# GRAPHFM: A GENERALIST GRAPH TRANSFORMER THAT LEARNS TRANSFERABLE REPRESENTATIONS ACROSS DIVERSE DOMAINS

## Anonymous authors

Paper under double-blind review

## ABSTRACT

Graph neural networks (GNNs) are often trained on individual datasets, requiring specialized models and significant hyperparameter tuning due to the unique structures and features of each dataset. This approach limits the scalability and generalizability of GNNs, as models must be tailored for each specific graph type. To address these challenges, we introduce GRAPHFM, a scalable multi-graph pretraining approach designed for learning across diverse graph datasets. GRAPHFM uses a Perceiver-based encoder with learned latent tokens to compress domain-specific features into a shared latent space, enabling generalization across graph domains. We propose new techniques for scaling up graph training on datasets of different sizes, allowing us to train GRAPHFM on 152 distinct graph datasets, spanning 7.4 million nodes and 189 million edges. This allows us to study the effect of scale on pretraining across domains such as molecules, citation networks, and product graphs, and show that training on diverse datasets improves performance over single-source pretraining. Our results demonstrate that pretraining on diverse real and synthetic graphs enhances adaptability and stability, leading to competitive performance with state-of-the-art models across various node classification tasks. This approach reduces the burden of dataset-specific training and provides a single generalist model capable of performing across multiple diverse graph structures and tasks.

# 031 032

033 034

006

008 009 010

011 012 013

014

015

016

017

018

019

021

025

026

027

028

029

# 1 INTRODUCTION

Graphs are a fundamental data structure used across diverse fields such as biology, social networks, 035 and recommendation systems (Hamilton et al., 2017). However, most graph neural network (GNN) architectures are designed in a highly specialized way, optimized for specific types of graphs. For 037 example, architectures that work well on homophilic graphs, such as citation networks, often fail to generalize to heterophilic graphs, like certain social or biological networks, due to the differences in their topologies (Abu-El-Haija et al., 2019; Yan et al., 2022). This specialization leads to a 040 fragmentation in model development, where the optimal architecture for one type of graph must be 041 significantly altered or redesigned for another. As the use of GNNs grows across diverse applications, 042 this piecemeal approach limits scalability and generalization, highlighting the need for a generalist 043 model that can handle a wide variety of graph structures without manual tuning.

044 A core challenge in building a generalist graph model lies in integrating diverse graphs, each with unique topologies, node features, and sizes, while enabling knowledge transfer across them. Without a 046 shared "vocabulary" for graph structures, models struggle to generalize effectively, as the differences 047 between graph types hinder the transfer of learned patterns (Galkin et al., 2023). At the same time, 048 recent advances in large-scale language models have shown that scaling up both model size and data diversity is essential for unlocking emergent capabilities and improving generalization across tasks (Wei et al., 2022). This makes scaling an equally critical factor in graph models. Pretraining 051 on diverse graphs requires algorithms that can efficiently handle large, heterogeneous inputs, while ensuring the model can still capture robust, transferable patterns. Therefore, building a generalist 052 graph model necessitates solutions that not only integrate diverse graph structures but also scale effectively, allowing the model to learn from vast, varied datasets without sacrificing performance.

In this work, we introduce GRAPHFM, a multi-graph pretraining framework aimed at addressing
this gap. Instead of building specialized models for each graph type, GRAPHFM uses a Perceiverbased transformer encoder (Jaegle et al., 2021b) to create a shared latent space that abstracts away
graph-specific details while preserving core structural properties. This enables the model to process a
wide range of graph types within a unified framework, moving beyond the specialist architectures
that dominate current GNN design. Our approach seeks to answer a key question: can pretraining on
diverse, multi-graph datasets lead to effective generalization and transfer across unseen graphs?

061 When tested on a variety of homophilic and heterophilic datasets, we demonstrate that our model 062 achieves performance comparable to all of the best baseline models, each of which is individually 063 tuned for its respective dataset. Overall, we achieve the best rank when compared with these models, 064 demonstrating that our approach has strong generalist performance. By combining datasets from biology, social networks, and recommendation systems, we show that our model can generalize 065 across graphs with varying topologies and features, providing the flexibility that specialized models 066 often lack. Moreover, our framework efficiently handles large mixtures of diverse graph datasets, 067 leveraging distributed training techniques to manage graphs of different sizes and complexities. 068

069 Our results show that increasing both the scale of the model and the diversity of the training data leads to significant improvements in downstream performance on new, unseen graphs and node-level tasks. 071 This demonstrates that it is indeed possible to train a generalist model on diverse graphs, which can effectively learn from and adapt to a wide range of graph types. In total, we pretrain on 152 distinct 072 graph datasets, comprising over 7.4 million nodes and 189 million edges across a wide variety of 073 graph types—an unprecedented number of different graph datasets in the literature. This extensive 074 pretraining allows our model to capture and transfer knowledge across a broad spectrum of graph 075 structures, showcasing the feasibility and advantages of building a unified model that generalizes 076 well to unseen tasks. 077

- 078 The main contributions of this work are as follows:
  - Scalable Pre-training Approach: We introduce a scalable framework for pretraining on diverse graphs using a Perceiver-based encoder with latent tokens, which efficiently handles graphs with varying sizes and topologies. Our approach includes advanced multi-graph sampling techniques that optimize GPU utilization, enabling large-scale pretraining across a wide range of graph datasets.
  - **Demonstration of Benefits from Across-Graph Pretraining:** We show that pretraining on diverse graphs significantly improves the model's ability to generalize and transfer knowledge to unseen graphs. This demonstrates that a generalist model can leverage common structural features across different datasets to outperform specialized models.
  - Scaling Analysis and Impact of Multi-Graph Pretraining: We provide the first scaling analysis for multi-graph pretraining on different domains, showing that larger models pretrained on more diverse graph datasets result in better generalization. Our results highlight that increasing both the scale of the model and the diversity of the training data improves performance on downstream tasks.

# 2 Methods

In this section, we describe our method, including the model architecture and tokenization (Section 2.1.1), our proposed multi-task node decoder for jointly solving node classification and regression tasks by querying from the latent space (Section 2.1.2), and efficient tools for scaling (Section 2.2) that allowed us to build a large pretrained model that could integrate the extreme diversity in our pretraining set.

101 102 103

079

081

082

084

085

090

092

093 094

095 096

098

100

2.1 Model

104 2.1.1 TOKENIZING DIVERSE GRAPHS

Each graph is represented as a sequence of node-level tokens, where each token embedding encodes both the node features and a positional embedding of the node. Let  $\mathcal{D} = \{\mathcal{G}_g\}_{g=1}^G$  denote a dataset containing G graphs, where each graph can be expressed in terms of its node and edges as



Figure 1: Overview of GRAPHFM architecture and multi-graph training approach: The input node-level tokens are passed through a cross-attention layer, followed by multiple self-attention layers to generate a compressed graph-level representation (latents). We decode node-level properties by creating a spatial sequence with features from a query node, a subset of its neighbors and the latents, which is then processed by a node decoder that uses self attention across the sequence.

125  $\mathcal{G}_g = (V_g, E_g)$ , with node features  $\{\mathbf{u}_i\}_{i=1}^{N_g}$ . To process a graph with a transformer, we start by 126 building a sequence of tokens as  $\mathbf{X}_g = [\mathbf{x}_1, \dots, \mathbf{x}_{N_g}]$ , where  $\mathbf{x}_i$  concatenates a projection of the 127 node features using a Multi Layer Perceptron (MLP),  $\mathbf{\tilde{u}}_i = \text{MLP}_g(\mathbf{u}_i)$ , and the positional encoding 128 (PE),  $\mathbf{p}_i$ , of the *i*<sup>th</sup> node. We use SignNet (Lim et al., 2022) which computes sign-invariant features 129 from the eigenvectors of the graph Laplacian and uses this as a basis for alignment of PE tokens 130 across all the graphs.

To build a model that can be trained across diverse graphs, we propose to tokenize each graph into a fixed and common latent space using a Perceiver encoder (Jaegle et al., 2021a). This encoder learns a set of latent query tokens which, using a cross-attention operation, query the nodes in the input graph and produce a compressed representation of it in the latent space. In the context of graphs, we can think about this as a way of routing communication between distant nodes by first going through a small number of learnable "virtual nodes" (Figure 1) that are compressed from the input graph.

For all graphs, we maintain a shared sequence of K learned latent tokens  $\mathbf{Z}_0 = [\mathbf{z}_{0,1}, \dots, \mathbf{z}_{0,K}]$ , with  $\mathbf{z}_{0,i} \in \mathbb{R}^D$  and K considerably smaller than the size of most graphs, in this work K = 512. Node embeddings in the input graph are then compressed via a cross-attention operation:

$$\mathbf{Z}_{g}^{(1)} \leftarrow \operatorname{Cross-Attn}(\mathbf{Q}_{g}, \mathbf{K}_{g}, \mathbf{V}_{g}) = \mathbf{Z}^{(0)} + \operatorname{softmax}\left(\frac{\mathbf{Q}\mathbf{K}_{g}^{T}}{\sqrt{d_{k}}}\right) \mathbf{V}_{g}, \tag{1}$$

where the queries,  $\mathbf{Q} = \mathbf{W}_q \mathbf{Z}_0$ , are projections of the learnable virtual node tokens, while the keys and values are projections of the graph's token embeddings:  $\mathbf{K}_g = \mathbf{W}_k \mathbf{X}_g$  and  $\mathbf{V}_g = \mathbf{W}_v \mathbf{X}_g$ , where the key and value weight matrices are shared by all the graphs. This operation is followed by a series of *L* self-attention blocks in the latent space to obtain a sequence of *K* latent tokens,  $\mathbf{Z}_g^{\text{out}}$ . We use the standard transformer block with pre-normalization layers and feed-forward nets (Vaswani, 2017). Note that the complexity here is  $KN_g + LK^2 \ll N_g^2$ ; When the number of latent tokens  $K \ll N_g$ , this results in a significant reduction in compute and memory.

*Remark.* Compressing every graph into a fixed set of virtual node embeddings, allows us to build a learnable "shared vocabulary" across graphs, and leverage common semantic and topological patterns across datasets and domains. Additionally, this approach also allows us to better integrate graphs of variable sizes, since most of the computation happens in the self-attention blocks, where all graphs are represented by an equally sized sequence of latent tokens.

156 2.1.2 NODE DECODER

124

141

142 143

157

Our encoder model is designed to do the bulk of the computation when processing the graph. To be able to readout node-level features, we developed a multi-task node decoder that combines the virtual node embeddings learned by our encoder  $\mathbf{Z}_g^{\text{out}}$  with local information from a node and its neighbors to create a sequence  $\mathbf{S}_g^i$  that can be processed by a transformer to produce a final node-level estimate of it's class information.



Figure 2: Characteristics of graph datasets used to train GraphFM: From left to right, we compute the histograms of the homophily ratio, average degree, number of nodes and number of edges of all 152 graphs used during training. The homophily ratio provides a measure of how frequently a node is directly connected to other nodes from the same class.

The sequence  $\mathbf{S}_{q}^{i}$  for the *i*th node can be represented as:

$$\mathbf{S}_{g}^{i} = \left[\underbrace{(\mathbf{x}_{i}; \tau_{\text{self}}), (\mathbf{x}_{\mathcal{N}_{i}^{1}}; \tau_{\text{neighbor}}) \dots (\mathbf{x}_{\mathcal{N}_{i}^{T}}; \tau_{\text{neighbor}}), (\mathbf{Z}_{1}^{\text{out}}; \tau_{\text{latent}}) \dots (\mathbf{Z}_{K}^{\text{out}}; \tau_{\text{latent}})}_{\text{virtual latent nodes}}\right], \quad (2)$$

where **x** and  $\tau_{type}$  denote the features and their token type (latent, self, or neighbor), respectively, and  $\mathcal{N}_i^j$  denotes the  $j^{\text{th}}$  neighbor selected in the neighborhood of node *i*. We use a small encoder-only transformer with a depth of *M* to obtain a final set of embeddings  $\mathbf{S}_g^{\text{out}_i}$  for node *i*. Note that the complexity is  $N_g M (K + T + 1)^2 \ll N_a^2$ .

# 2.1.3 MULTI-TASK PRETRAINING ON A VARIETY OF NODE CLASSIFICATION AND REGRESSION TASKS

In the end, a per-dataset linear classifier (or regressor)  $\mathbf{W}_g$  is tasked with producing the final predictions  $\hat{y}_i$  for node *i*, mapping the final embedding of node *i*, the first token in the  $\mathbf{S}_g^i$  sequence, to the output space as:  $\hat{y}_i = \mathbf{W}_g^T \mathbf{S}_g^{i,\text{out}}$ . The linear projection effectively translates the node-level embeddings into task-specific outputs, such as class labels for classification or continuous values for regression. The model handles a wide variety of tasks across different datasets, such as citation graphs are trained to predict academic fields and co-purchasing graphs are used to predict product categories. Each dataset has an arbitrary label space, varying not only in the number of labels but also in the nature and semantics of the output classes.

*Remark.* Since this model is trained end-to-end, the model learns how to optimally route and query
 information on graphs to maximize the performance on the various pre-training tasks. The virtual
 nodes allow for longer-range and global interactions to be encoded in the virtual node embeddings,
 and uses this information along with the local information provided by the node's neighbors.

201 202

175 176

177 178 179

181

182

183

185 186

187

188

2.2 IMPORTANT INGREDIENTS FOR TRAINING ON DIVERSE GRAPHS

203 204 2.2.1 MULTI-GRAPH PACKING

205 Typically when creating batches for training graph transformers, padding is used to extend the smaller 206 graphs to have the same size as the largest graph in the batch (Rampášek et al., 2022; Ying et al., 207 2021). This approach is likely inherited from the transformer architectures found in other domains where the context window (or sequence length) is usually fixed. But for graphs, the problem with 208 padding is particularly pronounced when there is a significant size disparity among different graphs 209 in the same batch. Alternative solutions exist, and in particular, the graph community have been 210 pioneers in batching variable-sized graphs. Message-passing frameworks combine multiple graphs 211 into a single large graph over which message passing is conducted (Fey & Lenssen, 2019b; Krell 212 et al., 2022). However, these out-of-the-box implementations are not suited for transformers which 213 use fully-connected attention. 214

215 We implement a custom data collator, which merges all graphs in the batch into a single large sequence of tokens, and adapts the attention mask to restrict each graph to itself. In particular, we

leverage Flash Attention (Dao, 2023) which makes computing attention over very large sequences
 extremely efficient. By doing so, we avoid any superfluous padding, and this in turn improves the
 computational efficiency during training.

# 2.2.2 BALANCED GPU UTILIZATION WITH THE DISTRIBUTEDSSSAMPLER

221 During multi-GPU distributed training, a global batch is formed by randomly sampling graphs from 222 different datasets, which is then equally split among the GPUs. Naively splitting the batch can lead 223 to unbalanced GPU utilization. On one hand, we can have a large batch of relatively small graphs, 224 and another where we can only have a batch with one or two very large graphs. This means that we would be forced to lower the batch size, to avoid going out of memory when multiple large graphs 225 are batched together. Our Distributed Snake Strategy Sampler (DistributedSSSampler) employs a 226 bidirectional filling strategy, where graphs, sorted by their size, are distributed in a snake-like pattern, 227 initially assigned to GPUs from right to left, then left to right and so on. This method effectively pairs 228 large graphs with small ones in subsequent passes, preventing the concentration of multiple large 229 graphs on the same GPU, thus achieving efficient load balancing and uniform GPU utilization. A 230 detailed algorithm and more details are provided in Appendix C.1. 231

We show the effectiveness of this approach in Figure 3, 232 where we demonstrate significantly lower variance in GPU 233 load compared to the default PyTorch batch sampler and 234 near 100% utilization. The effectiveness is more pro-235 nounced the more GPUs are used<sup>1</sup>. This subsequently 236 allows us to use substantially larger batch sizes, resulting 237 in further improvement in stability and a significant 2 to 238 4x speed-up in training time. 239

# 2.2.3 OVERALL TIME AND MEMORY SAVINGS

In total, our largest model, trained on all the pretraining data, takes ~6 days to train on 8 A40 GPUs for 300 epochs. With our distributed sampler, each epoch takes approximately 56 minutes (0.93 hours), compared to 299.04 minutes (~5 hours) without it. By using the distributed sampler, we observe a speedup of approximately 5.53x, reducing the total training time from 33 days to 6 days.



Figure 3: The computational benefits of using our multi-graph sampling approach: GPU memory utilization during distributed training when using the default batch sampler with 8 GPUs (left) vs. our DistributedSS-Sampler for N=4 (middle) and N=64 (right) GPUs. The total batch size is  $N \times b$ .

Please refer to Appendix C for an ablation study on the proposed sampler and multi-graph packing methods.

# 3 DATASETS

240

241

251

252 253

254

255

In standard practice, one would train on individual datasets, one at a time. However, to build our large multigraph model, we needed to curate a large dataset of graphs that have varied structures, features, and tasks.

256 Datasets used for pretraining. For pre-training, we curated a large set of 80 real-world graph 257 datasets from the PyTorch Geometric library (Fey & Lenssen, 2019a) and Network Repository (Rossi 258 & Ahmed) (Figure 2). These datasets span a wide range of domains, including: citation networks, 259 product recommendation graphs, webpage traffic graphs, biological protein-protein interactions, and 260 molecular graphs, and vary in their degree of heterophily (extent to which neighbors share the same 261 class or node-level labels). Each dataset contributes unique structural patterns and tasks, providing 262 a rich source for our model to learn diverse graph representations. In addition to these realworld 263 datasets, we generated 72 synthetic graphs (Tsitsulin et al., 2022) that vary in their hetero- and 264 homophily ratios and overall size and density (see Appendix B.1). We note that most datasets used in 265 popular benchmarks were left out of pretraining in order to test the pretrained model on these datasets in out-of-distribution (OOD) finetuning. 266

In Figure 2, we show a summary of various graph statistics, including the number of nodes and edges, the average degree of each node, and the homophily ratio of the graph. The homophily ratio ranges

<sup>&</sup>lt;sup>1</sup>The same effect can be obtained using gradient accumulation when resource bound. See Appendix C.1

from 0 to 1 and encodes the average amount of nodes with nearest neighbors from the same class. When comparing our realworld datasets with the synthetic graphs added to the mix (Figure 2), we see a good amount of overlap between most features except for the average degree. The average degree of realworld graphs spans a larger range, and the synthetic graphs have a more limited range. We also find an enrichment of heterophilic graphs with low homophily ratio in the added synthetic data. In total, we counted more than 7.4M nodes and 163.9M edges across all 152 datasets used for pretraining We point the reader to Appendix B.1 for a detailed description of all datasets.

277 Datasets for testing out-of-distribution transfer To demonstrate the adaptability of our pretrained 278 model through fine-tuning on unseen data (out-of-distribution datasets), we leverage a smaller, but 279 equally diverse set of graph datasets that are commonly used as benchmarks (see Appendix B.4). We 280 use 10 different datasets that range from academic collaboration networks like "Coauthor-CS" and 281 "Coauthor-Physics" (Sinha et al., 2015) to webpage link datasets such as "Chameleon" and "Squirrel" 282 (Rozemberczki et al., 2021), which are particularly challenging due to their low homophily ratios, 283 indicating less connectivity within the same class. These datasets not only test the transferability of 284 the learned representations but also highlight the model's capability in handling graphs with varied node degrees and class distributions. 285

286 287

288

- 4 Results
- 4.1 EXPERIMENTAL SETUP

**Training:** To train all of our models, we employed the LAMB optimizer (You et al., 2019) with a learning rate of  $10^{-4}$ . The learning rate is scheduled based on a linear warmup of 2 epochs, followed by cosine decay until the end of training. We use bfloat16 mixed-precision and flash attention (Dao, 2023) for higher compute efficiency while training. We trained our largest model (75M parameters) for 6.4 days on 8 NVIDIA A40 GPUs. We point the reader to further details on the architecture and model training in Appendix A.1.

297 **Baselines:** We compared GRAPHFM against six baseline models that were consistently reported in 298 both heterophilic and homophilic benchmarks. This included two GNN-based models: GCN(Kipf 299 & Welling, 2016) and GAT(Velickovic et al., 2017), two transformer-based models: SAN (Kreuzer et al., 2021) and NAGphormer (Chen et al., 2022b), and two heterophily-based models: MLP and 300 H2GCN (Zhu et al., 2020). For all of the baseline models, we include the best reported accuracy, 301 and when there are no reported results for a dataset, we extensively tuned each model as in standard 302 practice (see Appendix B.4). We also provide additional baselines in Appendix D.4 reported for 303 subsets of the datasets tested. 304

305 **Evaluation:** To evaluate the generalization of our pretrained model on new datasets that it hasn't 306 encountered during pretraining, we employed two fine-tuning strategies: (i) Low-resource MLP fine-tuning (MFT), where we freeze the encoder and node decoder weights and only update the 307 feature MLP weights, and (ii) combined MLP and node decoder fine-tuning (NFT), where we also 308 adapt the node decoder weights. MFT is aimed at evaluating near out-of-the-box performance by 309 leveraging the model's pretrained knowledge, with minimal additional training, whereas NFT allows 310 for more flexibility by adjusting weights of the pretrained node decoder to better align with the OOD 311 data. For all the fine-tuning experiments, we used a learning rate of  $10^{-3}$  and a weight decay of  $10^{-5}$ , 312 optimized using the AdamW optimizer (Loshchilov & Hutter, 2017), and use a gradual unfreezing 313 strategy to update the node decoder weights in our NFT experiments. Further details are provided in 314 the Appendix A.3.

- 315 316
- 4.2 EXPERIMENTS
- 317 318 319
- C

Q1: IS IT POSSIBLE TO BUILD A LARGE MODEL SPANNING MANY DOMAINS?

Recent efforts in graph neural networks (GNNs) have shown success in training models on many
 graphs (Beaini et al., 2023; Mao et al., 2024). However, these approaches primarily focus on graphs
 with homogeneous structures, limiting their ability to generalize across different types of graphs. In
 this experiment, we aim to address a more ambitious question: can we effectively train a large model
 on diverse, multi-graph datasets that vary significantly in their topologies, features, and downstream

classification tasks? Our goal is to determine whether a generalist model can span multiple graph
 domains and improve out-of-distribution (OOD) performance through diverse pretraining.

We trained three different model sizes: a small model with 327 389K parameters, a medium model with 18M parameters, 328 and a large model with 75M parameters. Each model 329 was pretrained on progressively larger datasets containing 330 different amounts of graph data, ranging from 200K to-331 kens (small), to 2M tokens (medium), and finally to 7.3M 332 tokens (large), created by taking random subsets of the 333 largest dataset (refer to Appendix B.2 for more details). 334 The datasets span a variety of real-world graph types and structures, as described in Section 3. For the largest scale 335 of data, we also introduced synthetic graphs into the mix to 336 further test the model's ability to generalize across highly 337 diverse graph structures. The synthetic graphs provided 338 additional variability in both topologies and node features, 339 allowing us to assess how well the model can handle graph 340 data that extends beyond typical real-world scenarios. 341

To evaluate how well the pretrained models generalize to 342 new, unseen data, we applied our lightweight MLP fine-343 tuning approach (MFT) on a set of nine held-out datasets. 344 These include four homophilic datasets (Coauthor-cs, 345 Coauthor-physics, Amazon-photos, and Amazon-comp) 346 and five heterophilic datasets (Texas, Wisconsin, Actor, 347 Squirrel, and Chameleon). As illustrated in Figure 4A, we 348 observe that performance on these OOD datasets improves 349 consistently as the data size increases. Notably, the largest 350 model, trained on the full 7.3M tokens, achieves a 2.1% 351 improvement in accuracy compared to the smaller models.

We further stratified our pretraining dataset to investigate
the effects of cross-domain training by creating three models: (i) "Soc" with social domain graphs (1.3M tokens),
(ii) "Soc + Bio" with social and biological graphs (2M
tokens), and (iii) "All" with all data, including synthetic



Figure 4: Scaling Analysis: (A) Average accuracy across OOD datasets (MFT) for model sizes (389K, 18M, 75M) and token counts (200K, 2M, 7.3M) seen during pre-training, using random splits of the pre-training data. (B) Accuracy on Coauthor-CS (citation domain) and Amazon-Photo (co-purchasing network) for the 75M model across different domain-wise pre-training splits.

graphs (7.3M tokens). As shown in Figure 4B, adding biological datasets improved performance on
 both Coauthor-CS (citation domain) and Amazon-Photo (co-purchasing network). This suggests that
 performance continues to scale even if the additional data is from seemingly unrelated domains (refer
 to Appendix D.1 for additional results).

These results underscore the importance of both model scale and data diversity. With more data diversity and larger models, the pretrained model demonstrates stronger generalization capabilities.
 This scaling analysis provides strong evidence that cross-domain pretraining enables better OOD performance, further validating the benefits of training on diverse datasets. Detailed configurations for each model size are provided in Appendix A.1.

366 367

## Q2: How does our generalist approach compare with others?

368

Next, we compared the performance of our largest model (75M) trained on all of the data, alongside a number of message passing architectures and state-of-the-art transformer-based models. To adapt our approach to new datasets, we freeze our pretrained encoder and then finetune either the feature MLPs (MFT) or the feature MLPs and node decoder parameters (NFT). In both of these cases, we use the same hyperparameters (learning rate =  $10^{-3}$ ) to finetune the model. In the case of NFT, we also incorporate a simple unfreezing schedule for updates which adds an additional two hyperparameters that we need to tune (refer to Appendix A.3.2).

On both homophilic and heterophilic benchmarks (Table 1), GRAPHFM performs on par with state of-the-art specialist methods that are trained from scratch on each dataset. While the best-performing model among the baseline methods varies across datasets, GRAPHFM consistently ranks among

Table 1: Results on a variety of homophilic and heterophilic node classification benchmarks. From left to
right, we show different message passing and graph transformer architectures, and then GRAPHFM in both the
lightweight MLP-only finetuning (MFT) and node decoder finetuning (NFT). The top three numbers are bold,
with the highest in bright red fading to black. Models are ranked on all 10 datasets and the average and standard
deviation ranking is at the bottom.

	GCN	MLP	GAT	H2GCN	SAN	NAG	GraphFM-MFT	GraphFM-NFT	
.9 Physics	95.38±0.20	95.12±0.26	$95.14{\pm}0.28$	96.28±0.13	96.83±0.18	96.66±0.16	96.64±0.17	96.77±0.12	
E CS	94.06±0.16	$92.99 \pm 0.51$	93.61±0.14	$94.02 \pm 0.31$	94.16±0.36	95.00±0.14	95.19±0.21	95.24±0.18	
o Photo	85.94±1.18	$88.66 {\pm} 0.85$	$87.13 {\pm} 1.00$	$91.56 {\pm} 0.80$	94.17±0.65	94.64±0.60	93.01±1.82	94.37±0.35	
5 Computer	89.47±0.46	84.63	90.78±0.13	$89.33{\pm}0.27^{*}$	89.83±0.16	$91.22 {\pm} 0.14$	$89.95 {\pm} 0.83$	90.07±0.21	
<sup>≖</sup> Ogbn arxiv	$\textbf{70.40{\pm}0.10}$	$52.63{\pm}0.12$	$67.56{\pm}0.12$	$68.29{\pm}0.67$	69.17±0.15	$68.21{\pm}0.02^{*}$	$69.96{\pm}0.21$	$\textbf{70.01} \pm \textbf{0.18}$	
.일 Texas	55.14±5.16	80.81±3.31	$52.16{\pm}6.63$	84.86±7.23	60.17±6.66	$68.37{\pm}5.27^{*}$	80.81±2.76	82.16±3.24	
. Wisconsin	51.76±3.06	85.29±3.31	$49.41 {\pm} 4.09$	87.65±4.98	51.37±2.08	$68.23{\pm}5.99^{*}$	$83.13 {\pm} 2.35$	$83.62 \pm 3.21$	
2 Actor	27.32±1.10	36.63±0.70	$27.44{\pm}0.89$	$35.70 \pm 1.00$	27.32±1.10	$34.33{\pm}0.94^{*}$	36.29±0.63	38.01±1.07	
E Chameleon	38.44±1.92	$46.21 \pm 2.99$	$38.44{\pm}1.92$	60.11±2.15	$44.32{\pm}1.73^*$	$57.39{\pm}0.02^*$	58.64±1.24	59.12±1.64	
<sup>≖</sup> Squirrel	$31.52{\pm}0.71$	$28.77 {\pm} 1.56$	$36.77 {\pm} 1.68$	$36.48{\pm}1.86$	$30.92{\pm}2.14^*$	$49.93{\pm}0.07^{*}$	$42.80{\pm}1.54$	$42.98{\pm}1.62$	
Avg Rank	5.8 ± 1.9	$6.3\pm2.5$	$6.2\pm1.8$	$4.0\pm2.5$	4.7 ± 2.2	$3.3\pm1.8$	$3.3\pm0.6$	$2.1\pm0.7$	
* This result wa	* This result was missing from existing literature and was obtained through extensive hyperparameter tuning.								



Figure 5: Analysis of the learning dynamics. Learning dynamics for 100 (A) random GCN and (B) NAG
(NAGphormer) models compared against our lightweight finetuned model GraphFM (MFT) for four datasets.
GRAPHFM works out of the box and achieves rapid learning on new datasets with few training steps, while the other approaches are less stable and often require early stopping with decreased performance over training.

the top three: GRAPHFM (NFT) achieves the highest average rank overall and GRAPHFM (MFT)
is tied for the second position with NAG. Note that GRAPHFM (MFT) demonstrates significantly
lower variance in rank, indicating more stable performance compared to NAG, which exhibits a more
bimodal distribution in its ranking.

In contrast, baseline models may excel on a few datasets but perform poorly on others. For example, H2GCN is a top performer on heterophilic datasets but struggles with homophilic ones, whereas NAG shows the opposite behavior. These baseline models are highly specialized and designed for specific types of datasets, limiting their generalization across diverse graph types. Additionally, it's important to note that all baseline models underwent extensive hyperparameter tuning, whereas GRAPHFM performs consistently well without any further tuning. Furthermore, NFT provides significant benefit for datasets like photos and actor. By making part of the pre-trained model learnable, we are able to better adapt to the OOD datasets. 

Q3: How does our model generalize out-of-the-box?

A major challenge in applying graph-based models is the extensive tuning often required to achieve
competitive performance. Most models are highly sensitive to hyperparameters like learning rate,
depth, and weight decay. Tuning these hyperparameters, especially across datasets with different
graph topologies and sizes, requires significant time and computational resources, and even then,
finding a good configuration can be difficult (Guo et al., 2022; Tsitsulin et al., 2022).

In contrast, GRAPHFM offers strong out-of-the-box performance without requiring any significant tuning. To demonstrate this, we evaluated GRAPHFM using the same fixed learning rate and weight decay across multiple datasets (learning rate =  $10^{-3}$ , weight decay =  $10^{-5}$ ) and observed stable and high performance across all datasets (Figure 5). Fine-tuning GRAPHFM with our simple MFT strategy resulted in low variance and rapid convergence, without the need for extensive hyperparameter exploration. This makes GRAPHFM highly efficient and cost-effective compared to models that require substantial tuning.

To highlight this contrast, we compared the performance of GRAPHFM with 100 randomly configured
versions of GCN and the best performing transformer-based NAGformer (Chen et al., 2022b). Both
baseline models exhibit a wide range of performance depending on the hyperparameter choices,
with some configurations leading to significant instability or poor results. For instance, in the Texas
dataset, GCN required exhaustive exploration of hyperparameter settings to find a stable and effective
configuration. Similarly, NAGformer's performance fluctuated greatly depending on the dataset and
the selected parameters, further emphasizing the cost of tuning.

446 Additionally, GRAPHFM demonstrates quick convergence, 447 reaching near-optimal performance within 10-20 training steps, 448 in stark contrast to GCN, which required considerably more iter-449 ations to converge. This efficiency is a direct result of leveraging a pretrained model, which allows GRAPHFM to start from a 450 robust initialization and quickly adapt to the target task. The 451 reduced need for hyperparameter tuning and faster convergence 452 further solidify the advantages of pretraining in minimizing com-453 putational overhead and time-to-solution. Ultimately, our results 454 position GRAPHFM as a cost-effective and reliable choice for 455 a wide range of node classification tasks.

456

458

# 457 Q4: How stable are the solutions?

459 Graph-based models are highly sensitive to hyperparameter configurations, where even small deviations from optimal settings 460 can lead to substantial performance degradation. This sensitivity 461 poses significant challenges for ensuring stable and robust de-462 ployment. Thus, we wanted to examine the stability of model 463 performance by exploring the performance landscape around the 464 optimal hyperparameter configuration. We analyze the perfor-465 mance of both a GCN and GRAPHFM (MFT) on Coauthor-CS 466 (homophilic) and Chameleon (heterophilic) datasets for different 467 hyperparameters around the optimal solution (Figure 6). The 468 set of hyperparameters is marked with a star, and other mod-469 els are colored based on the normalized  $\ell_2$ -distance of their hyperparameter vectors to the optimal hyperparameter vector. 470 For GRAPHFM, we observe that the distribution is concentrated 471 around the optimal point, suggesting low sensitivity to the choice 472



Figure 6: Comparison of model sensitivity. The performance of GCN and GRAPHFM for 100 different random hyperparameters on Chameleon and Coauthor-CS. Star denotes the model with the optimal hyperparameters, and the color indicates the  $\ell_2$ -distance between the optimal solution and each model's hyperparameters.

of the hyperparameters used for finetuning. We also observe that the relationship between hyperparameter deviation and performance is linear. On the other hand, for the GCN model, small deviations in hyperparameters can lead to large changes in performance, suggesting instability of the model with respect to the hyperparameters and a much noisier landscape around the optimal model.

477 478

479

# 5 RELATED WORK

Graph foundation modeling approaches. Foundation models have achieved significant success for language, vision and timeseries data (Radford et al., 2018; Dehghani et al., 2023; Das et al., 2023). These models are pre-trained on large datasets and can be adapted to a wide range of downstream tasks, effectively utilizing both prior knowledge from the pre-training stage and data from the downstream tasks to enhance performance (Brown et al., 2020). The concept of foundation models has recently extended into graph learning, leading to the development of Graph Foundation Models (GFMs) (Ibarz et al., 2022; Beaini et al., 2023; Galkin et al., 2023; Mao et al., 2024). These

models aim to generalize across diverse graph-structured data by leveraging large-scale pretraining,
similar to foundation models in vision and language domains. Recent efforts have primarily focused
on domain-specific GFMs, such as Mole-BERT for molecular graphs, which utilizes pretraining to
improve property prediction for molecules and materials (Xia et al., 2023). Additionally, large-scale
models like MatterSim (Yang et al., 2024), designed to predict molecular behaviors across different
elements and conditions.

492 Beyond domain-specific applications, GFMs are increasingly being developed for more general 493 tasks. Similarly, recent advancements have explored scaling laws in graph models, showing that 494 larger models can lead to improved transfer learning and generalization (Liu et al., 2024). Similar to 495 theirs, our work shows that scale improves performance. However, unique from all of these works is 496 our result for cross-domain pretraining to enhance generalization across diverse graph topologies. Triplet-GMPNN (Ibarz et al., 2022) which is a foundational GNN for algorithmic reasoning tasks that 497 is trained to perform various tasks from the CLRS benchmark (Veličković et al., 2022), or ULTRA 498 (Galkin et al., 2023), a foundation model for knowledge graphs trained on graphs with arbitrary entity 499 and relation vocabularies. Recent work has also shown how to use language modeling to help unify 500 many graphs (Liu et al., 2023b). 501

- 502 **Scaling up graph transformers.** Graph transformers bypass standard local learning rules in GNNs 503 by allowing all nodes on the graph to interact through self-attention (Dwivedi & Bresson, 2020). 504 However, due to the high computational cost and benefits of the inductive bias in message passing, a 505 number of methods have been proposed to move beyond full self-attention or combine transformers 506 with GNNs. One class of methods combine transformer blocks with GNNs, including GraphTrans 507 (Wu et al., 2021), GraphGPS (Rampášek et al., 2022), and SAT (Chen et al., 2022a). Another 508 strategy is to reduce the computational complexity by using the transformer module on a coarsened 509 or compressed graph. For instance, ANS-GT (Zhang et al., 2022) introduced a node-sampling-based graph transformers, incorporating hierarchical attention and graph coarsening, and Gapformer (Liu 510 et al., 2023a) uses k-hops local pooling and global pooling to coarsen the large graph into a smaller 511 set of nodes. Exphormer (Shirzad et al., 2023) coarsens the graph by doing computation through 512 expander graphs (Deac et al., 2022). This idea of compression has also been studied through the lens 513 of "skeletonization" (Cao et al., 2024) where they learn to identify uninformative background nodes 514 (Cao et al., 2024) and use this information to compress them to achieve competitive performance 515 with as little as 1% of the nodes in the graph. Many of these approaches leverage virtual nodes to 516 facilitate message passing across large graphs, however, the compression techniques used in these 517 works are often based on heuristics like pooling layers or expander graphs, in contrast to our work 518 where the compression is fully learned.
- 519 520

# 6 CONCLUSION

521 522

In this paper, we introduced GRAPHFM, a novel approach for multi-graph pretraining that effectively handles diverse graph datasets across various domains. A key finding of our work is the positive effect of scaling both model size and data diversity. We show that cross-domain pretraining leads to better out-of-distribution performance, proving that the inclusion of diverse graph types significantly enhances the model's ability to adapt to new, unseen data. This reveals the potential for graph foundation models to benefit from combining datasets across domains, facilitating more efficient and robust training processes.

While our results are promising, there are several areas for future exploration. Our current work
primarily focuses on node-level classification tasks; extending GRAPHFM to support tasks like
graph-level classification, link prediction, and self-supervised learning could broaden its applicability.
Moreover, expanding the diversity of pretraining datasets, such as including point clouds, mesh
graphs, or knowledge graphs, may further enhance the model's generalization capabilities and impact
across domains.

Looking ahead, we believe that generalist graph models like GRAPHFM have the potential to
 transform a variety of fields, particularly in scenarios where data is scarce or incomplete. Our work
 represents an important step toward more universal and adaptable graph models, and we anticipate
 further research into cross-domain pretraining as a promising direction for the future of graph
 learning.

# 540 REFERENCES

563

565

566

569

570

571

577

Sami Abu-El-Haija, Bryan Perozzi, Amol Kapoor, Nazanin Alipourfard, Kristina Lerman, Hrayr
Harutyunyan, Greg Ver Steeg, and Aram Galstyan. Mixhop: Higher-order graph convolutional
architectures via sparsified neighborhood mixing. In *international conference on machine learning*,
pp. 21–29. PMLR, 2019. 1

- Mehdi Azabou, Venkataramana Ganesh, Shantanu Thakoor, Chi-Heng Lin, Lakshmi Sathidevi, Ran Liu, Michal Valko, Petar Veličković, and Eva L Dyer. Half-hop: A graph upsampling approach for slowing down message passing. In *International Conference on Machine Learning*, pp. 1341–1360.
  PMLR, 2023. 19
- Dominique Beaini, Shenyang Huang, Joao Alex Cunha, Gabriela Moisescu-Pareja, Oleksandr Dymov, Samuel Maddrell-Mander, Callum McLean, Frederik Wenkel, Luis Müller, Jama Hussein Mohamud, et al. Towards foundational models for molecular learning on large-scale multi-task datasets. *arXiv preprint arXiv:2310.04292*, 2023. 6, 9
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
   Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
   few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020. 9
- Linfeng Cao, Haoran Deng, Yang Yang, Chunping Wang, and Lei Chen. Graph-skeleton: 1% nodes are sufficient to represent billion-scale graph. In *Proceedings of the ACM on Web Conference 2024*, WWW '24, pp. 570–581, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400701719. doi: 10.1145/3589334.3645452. URL https://doi.org/10.1145/3589334.3645452. 10
  - Dexiong Chen, Leslie O'Bray, and Karsten Borgwardt. Structure-aware transformer for graph representation learning. In *International Conference on Machine Learning*, pp. 3469–3489. PMLR, 2022a. 10
- Jinsong Chen, Kaiyuan Gao, Gaichao Li, and Kun He. Nagphormer: A tokenized graph transformer
   for node classification in large graphs. *arXiv preprint arXiv:2206.04910*, 2022b. 6, 9
  - Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv* preprint arXiv:2307.08691, 2023. 5, 6
- Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for time-series forecasting. *arXiv preprint arXiv:2310.10688*, 2023. 9
- Andreea Deac, Marc Lackenby, and Petar Veličković. Expander graph propagation. In *NeurIPS* 2022 Workshop on Symmetry and Geometry in Neural Representations, 2022. URL https://openreview.net/forum?id=6cthqh2qhCT. 10
- Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, 578 Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, Rodolphe Jenat-579 ton, Lucas Beyer, Michael Tschannen, Anurag Arnab, Xiao Wang, Carlos Riquelme Ruiz, Matthias 580 Minderer, Joan Puigcerver, Utku Evci, Manoj Kumar, Sjoerd Van Steenkiste, Gamaleldin Fathy 581 Elsayed, Aravindh Mahendran, Fisher Yu, Avital Oliver, Fantine Huot, Jasmijn Bastings, Mark 582 Collier, Alexey A. Gritsenko, Vighnesh Birodkar, Cristina Nader Vasconcelos, Yi Tay, Thomas 583 Mensink, Alexander Kolesnikov, Filip Pavetic, Dustin Tran, Thomas Kipf, Mario Lucic, Xiaohua 584 Zhai, Daniel Keysers, Jeremiah J. Harmsen, and Neil Houlsby. Scaling vision transformers to 22 bil-585 lion parameters. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), Proceedings of the 40th International Conference on Machine 586 Learning, volume 202 of Proceedings of Machine Learning Research, pp. 7480–7512. PMLR, 23– 29 Jul 2023. URL https://proceedings.mlr.press/v202/dehghani23a.html. 588 589

593 Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019a. 5

Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs.
 *arXiv preprint arXiv:2012.09699*, 2020. 10

594 595 596	Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In <i>ICLR Workshop on Representation Learning on Graphs and Manifolds</i> , 2019b. 4
597 598	Mikhail Galkin, Xinyu Yuan, Hesham Mostafa, Jian Tang, and Zhaocheng Zhu. Towards foundation models for knowledge graph reasoning. <i>arXiv preprint arXiv:2310.04562</i> , 2023. 1, 9, 10
599 600 601	Lingbing Guo, Qiang Zhang, and Huajun Chen. Unleashing the power of transformer for graphs. <i>arXiv preprint arXiv:2202.10581</i> , 2022. 8
602 603	William L Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. <i>arXiv preprint arXiv:1709.05584</i> , 2017. 1
604 605 606	Van Thuy Hoang, O Lee, et al. Mitigating degree biases in message passing mechanism by utilizing community structures. <i>arXiv preprint arXiv:2312.16788</i> , 2023. 17
607 608	Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. <i>Social networks</i> , 5(2):109–137, 1983. 16
610 611 612	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. <i>Advances in neural information processing systems</i> , 33:22118–22133, 2020. 17
613 614 615 616 617 618 619	<ul> <li>Borja Ibarz, Vitaly Kurin, George Papamakarios, Kyriacos Nikiforou, Mehdi Bennani, Róbert Csordás, Andrew Joseph Dudzik, Matko Bošnjak, Alex Vitvitskyi, Yulia Rubanova, Andreea Deac, Beatrice Bevilacqua, Yaroslav Ganin, Charles Blundell, and Petar Veličković. A generalist neural algorithmic learner. In Bastian Rieck and Razvan Pascanu (eds.), <i>Proceedings of the First Learning on Graphs Conference</i>, volume 198 of <i>Proceedings of Machine Learning Research</i>, pp. 2:1–2:23. PMLR, 09–12 Dec 2022. URL https://proceedings.mlr.press/v198/ibarz22a.html. 9, 10</li> </ul>
620 621 622 623	Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A general architecture for structured inputs & outputs. <i>arXiv preprint arXiv:2107.14795</i> , 2021a. 3
624 625 626	Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In <i>International conference on machine learning</i> , pp. 4651–4664. PMLR, 2021b. 2
627 628 629	Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. <i>arXiv preprint arXiv:1609.02907</i> , 2016. 6
630 631 632	Mario Michael Krell, Manuel Lopez, Sreenidhi Anand, Hatem Helal, and Andrew William Fitzgibbon. Tuple packing: Efficient batching of small graphs in graph neural networks. <i>arXiv preprint</i> <i>arXiv:2209.06354</i> , 2022. 4
633 634 635 636	Devin Kreuzer, Dominique Beaini, Will Hamilton, Vincent Létourneau, and Prudencio Tossou. Rethinking graph transformers with spectral attention. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 34:21618–21629, 2021. 6
637 638 639 640	Benjamin Lefaudeux, Francisco Massa, Diana Liskovich, Wenhan Xiong, Vittorio Caggiano, Sean Naren, Min Xu, Jieru Hu, Marta Tintore, Susan Zhang, Patrick Labatut, Daniel Haziza, Luca Wehrstedt, Jeremy Reizenstein, and Grigory Sizov. xformers: A modular and hackable transformer modelling library. https://github.com/facebookresearch/xformers, 2022. 15
641 642 643 644	Derek Lim, Joshua Robinson, Lingxiao Zhao, Tess Smidt, Suvrit Sra, Haggai Maron, and Stefanie Jegelka. Sign and basis invariant networks for spectral graph representation learning. <i>arXiv</i> preprint arXiv:2202.13013, 2022. 3
645 646 647	<ul> <li>Chuang Liu, Yibing Zhan, Xueqi Ma, Liang Ding, Dapeng Tao, Jia Wu, and Wenbin Hu. Gapformer:</li> <li>Graph transformer with graph pooling for node classification. In <i>Proceedings of the 32nd International Joint Conference on Artificial Intelligence (IJCAI-23)</i>, pp. 2196–2205, 2023a. 10, 17, 19</li> </ul>

648 649 650	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. <i>arXiv preprint arXiv:2310.00149</i> , 2023b. 10
651 652 653	Jingzhe Liu, Haitao Mao, Zhikai Chen, Tong Zhao, Neil Shah, and Jiliang Tang. Neural scaling laws on graphs. <i>arXiv preprint arXiv:2402.02054</i> , 2024. 10
654 655 656	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. <i>arXiv preprint arXiv:1711.05101</i> , 2017. 6, 16
657 658	Yi Luo, Guangchun Luo, Ke Yan, and Aiguo Chen. Inferring from references with differences for semi-supervised node classification on graphs. <i>Mathematics</i> , 10(8):1262, 2022. 17
659 660	Haitao Mao, Zhikai Chen, Wenzhuo Tang, Jianan Zhao, Yao Ma, Tong Zhao, Neil Shah, Michael Galkin, and Jiliang Tang. Graph foundation models. <i>arXiv preprint arXiv:2402.02216</i> , 2024. 6, 9
662 663 664	Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. Image-based recommendations on styles and substitutes. In <i>Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval</i> , pp. 43–52, 2015. 17
665 666 667 668	John Palowitch, Anton Tsitsulin, Brandon Mayer, and Bryan Perozzi. Graphworld: Fake graphs bring real insights for gnns. In <i>Proceedings of the 28th ACM SIGKDD Conference on Knowledge</i> <i>Discovery and Data Mining</i> , pp. 3691–3701, 2022. 16
669 670	Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-gcn: Geometric graph convolutional networks. <i>arXiv preprint arXiv:2002.05287</i> , 2020. 19
671 672 673	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018. 9
674 675 676	Ladislav Rampášek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Do- minique Beaini. Recipe for a general, powerful, scalable graph transformer. <i>Advances in Neural</i> <i>Information Processing Systems</i> , 35:14501–14515, 2022. 4, 10
677 678 679	RA Rossi and NK Ahmed. Networkrepository: An interactive data repository with multi-scale visual analytics, 2014, eprint arxiv, 2014. 5
680 681	Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding. <i>Journal of Complex Networks</i> , 9(2):cnab014, 2021. 6, 19
682 683 684 685	Hamed Shirzad, Ameya Velingker, Balaji Venkatachalam, Danica J Sutherland, and Ali Kemal Sinop. Exphormer: Sparse transformers for graphs. In <i>International Conference on Machine Learning</i> , pp. 31613–31632. PMLR, 2023. 10
686 687 688	Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June Hsu, and Kuansan Wang. An overview of microsoft academic service (mas) and applications. In <i>Proceedings of the 24th</i> <i>international conference on world wide web</i> , pp. 243–246, 2015. <b>6</b> , 17
689 690 691	Anton Tsitsulin, Benedek Rozemberczki, John Palowitch, and Bryan Perozzi. Synthetic graph generation to benchmark graph learning. <i>arXiv preprint arXiv:2204.01376</i> , 2022. <b>5</b> , 8
692	A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017. 3
693 694 695	Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, Yoshua Bengio, et al. Graph attention networks. <i>stat</i> , 1050(20):10–48550, 2017. 6
696 697 698	Petar Veličković, Adrià Puigdomènech Badia, David Budden, Razvan Pascanu, Andrea Banino, Misha Dashevskiy, Raia Hadsell, and Charles Blundell. The clrs algorithmic reasoning benchmark. <i>arXiv preprint arXiv:2205.15659</i> , 2022. 10
700 701	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. <i>arXiv preprint arXiv:2206.07682</i> , 2022. 1

702 703 704	Zhanghao Wu, Paras Jain, Matthew Wright, Azalia Mirhoseini, Joseph E Gonzalez, and Ion Stoica. Representing long-range context for graph neural networks with global attention. <i>Advances in Neural Information Processing Systems</i> , 34:13266–13279, 2021. 10
705 706 707	Jun Xia, Chengshuai Zhao, Bozhen Hu, Zhangyang Gao, Cheng Tan, Yue Liu, Siyuan Li, and Stan Z Li. Mole-bert: Rethinking pre-training graph neural networks for molecules. 2023. 10
708 709 710	Yujun Yan, Milad Hashemi, Kevin Swersky, Yaoqing Yang, and Danai Koutra. Two sides of the same coin: Heterophily and oversmoothing in graph convolutional neural networks. In 2022 IEEE International Conference on Data Mining (ICDM), pp. 1287–1292. IEEE, 2022. 1
711 712 713 714	Han Yang, Chenxi Hu, Yichi Zhou, Xixian Liu, Yu Shi, Jielan Li, Guanzhi Li, Zekun Chen, Shuizhou Chen, Claudio Zeni, et al. Mattersim: A deep learning atomistic model across elements, tempera- tures and pressures. arXiv preprint arXiv:2405.04967, 2024. 10
715 716 717	Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation? <i>Advances in Neural Information Processing Systems</i> , 34:28877–28888, 2021. 4
718 719 720 721	Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. arXiv preprint arXiv:1904.00962, 2019. 6
722 723 724	<ul> <li>Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. Graph- saint: Graph sampling based inductive learning method. <i>arXiv preprint arXiv:1907.04931</i>, 2019.</li> <li>20</li> </ul>
725 726 727	Zaixi Zhang, Qi Liu, Qingyong Hu, and Chee-Kong Lee. Hierarchical graph transformer with adaptive node sampling. <i>Advances in Neural Information Processing Systems</i> , 35:21171–21183, 2022. 10
729 730 731	Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Beyond homophily in graph neural networks: Current limitations and effective designs. <i>Advances in neural information processing systems</i> , 33:7793–7804, 2020. 6
732	
733	
734	
735	
736	
737	
738	
739	
741	
742	
743	
744	
745	
746	
746 747	
746 747 748	
746 747 748 749	
746 747 748 749 750	
746 747 748 749 750 751	
746 747 748 749 750 751 752	
746 747 748 749 750 751 752 753	
746 747 748 749 750 751 752 753 754	

758

759

APPENDIX

#### ADDITIONAL MODEL DETAILS А

#### 760 MODEL CONFIGURATION DETAILS A.1 761

762 We used pretrained 3 configuration of models—small (398K), medium (18M) and large (75M)—for 763 our analysis. Details of the configuration for each model are given in Table A1. In the first cross-764 attention layer, we used flash attention, whereas for all subsequent attention layers, we used memory-765 efficient attention. Both implementations were sourced from the xFormers library (Lefaudeux et al., 766 2022).

767 768

769

Table A1: Architectural details of GraphFM for different parameter sizes used in Section 2.2

770	Parameter Count	75M	18M	389K
771				
772	Num Latents $(K)$	512	256	32
773	Latent Dimension	512	256	32
774	<b>Cross-Attention</b>			
775				
776	Heads	4	4	4
777	FFN hidden dim	2048	1024	128
778				
779	Self-Attention			
780	Depth $(L)$	12	10	4
781	Heads	8	4	4
782				
783	FFN hidden dim	2048	1024	128
784	Node Decoder			
785	Dopth $(M)$	4	4	2
786	Depth (M)	4	4	Z
787	Heads	8	4	4
788	FFN hidden dim	2048	1024	128
789				

790 791

792

# A.2 RESCALING THE LEARNING RATES FOR DIFFERENT GRAPH SIZES

793 When training on variable sized graphs, the MLP and linear decoder for each dataset receive updates based on the number of nodes from their respective datasets present in the batch and thus smaller 794 graphs get updated less frequently when compared to large graphs. To mitigate this imbalance, we 795 implemented dataset-specific learning rates for the feature MLP and linear decoders. Since they 796 receive updates less frequently, when they do, we use a larger learning rate to update them. Without 797 this adjustment, the weights of the common Perceiver encoder and node decoder would advance more 798 quickly than those of the dataset-specific components, potentially leading to suboptimal learning for 799 smaller datasets.

800 801 802

803

804

805

A.3 FINE-TUNING STRATEGIES

In our evaluation of GraphFM's generalization capability, we employed two fine-tuning strategies aimed at adapting the model to out-of-distribution (OOD) datasets.

806 A.3.1 LOW-RESOURCE MLP FINE-TUNING (MFT) 807

This approach is designed to assess how well the pretrained model performs out-of-the-box on 808 different OOD graphs without extensive training. In MFT, we freeze the pretrained model and only 809 fine-tune a lightweight multi-layer perceptron (MLP) on top of the learned representations. This strategy allows us to quickly adapt the model to new tasks while retaining the majority of the original learned parameters. MFT is particularly useful in low-resource settings, where computational power or time is limited, as it requires minimal additional training while still providing insight into how well the pretrained model generalizes. For all MFT experiments, we used a learning rate of  $10^{-3}$  and a weight decay of  $10^{-5}$ , optimized using the AdamW optimizer (Loshchilov & Hutter, 2017).

816 A.3.2 MLP AND NODE DECODER FINE-TUNING (NFT)

In contrast to MFT, the NFT strategy involves fine-tuning part of the and is recommended when 818 sufficient computational resources are available and the goal is to extract the maximum performance 819 from the model. In NFT, we gradually unfreeze the node decoder, enabling the model to more 820 effectively adapt to the new dataset. Specifically, we set a predefined epoch U at which unfreezing 821 begins, starting from the bottom layers of the node decoder. After every S epochs, additional layers 822 are unfrozen in a bottom-up manner, facilitating gradual transition to full finetuning of the model. 823 Concurrently, the learning rate is decayed by a factor of 1.5 each time a new layer is unfrozen, 824 ensuring controlled parameter updates. For all datasets, we tune the hyperparameters U and S, with 825 U set to 10, 20, or 30 epochs and S set to 5 or 10 epochs. This gradual unfreezing mitigates training 826 instability, as smaller perturbations are made to higher-level feature representations. As a result, NFT 827 allows for better adaptation, particularly for out-of-distribution (OOD) datasets, and is well suited for case when exploiting the capacity of pretrained models is critical. 828

829 830

831 832

833 834

835

836

837

838

839

844

845

846

847 848

849 850

851

852

857

815

817

# **B** ADDITIONAL DETAILS ON DATASETS

# **B.1** PRETRAINING DATASETS

The largest model (75M parameters) was trained on 80 real world and 72 synthetic datasets. The real world datasets and their characteristics are given in Table A3.

The synthetic datasets were created using the GraphWorld (Palowitch et al., 2022) using the Stochastic Block Model (Holland et al., 1983). The generator parameters are listed in Table A2. In the graph generation process, the *node homophily ratio* is varied. The homophily is given by the following formula:

$$\frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \frac{|\{(v, w) : w \in \mathcal{N}(v) \land y_v = y_w\}|}{|\mathcal{N}(v)|},$$

where  $\mathcal{V}$  denotes the set of all nodes in the graph,  $\mathcal{N}(v)$  denotes all the neighbors of an arbitrary node v, and  $y_v$  denotes the class membership of the node  $v \in \mathcal{V}$ . We classify datasets into *homophilic datasets* and *heterophilic datasets* based on the homophily score: datasets with homophily  $\geq 0.5$  are classified as *homophilic datasets* and *heterophilic datasets* otherwise.

# B.2 DETAILS ON SMALL AND MEDIUM SCALE DATASET

The small and medium scale datasets, as discussed in Section 4.2, were created by taking a random subset of the large dataset(80 real and 72 synthetic).

Bataset subset for small scale data: The following datasets were used to train models with small
 scale data: Wiki, BlogCatalog, Roman-empire, Minesweeper, Tolokers, Questions, Twitch-EN,
 Twitch-FR, Twitch-PT, Twitch-RU, DeezerEurope, GitHub, LastFMAsia, Airports-USA, Airports Europe, PolBlogs and EmailEUCore

**Dataset subset for medium scale data:** The following datasets were used to train models
with medium scale data: Wiki, BlogCatalog, Roman-empire, Minesweeper, Tolokers, Questions,
Twitch-EN, Twitch-FR, Twitch-PT, Twitch-RU, DeezerEurope, GitHub, LastFMAsia, Airports-USA,
Airports-Europe, PolBlogs and EmailEUCore, Reddit, Reddit2, Flickr, Yelp, PPI, Facebook, Amazonratings, Minesweeper, Twitch-DE, Twitch-ES, FacebookPagePage, Airports-Brazil, penn94, reed98,
amherst41, johnshopkins55, genius, CitationFull-CiteSeer, CitationFull-Cora-ML and CitationFull-PubMed

Parameter Name	Description	Values	
nvertex	Number of vertices in the graph.	[32, 500000]	
p/q ratio	The ratio of in-cluster edge probabil- ity to out-cluster edge probability.	[0.1, 10.0]	
avg. degree	The average expected degrees of the nodes.	[1.0, 20.0]	
feature center distance	The variance of feature cluster cen- ters, generated from a multivariate Normal.	[0.0, 5.0]	
num clusters	The number of unique node labels.	[2, 6]	
cluster size slope	The slope of cluster sizes when index-ordered by size.	[0.0, 0.5]	
power exponent	The value of the power law expo- nent used to generate expected node degrees.	[0.5, 1.0]	

# B.3 DETAILS ON SOCIAL AND BIOLOGY DOMAIN DATASETS

The social and biology datasets, as discussed in Section D.1 and Section 4.2, included the following subsets:

**Dataset subset for social domain:** The following datasets were used to train the social-specific model: fb-CMU-Carnegie49, Yelp,Wiki, BlogCatalog, Facebook, Twitch-DE, Twitch-EN, Twitch-ES, Twitch-FR, Twitch-PT, Twitch-RU, DeezerEurope, GitHub, FacebookPagePage, LastFMAsia, penn94, reed98, amherst41, johnshopkins55, genius and soc-pokec.

Bataset subset for biology domain: The following datasets were added as part of the biology domain to train the combined social and biology model: BZR, DD, DD199, DD21, DD242, DD244, DD349, DD497, DD6, DD68, DD687, DHFR, ENZYMES, ENZYMES118, ENZYMES123, EN-ZYMES295, ENZYMES296, ENZYMES297, ENZYMES8, KKI, OHSU, PROTEINS-full, Peking-1, Tox21\_p53, gene, proteins-all and PPI.

# **B.4** FINETUNING DATASETS

For our evaluations, we held out a number of datasets that are used for standard benchmarks in both larger scale node classification and heterophilic graphs.

B.4.1 HOMOPHILIC DATASETS

We use five real-world datasets, Amazon Computers and Amazon Photos (McAuley et al., 2015), Coauthor CS and Coauthor Physics (Sinha et al., 2015) and Obgn-Arxiv (Hu et al., 2020). Key statistics for the different datasets are listed in Table A3 in the finetuning-section. The experimental setup follows that of (Luo et al., 2022), where we split the dataset into development and test sets. All the hyperparameter tuning is done on the development set and the best models are evaluated on the test set. The runs are averaged over 20 random splits to minimize noise. We follow a 60:20:20% train/val/test split for the Amazon and Coauthor datasets. For Obgn-Arxiv we follow the experimental setup used in (Hu et al., 2020). The results for the Coauthor-Physics, Coauthor-CS, and Amazon-Photos obtained from in Table A6 have been sourced from (Liu et al., 2023a). The results for the Amazon-Comp dataset are taken from (Hoang et al., 2023) except for MLP which was obtained from (Luo et al., 2022).

9	1	9	
9	2	0	

Table A3: Pre-Training Datasets and their characteristics

010		Dataset	Number of Graphs	Nodes	Edges	Homophily Ratio	Average Degree	Node Features	Node Classes	Learning Rate
921		BA-1_10_60-L5	1	804	46410	0.2	115.45	1	5	0.0014
000		BA-2_24_60-L2	1	10693	639750	0.5	119.66	1	2	0.0087
922		BZR	405	35.75	76.71	0.42	0.07	1	53	0.0082
023		CL-100K-1d8-L9	1	92482	373989	0.11	8.09	1	9	0.00064
525		CL-10K-1d8-L5	1178	284.32	44896	0.2	8.98	1	5	0.00096
924		DD199	1	841	1902	0.067	4 52	1	20	0.00085
		DD21	1	5748	14267	0.07	4.96	1	40	0.00085
925		DD242	1	1284	3303	0.08	5.14	1	20	0.00042
026		DD244	1	291	822	0.074	5.65	1	20	0.00085
520		DD349	1	897	2087	0.05	4.65	1	20	0.00085
927		DD497	1	903	2453	0.06	5.43	1	20	0.0028
000		DD68	1	775	2093	0.072	4.97	1	20	0.0028
928		DD687	1	725	2600	0.06	7.17	1	20	0.0028
929		DHFR	756	42.43	89.09	0.32	0.04	3	53	0.0018
020		ENZYMES	600	32.63	124.27	0.67	0.09	18	3	0.0020
930		ENZYMES118	1	96	121	0.58	2.52	1	2	0.00087
004		ENZYMES123	1	90	127	0.52	2.82	1	2	0.0076
931		ENZYMES295	1	124	139	0.72	2.24	1	2	0.0078
932		ENZYMES297	1	122	149	0.65	2.44	1	2	0.0020
002		ENZYMES8	1	88	133	0.77	3.02	1	2	0.0076
933		ER-AvgDeg10-100K-L2	1	99997	499332	0.50	9.99	2	2	0.0049
024		ER-AvgDeg10-100K-L5	1	99997	499332	0.20	9.99	1	5	0.0013
934		KKI MSBC 21	83	26.96	90.84	0	0.39	1	189	0.0012
935		MSRC-21 MSRC-21C	209	40.28	193.20	0.61	0.13	1	24	0.0003
		MSRC-9	221	40.58	193.21	0.69	0.26	1	10	0.009
936		OHSU	79	82.01	399.32	0	0.56	1	189	0.0095
027		PLC-40-30-L5	1	11025	437979	0.2	79.45	1	5	0.0086
937		PLC-60-30-L2	1	117572	7045181	0.5	119.84	1	2	0.0013
938		PROTEINS-Tull Peking 1	1113	39.06	145.63	0.97	0.05	2	8	0.0063
		SW-10000-6-0d3-L2	1	10000	30000	0.5	6	1	2	0.00096
939		SW-10000-6-0d3-L5	1	10000	30000	0.2	6	1	5	0.0088
940		SYNTHETIC	300	100	392	0.18	0.16	1	8	0.0018
540	50	TerroristRel	1	881	8592	0.92	19.51	1	2	0.0033
941	in i	10X21_p33 fb_CMU_Carpagie40	1	153503	240067	0.62	4.09	1	46	0.00054
0.40	Tra	gene	1	1103	1672	0.4	3.03	1	2	0.012
942	-re-	proteins-all	1	43471	162088	0.66	7.46	1	3	0.00075
943	-	reality-call	1	27058	51200	0.9	15	1	2	0.0071
		Reddit	1	232965	114615892	0.76	983.98	602	41	0.0035
944		Reddit2 Flialar	1	232965	23213838	0.78	199.29	500	41	0.0035
9/5		Yeln	1	716847	13954819	0.51	20.10	300	1001	0.00031
545		Wiki	1	2405	17981	0.71	14.95	4973	17	0.0012
946		BlogCatalog	1	5196	17981	0.40	132.21	8189	6	0.0099
0.47		PPI	1	56944	1612348	0.63	56.63	50	121	0.0016
947		Facebook	1	4039	88234	0.99	43.69	1283	193	0.0011
948		Amazon-ratings	1	22002	186100	0.03	3.81 15.2	300	18	0.0074
0.10		Minesweeper	1	10000	78804	0.68	15.76	7	2	0.0088
949		Tolokers	1	11758	1038000	0.59	176.56	10	2	0.0022
050		Questions	1	48921	307080	0.84	12.55	301	2	0.0061
950		Twitch-DE	1	9498	315774	0.64	66.49	128	2	0.0023
951		Twitch-EN	1	/126 4648	123412	0.59	21.82	128	2	0.0010
050		Twitch-FR	1	6551	231883	0.54	70.79	128	2	0.0010
952		Twitch-PT	1	1912	64510	0.58	67.47	128	2	0.0012
953		Twitch-RU	1	4385	78993	0.63	36.02	128	2	0.0011
000		DeezerEurope	1	28281	185504	0.52	13.11	128	2	0.0070
954		GitHub FacebookPagePage	1	22470	342004	0.84	30.66	128	2	0.0065
055		LastFMAsia	1	7624	55612	0.87	14.59	128	18	0.0092
900		Airports-Brazil	1	131	1074	0.46	16.39	131	4	0.0013
956		Airports-Europe	1	399	5995	0.40	30.05	399	4	0.0015
		Airports-USA	1	1190	13599	0.69	22.85	1190	4	0.0092
957		PolBlogs	1	1490	19025	0.91	25.54	1	2	0.0013
958		nenn94	1	41554	25571	0.50	131 11	4814	42	0.0052
550		reed98	1	962	37624	0.52	78.22	745	2	0.0032
959		amherst41	1	2235	181908	0.53	162.78	1193	2	0.011
		johnshopkins55	1	5180	373172	0.55	144.08	2406	2	0.0025
960		genius Citation Full City C	1	421961	984979	0.62	4.67	12	2	0.00040
961		CitationFull-CiteSeer	1	4230	10674	0.95	5.04	2879	0 7	0.0011
001		CitationFull-PubMed	1	19717	88648	0.80	8,99	500	3	0.00087
962		soc-pokec	1	1632803	30622564	0.44	37.51	500	3	0.00019

<sup>1</sup> Multi label binary classification.

Table A4: Fine-Tuning Datasets and Their Characteristics

66		Table A4: 1	Table A4: Fine-Tuning Datasets and Their Characteristics					
67	Dataset	Number of Graphs	Nodes	Edges	Homophily Ratio	Average Degree	Node Features	Node Classes
68	Actor	1	7600	30019	0.21	7.89	932	5
00	Amazon-Computers	1	13752	4491722	0.77	71.51	767	10
69	Amazon-Photo	1	7650	238162	0.82	62.26	745	8
70	Coauthor-CS	1	18333	163788	0.80	17.86	6805	15
10	Coauthor-Physics	1	34493	495924	0.93	28.75	8415	5
71	Chameleon	1	2277	36101	0.23	31.70	2325	5

# 972 B.4.2 HETEROPHILIC DATASETS

We use five real-world datasets with graphs that have a homophily level  $\leq 0.30$ , Texas, Wisconsin and Actor (Pei et al., 2020) and Chameleon and Squirrel (Rozemberczki et al., 2021). Key statistics for the different datasets are listed in Table A3 in the finetuning-section. We follow the experimental setup in (Pei et al., 2020), and use the same 10 train/val/test splits that are provided. The results for GCN based methods and heterophily based methods in Table A7 have been taken from (Azabou et al., 2023), and the results for transformer based methods have been taken from (Liu et al., 2023a)

B.5 STANDARD HYPERPARAMETER SEARCH GRID FOR BASELINES

The hyperparameter search space grid used for tuning baselines for Table 1 is detailed in Table A5.

Hyperparameter	Туре	Range
Hidden Dim	Categorical	{16, 32, 64, 128}
Depth	Categorical	$\{1, 2\}$
Dropout	Uniform	[0.0, 0.9]
Learning Rate	Log uniform	[5e-5, 5e-1]
Weight Decay	Log uniform	[1e-5, 1e-2]

Table A5: Hyperparameter Search Space

# C ADDITIONAL DETAILS ON MULTI-GRAPH TRAINING

One key aspect of our work is testing scale. Thus, to build a model across large amounts of diverse graph data, we developed a number of approaches for efficient training and multi-GPU usage.

Figure A1 shows an ablation study the epoch time for various GPU optimizations we have proposed in Section 2.2. The epoch time was calculated using the medium-sized model with 18M parameters, as detailed in Appendix A.1.

Note: Removing chaining made it impossible to run the largest model (75M parameters) with our available computational resources (8 A40 GPUs). Therefore, we performed the ablation using the medium-sized model. This highlights the significance of our optimization techniques, which enabled us to scale up and run such large models efficiently.

1008

1015

980

981 982

983 984

997 998

999

# 1009 C.1 DISTRIBUTEDSSSAMPLER

In designing this sampler, we prioritized ensuring that
it neither introduces bias into the data sampling process nor alters the distribution of the graphs from the
datasets. Its primary function is to enhance batch construction and distribution across GPUs.

First, the sampler defines a set of N buckets with a 1016 fixed node budget B, where N can be the number of 1017 GPUs and B is the node-level batch size. The graphs 1018 (across all GPUs) are sorted in descending order based 1019 upon their size. The sampler then employs a bidirec-1020 tional filling strategy within the buckets. The distri-1021 bution process, as described in Algorithm 1 involves distributing graphs in a snake-like pattern, initially fill-1023 ing from right to left, then switching to left to right and so on. When a graph is added to a bucket, it uses 1024 up part of the budget, equal to its size. This method 1025 effectively pairs larger graphs with smaller ones in



Figure A1: Ablation for GPU optimizations: Epoch time in minutes on removing various gpu optimizations proposed for GRAPHFM

subsequent passes, preventing the concentration of multiple large graphs on the same GPU, thus achieving efficient load balancing and uniform GPU utilization. Figure A2A shows an overview of how the sampler distributes the graphs into buckets. We find that stability is improved with a larger number of buckets N (Figure A2B). When the number of GPUs is fixed, we can achieve a larger Nby using gradient accumulation, which artificially increases the number of buckets by a factor equal to the number of accumulation steps, without biasing the sampling process.

1032 1033 Algorithm 1 Distribute graph nodes into virtual GPU buckets 1034 1: input: Batch size B, Bucket count N, Graphs in the dataset  $\mathcal{G} = \{\mathcal{G}_0, \mathcal{G}_1, \ldots\}$ , Subgraphs 1035 sampled for this minibatch  $\mathcal{G}^m = \{\mathcal{G}_0^m, \mathcal{G}_1^m, \ldots\}$ 2: precondition:  $\sum_{i} |\mathcal{G}_{i}^{m}| = N \times \overset{\sim}{B}$ 1036 1037 3: initialize: 4:  $buckets \leftarrow array of N empty arrays$ # will store subgraphs in each bucket 1038 5:  $counts \leftarrow array of N zeroes$ # will store number of nodes in each bucket 1039 6:  $b \leftarrow 0$ # bucket index 1040 7: # direction  $d \leftarrow 1$ 1041 Sort  $\mathcal{G}^m$  according to node-counts in  $\mathcal{G}$ , largest graph goes first 8: 1042 9: for all  $\mathcal{G}_i^m$  in  $\mathcal{G}^m$  do 1043 while  $|\mathcal{G}_i^m| > 0$  do 10: 1044 if counts[b] < B then 11: 1045 12: # insert a part of  $\mathcal{G}_i^m$  into bucket b  $n \leftarrow \min(|\mathcal{G}_i^m|, \ B - counts[b])$ 1046 13: 1047 14:  $counts[b] \leftarrow counts[b] + n$ append first n nodes of  $\mathcal{G}_i^m$  to buckets[b]1048 15: 1049 16: remove first n nodes from  $\mathcal{G}_i^m$ end if 17: 1050 # go to the next bucket, switching direction at the boundaries 18: 1051 19:  $b \leftarrow b + d$ 1052 if  $b \ge N$  or b < 0 then 20: 1053 21:  $d \leftarrow -d$ 1054  $b \leftarrow b + d$ 22: 1055 23: end if 1056 24: end while 1057 25: end for 1058 26: return buckets 1059 в 1061 Mini-batch 100 GPU Memory Usage (%) 1062 ..... 4 90 1063 3 80 1064 70 Largest and smallest aphs share GPU buckets 6 3 60 4 1067 50 N = 64 b = 1376 N = 4 b = 960 1068 baseline b = 640 1069 Figure A2: Multi-GPU utilization: A: A diagram visualizing our sample distribution strategy. B: GPU memory 1070 utilization during distributed training when using the default batch sampler vs. our DistributedSSSampler for 1071 N=4 and N=64 buckets.

1072 1073

# 1074 C.2 GRAPHSAINT RANDOM WALK SAMPLER

Efficient neighborhood sampling for large graphs is crucial for our node decoder, as traditional methods for k-hop neighborhood sampler often become computationally prohibitive with the increasing size and complexity of the graph data. To overcome these limitations, we have adopted the GraphSAINT Random Walk Sampler (Zeng et al., 2019), specifically designed for efficient sampling in large-scale graphs.

# 1080 C.3 RAM OPTIMIZATION IN MULTI-GPU ENVIRONMENTS

In multi-GPU training environments, efficient use of system memory is crucial, especially when handling large graph datasets. Traditional approaches lead to substantial memory redundancy, as each GPU process typically loads a complete dataset into system RAM. This results in each process duplicating the dataset in system memory, leading to inefficient memory usage and potential system overload.

To address this, we utilize a shared memory management approach using Python's multiprocessing.Manager() to coordinate dataset access across multiple GPU processes.
 This method ensures that each dataset is loaded into RAM only once, regardless of the number of GPUs, thereby avoiding duplication and conserving memory resources.

1091 1092 1093

# D ADDITIONAL EXPERIMENTS

1094 1095 1096

# D.1 SEPARATING PRETRAINING DATASETS INTO DIFFERENT DOMAINS

We further stratified our pretraining dataset to invetigate the effects of cross-domain training, and created three models that contained: (i) graph datasets from "social domains" including product graphs and citation networks (1.3M tokens), (ii) both the social datasets and all biological graphs in the dataset (Bio+Soc, 2M tokens), and (iii) compare with our model trained on all data including synthetic graphs (7.3M tokens).

When comparing graph features across social and biological domains, we found distinct structural differences: biological datasets generally exhibited higher levels of heterophily, lower average degree, and fewer edges, whereas social graphs showed more homophily, higher degrees, and denser connections (Figure A4B). Synthetic graphs added a wide range of characteristics, particularly increasing the number of heterophilic graphs used in pretraining, which contributed to a broader diversity of features (Figure A4A).

All three models were then fine-tuned on four homophilic datasets (coauthor-CS, coauthor-physics, amazon-photos, and amazon-computers) and five heterophilic datasets (Texas, Wisconsin, Actor, Squirrel, and Chameleon) held out for fine-tuning.

1114 As shown in Figure A3 we find that incorporating biology datasets despite being seemingly unrelated to the 1115 target domain-improved performance on the OOD 1116 datasets. This suggests that knowledge learned from 1117 the biology domain positively impacts performance in 1118 seemingly unrelated domains. Furthermore, adding all 1119 available datasets, including synthetic graphs, boosted 1120 performance even more, indicating that diversity (not 1121 just domain specific data), is the key to improving 1122 generalization.



Figure A3: **Domain Scaling**: Average accuracy across OOD datasets (using MFT) for models trained on different subsets of data

- 1123 1124 1125
- 1126

D.2 SCALING ANALYSIS BREAKDOWN FOR DIFFERENT TEST DATASETS

The main text shows the average scaling. We break down the scaling performance for different datasets (Figure A5). All of the datasets benefit from scale, with more difficult datasets benefiting more from scaling the model size and dataset.

1130

1131 1132 D.3 RANKING OF DIFFERENT MODELS

1133

We visualize the ranking of the different models (Figure A6).



Figure A4: Characteristics of graph datasets used to train GraphFM: From left to right, we compute the histograms of the homophily ratio, average degree, number of nodes and number of edges of all 152 graphs used during training. The homophily ratio provides a measure of how frequently a node is directly connected to other nodes from the same class.



Figure A5: Accuracy as the model and dataset size are increased. Results are shown for four datasets, Chameleon and Wisconsin (heterophilic), and Coauthor Physics and Amazon Photo (homophilic).

1180 D.4 ADDITIONAL BASELINES

The main text presents a comparison of GRAPHFM with baselines that are more consistently reported across the literature. Table A6 and Table A7 provides additional baselines for all the OOD datasets.
1184
1185
1186
1187



Figure A6: Mean rank of various models accross 10 OOD datasets (lower is better).

Table A6: *Results on node classification tasks for large graph datasets*. We report the accuracy (%) with standard deviation over 10 splits (OOM indicates Out of Memory).

Method	Photo	Physics	CS	ogbn-arxiv	Comp
GCN	85.94±1.18	95.38±0.20	94.06±0.16	$70.40 {\pm} 0.10$	$89.47 \pm 0.46$
GatedGCN	$57.84{\pm}14.6$	$95.89 {\pm} 0.21$	$89.94{\pm}2.24$	62.71±1.76	-
APPNP	$84.71 \pm 1.25$	$95.04{\pm}0.31$	$87.49 {\pm} 0.48$	$70.20 {\pm} 0.16$	$90.18\pm0.17$
GCNII	$67.06 \pm 1.74$	$94.88 {\pm} 0.32$	$84.23 \pm 0.78$	$69.78 \pm 0.16$	-
GAT	$87.13 \pm 1.00$	$95.14 \pm 0.28$	$93.61 \pm 0.14$	$67.56 \pm 0.12$	$90.78 \pm 0.13$
GATv2	$81.52 \pm 3.23$	$95.02 \pm 0.32$	$88.46 \pm 0.61$	$68.84 \pm 0.13$	-
SuperGAT	85.83±1.29	95.11±0.26	88.11±0.43	$66.99 \pm 0.07$	-
		Heterophily-ba	ased methods		
MLP	$88.66 {\pm} 0.85$	$95.12{\pm}0.26$	$92.99 {\pm} 0.51$	$52.63 {\pm} 0.12$	84.63
MixHop	$93.24 {\pm} 0.59$	$96.34 \pm 0.22$	$93.88 {\pm} 0.63$	$70.83 \pm 0.30$	-
H2GCN	$91.56 {\pm} 0.70$	96.28±0.13	$94.02 \pm 0.31$	$68.29 \pm 0.67$	$89.33 \pm 0.27$
FAGCN	$87.53 \pm 0.75$	$95.86 \pm 0.12$	$91.82{\pm}0.54$	$66.12 \pm 0.02$	-
GPRGNN	$92.27 \pm 0.44$	$96.06 \pm 0.21$	$93.60 \pm 0.36$	$68.28 \pm 0.21$	$89.32 \pm 0.29$
	Gra	ph Transform	er-based metho	ods	
SAN	94.17±0.65	$96.83 {\pm} 0.18$	94.16±0.36	69.17±0.15	$89.83\pm0.16$
Graphormer	$85.20 \pm 4.12$	OOM	OOM	OOM	OOM
LiteGT	-	OOM	$92.16 \pm 0.44$	OOM	-
UniMP	$92.49 \pm 0.47$	$96.82 \pm 0.13$	$94.20 \pm 0.34$	$73.19 \pm 0.18$	-
DET	$91.44 \pm 0.49$	$96.30 \pm 0.18$	$93.34 \pm 0.31$	$55.70 \pm 0.30$	-
NAGphormer	$94.64 \pm 0.60$	$96.66 \pm 0.16$	$95.00 \pm 0.14$	$68.21 \pm 0.021$	$91.22 \pm 0.14$
GRAPHFM -MFT	$93.01 \pm 1.82$	$96.64 \pm 0.17$	$95.19 \pm 0.21$	$65.29 \pm 0.16$	$89.95 \pm 0.83$
GRAPHFM -NFT	94.37±0.35	96.77±0.12	$95.24 \pm 0.18$	$70.01 \pm 0.18$	$90.07 \pm 0.21$

Table A7: *Results on node classification tasks for heterophilic graphs.* We report the test accuracy across many heterophilic graph benchmark datasets. The standard deviation is reported across 10 train/test splits.

Method	Texas	Wisconsin	Actor	Squirrel	Chameleon				
		GCN-based	nethods						
GCN	$55.14\pm5.16$	$51.76\pm3.06$	$27.32 \pm 1.10$	$31.52\pm0.71$	$38.44 \pm 1.9$				
GAT	$52.16\pm 6.63$	$49.41 \pm 4.09$	$27.44 \pm 0.89$	$36.77 \pm 1.68$	$48.36 \pm 1.5$				
GraphSAGE	$82.43 \pm 6.14$	$81.18\pm5.56$	$34.23\pm0.99$	$41.61\pm0.74$	$58.73 \pm 1.6$				
Heterophily-based methods									
MLP	$80.81 \pm 4.75$	$85.29 \pm 3.31$	$36.63\pm0.70$	$28.77 \pm 1.56$	$46.21 \pm 2.9$				
HH-GCN	$71.89 \pm 3.46$	$79.80 \pm 4.30$	$35.12 \pm 1.06$	$47.19 \pm 1.21$	$60.24 \pm 1.9$				
HH-GAT	$80.54 \pm 4.80$	$83.53 \pm 3.84$	$36.70\pm0.92$	$46.35 \pm 1.86$	$61.12 \pm 1.8$				
HH-GraphSAGE	$85.95 \pm 6.42$	$85.88 \pm 3.99$	$36.82\pm0.77$	$45.25 \pm 1.52$	$62.98 \pm 3.3$				
MixHop	$77.84 \pm 7.73$	$75.88 \pm 4.90$	$32.22\pm2.34$	$43.80 \pm 1.48$	$60.50\pm2.5$				
GGCN	$84.86 \pm 4.55$	$86.86 \pm 3.29$	$37.54 \pm 1.56$	$55.17 \pm 1.58$	$71.14 \pm 1.8$				
$H_2GCN$	$84.86\pm7.23$	$87.65 \pm 4.98$	$35.70 \pm 1.00$	$36.48 \pm 1.86$	$60.11\pm2.1$				
	Graj	ph Transformer	-based methods						
SAN	$60.17 \pm 6.66$	$51.37\pm3.08$	$27.12\pm2.59$	$39.92\pm2.14$	$44.32 \pm 1.7$				
UniMP	$73.51\pm8.44$	$79.60\pm5.41$	$35.15\pm0.84$	-	-				
ET	$56.76 \pm 4.98$	$54.90\pm6.56$	$28.94 \pm 0.64$	-	-				
NAGphormer	$63.51 \pm 5.85$	$62.55\pm 6.22$	$34.33\pm0.94$	$49.93 \pm 0.07$	$57.39\pm0.0$				
Gapformer	$80.27 \pm 4.01$	$83.53\pm3.42$	$36.90\pm0.82$	-	-				
GRAPHFM -MFT	$80.81 \pm 2.76$	$83.13\pm2.35$	$36.29\pm0.63$	$42.80 \pm 1.54$	$58.64 \pm 1.2$				
GRAPHFM -NFT	$82.16\pm3.24$	$83.62\pm3.21$	$38.01 \pm 1.07$	$42.98 \pm 1.62$	$59.12 \pm 1.6$				