocent person is fraudulent, have an undesirable negative impact. Thus, particular care should be taken to minimize the probability of such errors. We believe that the release of high-quality and properly anonymized datasets for these tasks such as the ones in our benchmark will encourage the community to develop better models, since the community will be able to use these datasets as a realistic and reliable testbed to investigate which methods lead to reductions in undesirable model errors.

**Table 5:** Experimental results under the TH (temporal high) data split in the transductive setting. The best result and those statistically indistinguishable from it are highlighted in violet. TLE stands for time limit exceeded (24 hours); RTE stands for runtime error in the official code of GFMs.

	mul	lticlass classification	binary classification			
	hm-categories	pokec-regions	web-topics	tolokers-2	artnet-exp	web-fraud
best const. pred. ResMLP	$19.57 \pm 0.00 \\ 32.44 \pm 0.54$	$2.98 \pm 0.00$ $4.18 \pm 0.01$	$23.28 \pm 0.00$ $35.49 \pm 0.03$	$8.61 \pm 0.00$ $21.72 \pm 6.69$	$7.84 \pm 0.00$ $37.48 \pm 0.51$	$0.15 \pm 0.00$ $2.83 \pm 0.26$
GCN GraphSAGE GAT GT	$\begin{array}{c} 56.91 \pm 0.55 \\ 59.62 \pm 0.51 \\ 61.28 \pm 0.97 \\ 63.31 \pm 0.45 \end{array}$	$11.88 \pm 0.31$ $16.60 \pm 0.28$ $20.43 \pm 0.55$ $25.09 \pm 0.58$	$38.20 \pm 0.19$ $39.00 \pm 0.09$ $39.24 \pm 0.23$ $39.19 \pm 0.15$	$46.72 \pm 1.19$ $17.05 \pm 7.65$ $38.59 \pm 6.19$ $34.15 \pm 4.81$	$40.64 \pm 0.40  40.50 \pm 0.84  41.85 \pm 0.63  40.10 \pm 0.60$	$4.85 \pm 1.42$ $16.01 \pm 2.22$ $16.50 \pm 1.14$ $11.97 \pm 1.13$
OpenGraph (ICL) AnyGraph (ICL) GCOPE (FT)	$\begin{array}{c} 5.76 \pm 1.03 \\ 9.47 \pm 1.13 \\ 19.14 \pm 0.58 \end{array}$	$\begin{array}{c} 0.80 \pm 0.45 \\ 9.20 \pm 0.67 \\ \text{TLE} \end{array}$	$\begin{array}{c} \text{RTE} \\ 11.14 \pm 5.16 \\ \text{TLE} \end{array}$	$\begin{array}{c} 9.12 \pm 1.74 \\ 13.52 \pm 4.74 \\ 8.46 \pm 1.15 \end{array}$	$16.19 \pm 1.36$ $11.80 \pm 1.01$ $20.83 \pm 1.18$	$\begin{array}{c} \text{RTE} \\ 0.16 \pm 0.01 \\ \text{TLE} \end{array}$

(b) Results for regression datasets.  $R^2$  is reported for all datasets.

	hm-prices	avazu-ctr	twitch-views	artnet-views
best const. pred. ResMLP	$-2.85 \pm 0.00$ $61.64 \pm 0.79$	$\begin{array}{c} 0.00 \pm 0.00 \\ 20.35 \pm 1.50 \end{array}$	$-22.31 \pm 0.00$ $11.91 \pm 8.00$	$-9.32 \pm 0.00 42.15 \pm 0.52$
GCN GraphSAGE GAT GT	$65.20 \pm 0.84$ $67.93 \pm 1.24$ $70.83 \pm 0.99$ $69.70 \pm 0.84$	$37.49 \pm 0.26$ $38.38 \pm 0.39$ $39.21 \pm 0.17$ $38.27 \pm 0.27$	$68.17 \pm 0.24$ $61.46 \pm 0.68$ $66.32 \pm 0.59$ $65.83 \pm 0.24$	$\begin{array}{c} 54.44 \pm 0.43 \\ 51.87 \pm 0.48 \\ 52.30 \pm 0.34 \\ 51.67 \pm 0.54 \end{array}$

**Table 6:** Experimental results under the THI (temporal high / inductive) setting. The best result and those statistically indistinguishable from it are highlighted in blue. TLE stands for time limit exceeded (24 hours); RTE stands for runtime error in the official code of GFMs.

(a) Results for classification datasets. Accuracy is reported for multiclass classification datasets and Average Precision is reported for binary classification datasets.

	mu	lticlass classification	binary classification			
	hm-categories	pokec-regions	web-topics	tolokers-2	artnet-exp	web-fraud
best const. pred. ResMLP	$19.57 \pm 0.00 \\ 32.44 \pm 0.54$	$2.98 \pm 0.00$ $4.18 \pm 0.01$	$23.28 \pm 0.00$ $35.49 \pm 0.03$	$8.61 \pm 0.00$ $21.72 \pm 6.69$	$7.84 \pm 0.00$ $37.48 \pm 0.51$	$0.15 \pm 0.00$ $2.83 \pm 0.26$
GCN GraphSAGE GAT GT	$47.05 \pm 1.69$ $48.11 \pm 2.08$ $59.34 \pm 1.09$ $59.54 \pm 1.59$	$6.88 \pm 0.51$ $8.04 \pm 0.26$ $13.38 \pm 0.35$ $17.22 \pm 0.42$	$37.76 \pm 0.06$ $38.04 \pm 0.18$ $38.77 \pm 0.35$ $38.78 \pm 0.08$	$32.43 \pm 8.03$ $30.86 \pm 9.48$ $24.53 \pm 9.55$ $22.89 \pm 10.40$	$41.28 \pm 0.28  40.53 \pm 0.40  41.64 \pm 0.32  40.26 \pm 0.82$	$3.44 \pm 0.33$ $13.88 \pm 1.32$ $11.98 \pm 1.54$ $7.84 \pm 2.35$
OpenGraph (ICL) AnyGraph (ICL) GCOPE (FT)	$5.76 \pm 1.03$ $9.47 \pm 1.13$ $15.69 \pm 3.44$	$\begin{array}{c} 0.80 \pm 0.45 \\ 9.20 \pm 0.67 \\ \text{TLE} \end{array}$	$\begin{array}{c} \text{RTE} \\ 11.14 \pm 5.16 \\ \text{TLE} \end{array}$	$9.12 \pm 1.74$ $13.52 \pm 4.74$ $10.73 \pm 1.48$	$16.19 \pm 1.36$ $11.80 \pm 1.01$ $19.10 \pm 0.96$	$\begin{array}{c} \text{RTE} \\ 0.16 \pm 0.01 \\ \text{TLE} \end{array}$

(b) Results for regression datasets.  $R^2$  is reported for all datasets.

	hm-prices	avazu-ctr	twitch-views	artnet-views
best const. pred. ResMLP	$-2.85 \pm 0.00$ $61.64 \pm 0.79$	$\begin{array}{c} 0.00 \pm 0.00 \\ 20.35 \pm 1.50 \end{array}$	$-22.31 \pm 0.00$ $11.91 \pm 8.00$	$-9.32 \pm 0.00 42.15 \pm 0.52$
GCN GraphSAGE GAT GT	$64.31 \pm 0.82$ $65.80 \pm 0.56$ $69.74 \pm 1.50$ $67.33 \pm 2.05$	$34.78 \pm 0.48$ $36.79 \pm 0.55$ $37.18 \pm 1.02$ $36.49 \pm 0.83$	$63.58 \pm 0.54$ $56.60 \pm 0.71$ $61.41 \pm 1.30$ $60.67 \pm 1.02$	$53.73 \pm 0.47$ $53.37 \pm 0.30$ $52.44 \pm 0.59$ $52.26 \pm 0.48$