

AnyGraph: Graph Foundation Model in the Wild

Lianghao Xia

Harbin Institute of Technology (Shenzhen)
The University of Hong Kong
xialh@hit.edu.cn

Chao Huang

The University of Hong Kong
chuang7@hku.hk

Abstract

The ubiquity of text-attributed graph data has highlighted the need for graph learning models with exceptional generalization across diverse textual and structural contexts. Current approaches struggle to extract generalizable insights from heterogeneous graph data, requiring extensive fine-tuning and limiting versatility across domains. In this work, we propose AnyGraph, a unified graph foundation model designed to handle key challenges: i) **Structure Heterogeneity** - addressing distribution shift in graph structural patterns; ii) **Feature Heterogeneity** - handling diverse textual representations; iii) **Fast Adaptation** - efficiently adapting to new graph-text domains. We build AnyGraph upon a Graph Mixture-of-Experts (MoE) architecture with a lightweight expert routing mechanism that effectively manages cross-domain distribution shift. Extensive experiments on 38 diverse datasets demonstrate AnyGraph’s strong zero-shot performance across domains with significant distribution shift, validating its fast adaptation ability and scaling law emergence. Our model is open-sourced and available at: <https://github.com/HKUDS/AnyGraph>.

1 Introduction

The growing ubiquity of relational data in the form of graphs has underscored the pressing need for advanced graph-based methods that excel at generalization, particularly in natural language processing tasks (Fey et al., 2024; Jin et al., 2020). As real-world applications of graph-structured data continue to proliferate across diverse domains, including social networks, academic networks, transportation systems, and biological networks, the ability of graph-based NLP models to effectively handle distribution shifts and adapt to new graph domains has become increasingly crucial (Zhang et al., 2023; Zhao et al., 2024; Mao et al., 2024). Developing models with robust zero-shot learning

performance and fast adaptation capabilities can unlock transformative opportunities for leveraging the rich insights encoded within graph data for language understanding tasks.

The field of graph learning has seen significant advancements in recent years, largely driven by the power of Graph Neural Networks (GNNs) (Liu et al., 2022; Xiao et al., 2021; Li et al., 2021). However, the state-of-the-art models often fall short when it comes to truly generalizable performance. Existing approaches are heavily reliant on arduous fine-tuning processes, making them ill-equipped to handle the diverse array of graph structures and distributions encountered in real-world applications. This inability to adapt swiftly and seamlessly to novel graph domains poses a critical barrier to the widespread adoption of graph learning technologies. Therefore, addressing this challenge is of high importance if we are to fully harness the transformative potential of graph-based insights.

Inspired by successful foundation models in vision and language (Wang et al., 2022, 2023), graph foundation models hold immense potential for learning transferable representations from diverse domains. However, building effective graph foundation models faces several key challenges:

(i) **Structure Heterogeneity**. Graph datasets exhibit diverse structural properties, from varying node degree distributions to complex hierarchical arrangements. These structural variations significantly impact model performance and generalization, requiring unified foundational models that can robustly handle diverse graph topologies.

(ii) **Feature Heterogeneity**. Graph features span textual, visual, categorical, numerical, and other multi-modal content with varying dimensionality and semantics across domains. Social graphs may include textual content and demographics, while molecular graphs feature atomic compositions and bond types, necessitating models capable of handling diverse feature representations.

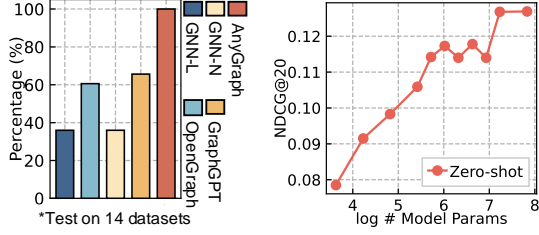


Figure 1: The zero-shot generalizability (left) and performance scaling trend (right) of our AnyGraph model.

(iii) **Fast Adaptation.** Graph foundation models must efficiently adapt to new datasets and domains without extensive retraining. The ideal model should quickly adjust to structural and distributional characteristics of unseen graph data, generalizing across diverse application scenarios from user behavior to biological systems.

(iv) **Scaling Laws.** Successful foundation models for visual data (Cherti et al., 2023) and natural language data (Muennighoff et al., 2024) exhibit scaling laws where performance improves with model size and training data. Graph foundation models should harness this scaling phenomenon to unlock unprecedented capability and generalization as data amount and model complexity grow.

The Presented Work. To tackle the above challenges, our AnyGraph model is built upon a Mixture-of-Experts (MoE) architecture, which allows for effective handling of both the in-domain and cross-domain distribution shift in structure-level and feature-level. The proposed graph MoE paradigm empowers AnyGraph to learn a diverse ensemble of graph experts, each tailored to specific structural characteristics. This enables the model to effectively manage the distribution shift in graph topologies. Furthermore, the MoE architecture facilitates fast adaptation of AnyGraph. Rather than relying on a single, fixed-capacity model, the Graph MoE can efficiently tailor some of its expert networks to capture distinct characteristics of new graph data. A lightweight graph expert routing mechanism also allows AnyGraph to quickly identify and activate the most relevant experts for a given input graph, without requiring extensive retraining or fine-tuning across the entire model. The key contributions include:

- **Graph MoE Architecture Design.** We propose a novel graph MoE architecture that addresses heterogeneity in real-world graph data through dynamic expert selection. Unlike fixed-capacity models (Chen et al., 2024; Liu et al., 2024; Li

et al., 2024), our approach flexibly adjusts to diverse datasets by selecting appropriate experts, avoiding interference and catastrophic forgetting.

- **Superior Cross-domain Generalization.** Extensive experiments demonstrate AnyGraph’s strong generalization across diverse graph tasks and domains. Results showcase significant improvements over existing models in both predictive performance and robustness to distribution shift.
- **Efficient Domain Adaptation.** AnyGraph’s dynamic expert selection mechanism enables swift adaptation to new graph domains without extensive retraining. By routing inputs through relevant experts, AnyGraph quickly activates specialized networks, demonstrating rapid convergence and exceptional performance.
- **Emergent Scaling Properties.** Extensive experiments on 38 diverse datasets reveal that AnyGraph follows scaling laws where performance improves with model size and training data. The model exhibits emergent generalization abilities with scaling, a critical property largely overlooked in prior graph learning research.

2 Preliminaries

Graph-Structured Data. A graph \mathcal{G} consists of a set of nodes $\mathcal{V} = \{v_i\}$ and a set of edges $\mathcal{E} = \{(v_i, v_j)\}$. In many cases, each node v_i is associated with a feature vector $\mathbf{f}_i \in \mathbb{R}^{d_0}$. To efficiently utilize such graph-structured data, the link information is typically recorded using an adjacency matrix $\mathbf{A} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$. Each element $a_{i,j}$ of \mathbf{A} is either 1 or 0, indicating whether there is an edge from node v_i to v_j . Additionally, the feature vectors of the nodes are usually represented by a feature matrix $\mathbf{F} \in \mathbb{R}^{|\mathcal{V}| \times d_0}$, where each row corresponds to a node’s feature vector.

Graph Foundation Models (GFMs) excel at generalizing to unseen graph data that differs from training datasets in feature spaces and node/edge semantics. Given training graphs $\mathbb{S} = \{\mathcal{G}_s\}$ with labels \mathcal{Y}_s and test graphs $\mathbb{T} = \{\mathcal{G}_t\}$ with labels \mathcal{Y}_t , a GFM f_{Θ} with parameters Θ can be trained using a differentiable objective \mathcal{L} and evaluated using criterion \mathcal{C} to measure downstream task accuracy.

$$\begin{aligned} & \arg \max_{f, \mathcal{L}} \sum_{\mathcal{G}_t} \mathcal{C}(f_{\Theta}(\mathcal{G}_t), \mathcal{Y}_t) \\ \Theta = & \arg \min_{\Theta} \sum_{\mathcal{G}_s} \mathcal{L}(f_{\Theta}(\mathcal{G}_s), \mathcal{Y}_s) \quad (1) \end{aligned}$$

The formulation shows two key requirements for GFMs: **i)** an architecture (f) capable of encoding diverse features and structures, and **ii)** a training process (\mathcal{L}) that effectively learns from diverse data to find optimal parameters Θ . Our AnyGraph addresses these through a MoE architecture with automated routing and graph augmentation, trained on diverse graphs with different features.

3 Methodology

AnyGraph aims to address both cross-domain and in-domain heterogeneity in graph structures and node features, while enabling fast adaptation to new data. The framework is depicted in Fig. 2.

3.1 MoE Architecture of AnyGraph

Addressing Cross-domain Graph Heterogeneity. To model heterogeneous graph patterns across domains, AnyGraph employs a MoE architecture consisting of multiple graph expert models, each responsible for handling graphs with specific characteristics. An automated routing algorithm is designed to assign input graph data to the most competent expert model for training and prediction. Specifically, the AnyGraph framework can be denoted as $\mathcal{M} = (f_{\Theta_1}, f_{\Theta_2}, \dots, f_{\Theta_K}, \psi)$, where K denotes the number of experts. For an input graph \mathcal{G} , the routing algorithm ψ firstly identifies the most competent expert model, which is then used for predicting the graph data, as follows:

$$\hat{y}_{i,j} = \hat{\mathbf{e}}_i^\top \hat{\mathbf{e}}_j, \quad \hat{\mathbf{E}} = f_{\Theta_k}(\mathcal{G}), \quad k = \psi(\mathcal{G}) \quad (2)$$

where each expert model f_{Θ_k} can be viewed as a projection from the graph space to a node embedding space with uniquely trained parameters Θ_k . And $\hat{y}_{i,j}$ represents the dot-product-based prediction of whether the entity v_i should be related to the entity v_j . Here, v_i and v_j could be vanilla graph nodes, class labels, or graph labels, enabling link prediction, and node/graph classification tasks.

Graph Expert Routing Mechanism. Inspired by the effectiveness of graph self-supervised learning tasks (Jin et al., 2022), we propose measuring the competence of expert models on specific graph datasets using the models’ self-supervised learning loss values. Specifically, for an input graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, the routing mechanism ψ calculates the dot-product-based relatedness scores for some positive edges $(v_{c_1}, v_{p_1}), \dots, (v_{c_S}, v_{p_S}) \in \mathcal{E}$ and analogously calculates the relatedness scores for some sampled negative node pairs

$(v_{c_1}, v_{n_1}), \dots, (v_{c_S}, v_{n_S}) \notin \mathcal{E}$. The following score difference is then calculated as the competence indicator φ_k for the k -th expert model regarding the input graph \mathcal{G} :

$$\varphi_k = \frac{1}{S} \cdot \sum_{s=1}^S \sigma(\hat{\mathbf{e}}_{c_s}^\top \hat{\mathbf{e}}_{p_s} - \hat{\mathbf{e}}_{c_s}^\top \hat{\mathbf{e}}_{n_s}) \quad (3)$$

where $\sigma(\cdot)$ represents sigmoid activation, which constrains the competence score to the range of (0, 1). This prevents the few outlier cases where the non-activated score difference is excessively large or small, which could otherwise distort the results. **Training Frequency Regularization.** Though being empirically accurate in measuring models’ competence using the above competence score, this method tends to result in a winner-takes-all sub-optimal situation. In this scenario, a single model, or very few models, is predominantly selected as the most competent expert and is used to handle almost all input graphs. These models generally receive more or better training samples in the early training stages, giving them an advantage over other experts. Consequently, subsequent training samples are also mostly assigned to them due to their performance advantages, ultimately causing other experts to remain largely untrained.

This situation contradicts our motivation of using different expert models to learn different subsets of graph modeling knowledge. To this end, we propose a training frequency regularization approach that recalibrates the competence score as follows:

$$\varphi'_k = \varphi_k \cdot \left(\left(1 - \frac{m_k}{\sum_{k'} m_{k'}} \right) \cdot \rho + 1.0 - \frac{\rho}{2} \right) \quad (4)$$

where φ'_k represents the recalibrated routing score for the k -th expert model f_{Θ_k} , based on the number of previously assigned training steps m_k for $k = 1, \dots, K$. The notation ρ refers to a hyper-parameter for the recalibration scale. A larger ρ results in a greater adjustment to the competence score φ_k . With this additional step, the expert routing mechanism will assign more training instances to the less trained expert models, thereby preventing the aforementioned winner-takes-all situation. **Fast Adaptation Capabilities of AnyGraph.** With the MoE architecture and routing mechanism, the training and inference process of AnyGraph is conducted by only one expert model. This approach consumes only $1/K$ of the computational and memory resources required for predictions and optimization, compared to other non-MoE graph

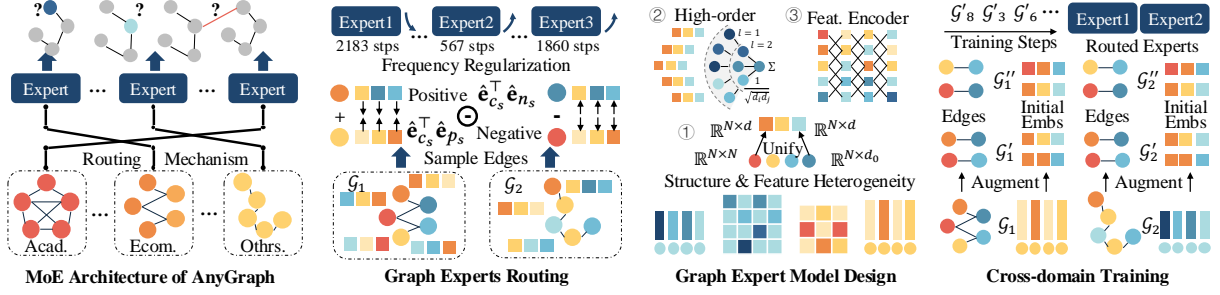


Figure 2: The overall model architecture of the proposed AnyGraph framework.

foundation models based on complex networks like transformers. This enables fast adaptation for AnyGraph when encountering new data.

3.2 Adaptive and Efficient Graph Experts

Addressing In-domain Graph Heterogeneity. To handle graph data with different adjacency and feature dimensionalities, the expert models of our AnyGraph employ a structure and feature unification process. Adjacency matrices and node features of varying sizes are both mapped into initial node embeddings of fixed dimensionality using a unified mapping process. Inspired by the effectiveness of singular value decomposition (SVD) in extracting important latent features (Cai et al., 2023), we utilize SVD for this unified mapping process:

$$\begin{aligned} \mathbf{U}_A, \Lambda_A, \mathbf{V}_A &= \text{SVD}(\tilde{\mathbf{A}}), \quad \mathbf{U}_F, \Lambda_F, \mathbf{V}_F = \text{SVD}(\mathbf{F}) \\ \mathbf{E}_0 &= \text{LN} \left(\mathbf{U}_A \sqrt{\Lambda_A} + \mathbf{V}_A \sqrt{\Lambda_A} + \text{Flip}(\mathbf{U}_F \sqrt{\Lambda_F}) \right) \end{aligned} \quad (5)$$

Here, $\mathbf{U}_A, \mathbf{U}_F \in \mathbb{R}^{|\mathcal{V}| \times d}$ and $\mathbf{V}_A \in \mathbb{R}^{|\mathcal{V}| \times d}, \mathbf{V}_F \in \mathbb{R}^{d_0 \times d}$ refer to the d -dimensional features obtained through SVD of the Laplacian-normalized adjacency matrix $\tilde{\mathbf{A}}$ and the node feature matrix \mathbf{F} , respectively. If the dimensionality of $\tilde{\mathbf{A}}$ or \mathbf{F} is less than d , SVD uses a smaller rank d' equal to the smallest dimensionality of $\tilde{\mathbf{A}}/\mathbf{F}$, and the remaining dimensions are padded with zeros up to d . $\text{LN}(\cdot)$ denotes layer normalization.

Due to the nature of SVD, the dimensions of these features ($\mathbf{U}_*, \mathbf{V}_*$) are ranked from the most important to the least important, corresponding to the descending eigenvalues in the diagonal matrices Λ_A and Λ_F . In light of this characteristic, we propose to better preserve the most important feature dimensions for both $\tilde{\mathbf{A}}$ and \mathbf{F} . In particular, the function $\text{Flip}(\cdot)$ reverses the d dimensions of each row for the SVD features of \mathbf{F} , such that the important features of $\tilde{\mathbf{A}}$ are aligned with the less important features of \mathbf{F} , and vice versa.

High-order Connectivity Injection. A non-trainable layer normalization $\text{LN}(\cdot)$ is applied for numerical stability. The initialized embeddings, denoted as $\mathbf{E}_0 \in \mathbb{R}^{|\mathcal{V}| \times d}$, have consistent representation dimensionality and relatively stable semantics across datasets. To better preserve the multi-hop connection information into the initial embeddings, AnyGraph adopts a simplified GCN without parameters (Wu et al., 2019) for \mathbf{E}_0 as follows:

$$\mathbf{E}_1 = \sum_{l=1}^L \mathbf{E}_0^{(l)}, \mathbf{E}_0^{(l)} = \tilde{\mathbf{A}} \cdot \mathbf{E}_0^{(l-1)}, \mathbf{E}_0^{(0)} = \mathbf{E}_0 \quad (6)$$

Efficient and Strong Feature Encoder. To achieve efficiency while retaining the capacity to encode graph features, our graph experts are configured by deep multi-layer perceptron (MLP) networks. Specifically, the final node embeddings given by an expert model is calculated as follows:

$$\bar{\mathbf{E}}^{(l+1)} = \text{LN} \left(\delta \left(\sigma \left(\bar{\mathbf{E}}^{(l)} \mathbf{W} + \mathbf{b} \right) \right) + \bar{\mathbf{E}}^{(l)} \right) \quad (7)$$

The final embeddings are denoted as $\hat{\mathbf{E}} = \bar{\mathbf{E}}^{(L)} \in \mathbb{R}^{|\mathcal{V}| \times d}$, where L' represents the number of fully-connected layers. And $\bar{\mathbf{E}}^{(0)}$ is initialized by the aforementioned embeddings \mathbf{E}_1 . Each layer of our MLP module comprises a linear transformation $\mathbf{W} \in \mathbb{R}^{d \times d}$ and bias $\mathbf{b} \in \mathbb{R}^d$, followed by a ReLU non-linear activation $\sigma(\cdot)$, a dropout layer $\delta(\cdot)$, a residual connection, and layer normalization.

Multiple Simple Experts as Strong Encoder. It is worth noting that each graph expert in AnyGraph adopts a very simple learnable network, foregoing the capacity to mine complex hidden relations like those in heavy graph neural networks such as GATs (Veličković et al., 2018) and Graph Transformers (Hu et al., 2020). This is because AnyGraph employs a MoE architecture, where each expert is expected to handle only a sub-domain of all graph data through simple feature transformations. Therefore, no complex models are needed

to accommodate different types of graphs within a single network. Compared to other graph foundation models relying on a single heavy network, this approach further accelerates training and inference.

3.3 Efficient Cross-domain Model Training

To maximize the cross-graph generalization capabilities of AnyGraph, the training samples from different datasets are mixed together and randomly shuffled during the model training process. Each batch of training samples is composed of the following information:

$$\begin{aligned} \mathcal{S} &= \left(\{(v_{c_b}, v_{p_b}) | b \in B\} \subset \mathcal{E}_{\mathcal{G}_s}, \right. \\ &\quad \mathbf{E}_1 = \text{InitialEmbed}(\mathcal{G}_s), \\ &\quad \left. f_{\Theta_k} \text{ where } k = \psi(\mathcal{G}_s) \right) \end{aligned} \quad (8)$$

Inspired by the effectiveness of link-wise graph pre-training tasks (Jin et al., 2022), we utilize link prediction as the training task. Here, (v_{c_b}, v_{p_b}) denotes the positive edges for link prediction, and B denotes the batch size. To facilitate batch training, each training batch involves only one training graph \mathcal{G}_s . The initial node embeddings \mathbf{E}_1 and the most competent expert model f_{Θ_k} are preprocessed in advance to accelerate the training. Specifically, the loss function for AnyGraph is as follows:

$$\mathcal{L} = \sum_S \sum_{b \in B} -\frac{1}{B} \log \frac{\exp(\hat{y}_{c_b, p_b} - \hat{y}_{\max})}{\sum_{v_n \in \mathcal{V}_{\mathcal{G}_s}} \exp(\hat{y}_{c_b, n} - \hat{y}_{\max})} \quad (9)$$

This objective maximizes the prediction scores for positive samples (v_{c_b}, v_{p_b}) and minimizes the predictions for all possible node pairs between v_{c_b} and all nodes v_n . To avoid numerical instability, we subtract the batch-specific maximum score, \hat{y}_{\max} , from all prediction scores. More training details are presented in Appendix A.1, including data augmentation methods and complexity analysis.

4 Evaluation

4.1 Experimental Settings

Experimental Datasets. For a comprehensive evaluation of the cross-domain graph generalizability, we employ a total of **38** datasets. These datasets span a wide range of domains, including e-commerce (e.g. user interactions and product-wise relations), academic graphs (e.g. citation and collaboration networks), biological information networks (e.g. relations among drugs and proteins),

and other domains like email networks, website networks, trust networks, and road networks.

Dataset Groups. We set up different dataset groups and conduct cross-dataset evaluations on these groups. Specifically, all datasets are divided into two cross-domain groups, **Link1** and **Link2**, which have a similar number of total edges and a similar number of domain-specific edges. Additionally, we have three domain-specific groups: **Ecommerce**, **Academic**, and **Others**. The **Others** group is primarily composed of biological networks, combined with other small domains that have fewer datasets. See Appendix A.2 for more information of our experimental datasets.

Experimental Settings. We follow previous works (He et al., 2020; Kipf and Welling, 2017) for dataset splitting and metrics. Our AnyGraph model and the graph foundation models are evaluated on a cross-graph zero-shot prediction task. For baselines that cannot handle cross-dataset transfer, we evaluate their few-shot performance. Details of the evaluation protocols are provided in Appendix A.3. The **Hyperparameter Settings** of AnyGraph are provided in Appendix A.4. The compared **Baseline Methods** are introduced in Appendix A.5.

4.2 AnyGraph’s Zero-Shot Prediction (RQ1)

We evaluated AnyGraph’s zero-shot capabilities across 38 diverse graph datasets using two model versions trained separately on Link1 and Link2 datasets. Each model made predictions on the other’s dataset, despite different feature spaces and data sources. Results are shown in Tables 1 and 2. We have the following observations:

i) Superior Generalizability across Diverse Datasets. • **Superior Prediction Accuracy.** Compared to the few-shot capabilities of existing GNN models, pre-training techniques, and foundation models, AnyGraph demonstrates exceptional zero-shot prediction accuracy across various domains. This superior performance spans both link prediction and node classification tasks. • **Effectively Handling Heterogeneity.** The enhanced generalizability can be attributed to the effective handling of structure-level and feature-level data heterogeneity through unified structure and feature representations in the expert models. This approach enables AnyGraph to develop comprehensive modeling functions that are universally applicable across different graph data scenarios. • **Comprehensive Training.** Additionally, the extensive training regimen, which incorporates a variety of large-scale

Table 1: Comparing AnyGraph (in zero-shot setting) with baseline models (with 5% and 10% training data) on link prediction task (in terms of Recall@20, NDCG@20), and node classification task (in terms of Accuracy, Macro F1).

Data	GIN				GAT				GPF				GraphPrompt				GraphCL				AnyGraph	
	Train 5%		Train 10%		Train 5%		Train 10%		Tune 5%		Tune 10%		Tune 5%		Tune 10%		Tune 5%		Tune 10%		0-shot	
Metric	R	N	R	N	R	N	R	N	R	N	R	N	R	N	R	N	R	N	R	N	R	N
Link1	6.46	3.06	11.80	5.45	13.52	6.65	13.45	6.78	6.04	2.92	6.80	3.27	4.33	2.24	5.42	3.11	17.23	9.00	20.55	10.76	23.94	12.68
Link2	6.72	4.50	21.62	13.41	9.83	5.91	15.30	8.84	7.44	4.25	16.58	9.84	6.06	3.36	6.10	3.62	29.18	17.62	31.42	19.91	46.42	27.21
Ecom.	3.36	2.58	13.41	8.06	3.79	2.94	9.64	5.78	7.25	3.84	18.72	10.94	4.90	2.59	6.06	3.36	22.13	13.19	26.05	14.59	26.92	15.05
Acad.	10.82	4.70	20.61	9.04	14.95	6.29	11.17	4.67	13.22	5.80	14.83	6.41	6.73	3.05	7.72	3.40	24.86	12.50	28.69	14.31	32.74	15.31
Othrs.	6.92	4.46	18.43	11.85	16.34	9.22	16.17	20.88	2.40	2.12	4.51	3.44	2.93	2.36	3.42	2.72	24.54	14.93	24.62	15.90	46.83	28.97
Metric	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
Node	20.79	19.46	36.04	30.60	53.76	40.14	54.83	41.61	12.77	11.45	16.29	16.00	18.01	20.59	23.15	22.89	43.70	33.72	48.75	36.15	64.31	43.24

Table 2: Comparing AnyGraph to existing graph foundation models in zero-shot prediction capabilities.

Method	GraphGPT				OpenGraph	
	Pubmed		Cora		Ecom. w/o GR	
Data	Acc	MacF1	Acc	MacF1	Recall	NDCG
Baseline	0.1813	0.1272	0.7011	0.6491	0.1444	0.1099
Ours-F	0.5852	0.5325	0.7134	0.6003	0.2281	0.1600
Ours	0.6088	0.5492	0.7809	0.7591	0.2382	0.1552

datasets, equips AnyGraph with a deep and broad expertise in graph learning.

ii) Limitation of existing pre-training GNNs. • Challenges of Cross-Domain Transfer. Existing pre-training and tuning methods, like GPF, GraphPrompt, and GraphCL, employ self-supervised learning and are pre-trained on half the datasets, then fine-tuned on the remaining datasets using few-shot data. However, this pre-training often fails to yield significant improvements due to substantial distribution disparities across data domains. For instance, datasets may exhibit vastly different link densities or utilize distinct node features, which significantly challenges the transfer of useful knowledge from divergent pre-training datasets during fine-tuning and prediction. **• AnyGraph’s Robust Adaptability.** To address this challenge, AnyGraph incorporates multiple graph experts tailored to various sub-domains of graph data. This MoE architecture effectively manages datasets from distinctly different domains, such as e-commerce user behaviors, academic networks, and road networks, demonstrating its robust adaptability.

4.3 Scaling Law of AnyGraph (RQ2)

We tested the performance scaling trend using 18 AnyGraph variants with different model sizes and training data volumes (see Appendix A.6 for their configurations). Figure 3 shows overall, domain-specific, zero- and full-shot performance results.

i) AnyGraph Generalizability Follows Scaling Law. As model size and training data increase, AnyGraph’s full-shot performance reaches satura-

tion while zero-shot accuracy continues improving, supporting the scaling law of graph foundation models. Two factors explain this: **• Task Difficulty** - full-shot performance saturates because evaluation tasks may be insufficiently challenging, as in-domain generalization is simpler, suggesting the need to test larger models on more complex tasks; and **• MoE Architecture** - the Mixture of Experts architecture enables AnyGraph to better handle diverse knowledge, particularly beneficial for zero-shot scenarios with distribution differences.

ii) Emergent Abilities of AnyGraph. The overall zero-shot performance curve illustrates that as the model size increases, the performance sometimes experiences periodic stagnation. With further increments in parameters, AnyGraph’s performance undergoes a sudden significant improvement. This phenomenon indicates the emergent abilities of AnyGraph, demonstrating the effectiveness of scaling up in enhancing its generalization capabilities. **iii) Insufficient training data may bring bias.** In the initial stages of increasing the training data, the introduction of new datasets might negatively impact performance due to their differences from the test graphs. However, this issue can be mitigated by further expanding the training data. By providing the model with a more comprehensive set of training samples, it helps prevent overfitting and reduces bias stemming from dataset disparities.

4.4 Ablation Study (RQ3)

We evaluated AnyGraph’s sub-modules through ablation studies on both cross-domain and domain-specific datasets, with results shown in Figure 4. We make the following observations:

- MoE Significantly Enhances Zero-Shot Performance.** The **-MoE** variant with a single expert model shows good full-shot but poor zero-shot performance, highlighting how multiple experts enhance AnyGraph’s generalization by managing domain disparities through separated models.

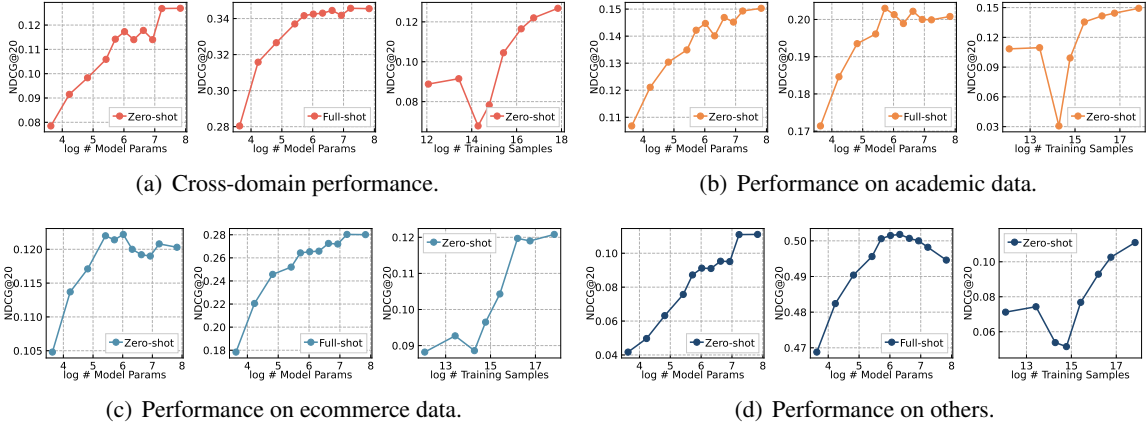
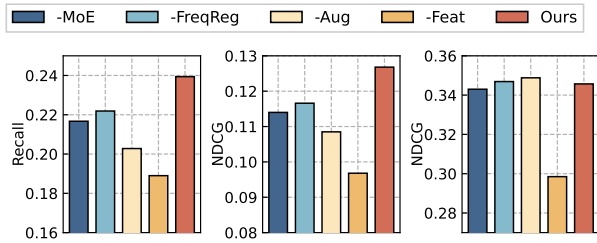
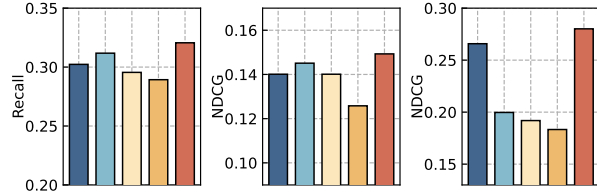


Figure 3: Zero-shot and full-shot performance *w.r.t.* the amount of parameters and training samples.



(a) Zero-shot (left & middle) and full-shot (right) performance tested across multiple domains.



(b) Zero-shot (left & middle) and full-shot (right) performance tested on Academic datasets only.

Figure 4: Impact of AnyGraph’s sub-modules on zero-shot and full-shot prediction capabilities.

• **Feature Modeling is Crucial in AnyGraph.**

The -Feat variant shows the largest performance drop in both zero-shot and full-shot tests, demonstrating the importance and effectiveness of AnyGraph’s unified feature representation for handling heterogeneous graph data across domains.

- **Frequency Regularization and Graph Augmentation.** The -FreqReg and -Aug variants, created by removing routing frequency adjustment and data augmentations respectively, show reduced performance, confirming their importance for model robustness across datasets.

4.5 Investigation on Expert Routing (RQ4)

This section delves into the expert routing mechanism of AnyGraph. Figure 5 displays the com-

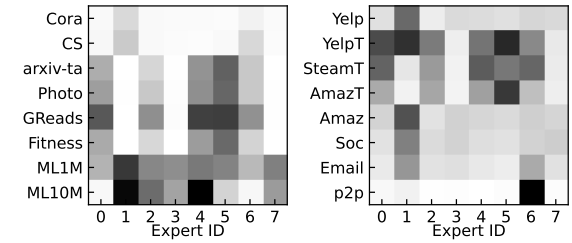


Figure 5: Matching scores between datasets and experts, given by the routing mechanism.

petence scores of various expert models for the input datasets, as determined by AnyGraph’s routing algorithm based on self-supervised loss. The figure illustrates that datasets sharing common characteristics—such as source of collection or feature construction method—are often routed to the same expert models by AnyGraph. For instance, datasets like arxiv-ta, Photo, GoodReads, and Fitness, which utilize a common text-embedding-based feature space, are assigned to highly similar experts. Additionally, ML1M and ML10M, both sourced from the movie-rating platform MovieLens, are predominantly associated with expert 1. It is also notable that this routing pattern extends to zero-shot datasets, as shown on the right part of Figure 5. Here, YelpT, SteamT, and AmazonT, which share the same feature space, are assigned to very similar models. This outcome highlights the effectiveness and the explainability of AnyGraph’s routing mechanism.

4.6 Efficiency Study (RQ5)

Tuning Curve Comparison. We compared AnyGraph’s fine-tuning efficiency against GraphCL and from-scratch GCN training. Figure 6 shows AnyGraph quickly reaches high performance on new datasets, sometimes significantly surpassing

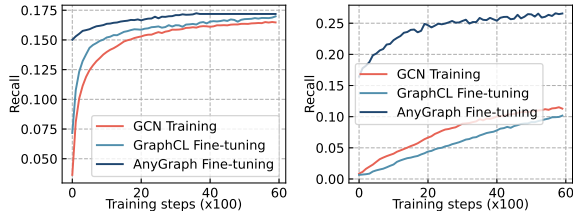


Figure 6: Performance v.s. training/tuning steps, on Citation-2019 (left) and PPA (right) datasets.

Table 3: Training time for each 100 steps.

Dataset	CS	ML1M	Yelp	Email	Cite19	roadNet	PPA
GCN	1.5s	4.2s	6.0s	2.5s	19.2s	27.8s	101.1s
GraphCL	1.1s	4.9s	9.4s	2.8s	43.1s	57.1s	130.8s
Ours	1.5s	3.5s	6.1s	3.0s	31.6s	37.3s	41.1s

the others (*e.g.*, on PPA dataset). This advantage stems from AnyGraph’s strong cross-domain generalization providing a high starting point, and its efficient MoE architecture requiring only one MLP network for effective modeling and tuning.

In addition, it is observed that pre-training GraphCL does not always benefit its fine-tuning on new datasets, as evidenced by GraphCL’s underperformance relative to GCN in Figure 6 (right). This is due to the large distribution gap between the pre-training data Link2 and the test data PPA.

Training Time Comparison. To evaluate the model efficiency, we compared the training times of the three models. As indicated in Table 3, AnyGraph, despite having significantly more parameters, has training times that are comparable to, or even less than, the other two models. This underscores the efficiency of our model design.

Specifically, AnyGraph avoids the cumbersome full-graph propagation. Instead, it utilizes structure-aware embeddings derived through a non-trainable pre-processing method. This significantly reduces both the time and memory requirements. Furthermore, the MoE architecture enables AnyGraph to use only $1/K$ of the computational resources for most prediction and optimization processes, thereby greatly reducing computational costs.

5 Related Works

Text-aware Graph Neural Models (GNNs). Graph learning has gained significant interest across fields like user behavior modeling and biology applications (Chang et al., 2021; Hao et al., 2020). GNNs learn node representations through iterative message passing, capturing both node-specific information and topological structures. Notable techniques include GCNs (Jin et al., 2021),

GATs (Brody et al., 2022), GIN (Xu et al., 2018), and Graph Transformers (Hu et al., 2020). Despite advancements, these methods rely on high-quality training data and struggle with generalization.

Self-Supervised Graph Learning. To address GNN generalization challenges, considerable efforts (Xie et al., 2022) have focused on self-supervised learning objectives to capture invariant graph features. GraphCL (You et al., 2020) introduced contrastive pre-training robust to perturbations. JOAO (You et al., 2021) and GCA (Zhu et al., 2021) developed adaptive augmentation strategies, while GPF (Fang et al., 2023) and Graph-Prompt (Liu et al., 2023) focused on fast adaptation to downstream tasks. However, generalizability remains confined to similar structural and feature patterns, overlooking cross-domain challenges.

Large-scale Graph Pre-training. Recent advances have explored pre-training large-scale graph models across multiple datasets, inspired by LLMs’ generalizability. OFA (Liu et al., 2024) and ZeroG (Li et al., 2024) utilize text embeddings to standardize feature spaces across datasets. InstructGLM (Ye et al., 2024), GraphGPT (Tang et al., 2024a), and LLaGA (Chen et al., 2024) synchronize graph representations with LLM hidden spaces, while HiGPT (Tang et al., 2024b) extends capabilities to heterogeneous graphs. However, these methods require substantial text features, limiting use to text-abundant environments, and typically train within specific domains without addressing cross-domain variance.

6 Conclusion

This work presents AnyGraph, an effective and efficient graph foundation model designed to address the multifaceted challenges of structure and feature heterogeneity across diverse graph datasets. AnyGraph’s innovative Mixture-of-Experts (MoE) architecture, coupled with its dynamic expert routing mechanism, positions it at the state-of-the-art of cross-domain generalization capabilities. Extensive experiments on 38 varied graph datasets have not only underscored AnyGraph’s superior zero-shot learning performance but also its robustness to distribution shifts and its adherence to scaling laws, thereby enhancing its predictive accuracy with increased model size and data volume. The model’s efficiency in training and inference, validated through comparison with existing methods, further cements its practical applicability.

7 Limitations

We identify the following limitations of AnyGraph:

Static Graph Focus. The current AnyGraph framework operates on static graph snapshots and does not explicitly model temporal dynamics. For applications involving evolving graph structures such as social networks or dynamic recommendation systems, incorporating temporal modeling capabilities could enhance performance on time-sensitive prediction tasks.

Limited Heterogeneous Relation Modeling. The current architecture treats all edges uniformly within each expert model. For graphs with multiple edge types or complex heterogeneous relations, the model may not fully capture the nuanced patterns that relation-specific modeling could provide.

References

- Shaked Brody, Uri Alon, and Eran Yahav. 2022. How attentive are graph attention networks? In *International Conference on Learning Representations*.
- Xuheng Cai, Chao Huang, Lianghao Xia, and Xubin Ren. 2023. Lightgcl: Simple yet effective graph contrastive learning for recommendation. *International Conference on Learning Representations (ICLR)*.
- Jianxin Chang, Chen Gao, Yu Zheng, Yiqun Hui, Yanan Niu, Yang Song, Depeng Jin, and Yong Li. 2021. Sequential recommendation with graph neural networks. In *The International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 378–387.
- Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. 2022. Graph unlearning. In *IEEE Symposium on Security & Privacy (SIGSAC)*, pages 499–513.
- Runjin Chen, Tong Zhao, Ajay Jaiswal, Neil Shah, and Zhangyang Wang. 2024. Llaga: Large language and graph assistant. In *International Conference on Machine Learning (ICML)*.
- Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. 2023. Reproducible scaling laws for contrastive language-image learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2829.
- Taoran Fang, Yunchao Zhang, Yang Yang, Chunping Wang, and Lei Chen. 2023. Universal prompt tuning for graph neural networks. *Conference on Neural Information Processing Systems (NeurIPS)*.
- Matthias Fey, Weihua Hu, Kexin Huang, Jan Eric Lenssen, Rishabh Ranjan, Joshua Robinson, Rex Ying, Jiaxuan You, and Jure Leskovec. 2024. Position: Relational deep learning-graph representation learning on relational databases. In *International Conference on Machine Learning (ICML)*.
- Zhongkai Hao, Chengqiang Lu, Zhenya Huang, Hao Wang, Zheyuan Hu, Qi Liu, Enhong Chen, and Cheekong Lee. 2020. Asgn: An active semi-supervised graph neural network for molecular property prediction. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 731–752.
- Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *The International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 639–648.
- Ziniu Hu, Yuxiao Dong, Kuansan Wang, and Yizhou Sun. 2020. Heterogeneous graph transformer. In *ACM The Web Conference (WWW)*, pages 2704–2710.
- Di Jin, Zhizhi Yu, Cuiying Huo, Rui Wang, Xiao Wang, Dongxiao He, and Jiawei Han. 2021. Universal graph convolutional networks. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 10654–10664.
- Wei Jin, Xiaorui Liu, Xiangyu Zhao, Yao Ma, Neil Shah, and Jiliang Tang. 2022. Automated self-supervised learning for graphs. In *International Conference on Learning Representations (ICLR)*.
- Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. 2020. Graph structure learning for robust graph neural networks. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 66–74.
- Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*.
- Guohao Li, Matthias Müller, Bernard Ghanem, and Vladlen Koltun. 2021. Training graph neural networks with 1000 layers. In *International Conference on Machine Learning (ICML)*, pages 6437–6449.
- Yuhan Li, Peisong Wang, Zhixun Li, Jeffrey Xu Yu, and Jia Li. 2024. Zerog: Investigating cross-dataset zero-shot transferability in graphs. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- Hao Liu, Jiarui Feng, Lecheng Kong, Ningyue Liang, Dacheng Tao, Yixin Chen, and Muhan Zhang. 2024. One for all: Towards training one graph model for all classification tasks. In *International Conference on Learning Representations (ICLR)*.
- Yixin Liu, Ming Jin, Shirui Pan, Chuan Zhou, Yu Zheng, Feng Xia, and S Yu Philip. 2022. Graph self-supervised learning: A survey. *Transactions on Knowledge and Data Engineering (TKDE)*, pages 5879–5900.

- Zemin Liu, Xingtong Yu, Yuan Fang, and Xinming Zhang. 2023. Graphprompt: Unifying pre-training and downstream tasks for graph neural networks. In *ACM The Web Conference (WWW)*, pages 417–428.
- Haitao Mao, Zhikai Chen, Wei Jin, Haoyu Han, Yao Ma, Tong Zhao, Neil Shah, and Jiliang Tang. 2024. Demystifying structural disparity in graph neural networks: Can one size fit all? In *Conference on Neural Information Processing Systems (NeurIPS)*, volume 36.
- Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. 2024. Scaling data-constrained language models. In *Conference on Neural Information Processing Systems (NeurIPS)*, volume 36.
- Mingchen Sun, Kaixiong Zhou, Xin He, Ying Wang, and Xin Wang. 2022. Gppt: Graph pre-training and prompt tuning to generalize graph neural networks. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 1717–1727.
- Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Lixin Su, Suqi Cheng, Dawei Yin, and Chao Huang. 2024a. Graphgpt: Graph instruction tuning for large language models. In *The International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 491–500.
- Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Long Xia, Dawei Yin, and Chao Huang. 2024b. Higtpt: Heterogeneous graph language model. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. In *International Conference on Learning Representations (ICLR)*.
- Junke Wang, Dongdong Chen, Zuxuan Wu, Chong Luo, Luwei Zhou, Yucheng Zhao, Yujia Xie, Ce Liu, Yungang Jiang, and Lu Yuan. 2022. Omnivl: One foundation model for image-language and video-language tasks. *Conference on Neural Information Processing Systems (NeurIPS)*, 35:5696–5710.
- Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. 2023. Internimage: Exploring large-scale vision foundation models with deformable convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14408–14419.
- Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, pages 6861–6871.
- Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised graph learning for recommendation. In *The International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 726–735.
- Lianghao Xia, Ben Kao, and Chao Huang. 2024. Opengraph: Towards open graph foundation models. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Teng Xiao, Zhengyu Chen, Donglin Wang, and Suhang Wang. 2021. Learning how to propagate messages in graph neural networks. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 1894–1903.
- Yaochen Xie, Zhao Xu, Jingtun Zhang, Zhengyang Wang, and Shuiwang Ji. 2022. Self-supervised learning of graph neural networks: A unified review. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 45(2):2412–2429.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2018. How powerful are graph neural networks? In *International Conference on Learning Representations (ICLR)*.
- Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. 2024. Language is all a graph needs. In *Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 1955–1973.
- Yuning You, Tianlong Chen, Yang Shen, and Zhangyang Wang. 2021. Graph contrastive learning automated. In *International Conference on Machine Learning (ICML)*, pages 12121–12132. PMLR.
- Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. 2020. Graph contrastive learning with augmentations. *Conference on Neural Information Processing Systems (NeurIPS)*, 33:5812–5823.
- Qiannan Zhang, Shichao Pei, Qiang Yang, Chuxu Zhang, Nitesh V Chawla, and Xiangliang Zhang. 2023. Cross-domain few-shot graph classification with a reinforced task coordinator. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 4893–4901.
- Haihong Zhao, Aochuan Chen, Xiangguo Sun, Hong Cheng, and Jia Li. 2024. All in one and one for all: A simple yet effective method towards cross-domain graph pretraining. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*.
- Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. 2021. Graph contrastive learning with adaptive augmentation. In *ACM The Web Conference (WWW)*, pages 2069–2080.

A Appendix

A.1 Details of Model Training

Feature and Structure Augmentation. To further enrich the training data, the training of AnyGraph undergoes periodic reprocessing of, **firstly**, the initial graph embeddings \mathbf{E}_1 , and **secondly**, the graph routing results. We demonstrate that such reprocessing augments the features and structures of the original graph data, thereby training AnyGraph using more diversified input data.

For the initial graph embeddings, we periodically reconduct the SVD and simplified GCN processes after a certain number of training steps. This helps generate different embedding spaces for the same data, thereby greatly improving the generalizability of AnyGraph regarding representation heterogeneity (Li et al., 2024). To prevent this process from consuming excessive computational time, we propose adopting different augmentation frequencies adaptive to the size of different datasets. Specifically, each dataset undergoes this representation augmentation after $|\mathcal{E}|/(10B)$ training steps.

For the graph routing results, we also periodically recalculate the recalibrated competence scores. Specifically, the positive sample pairs (v_{c_s}, v_{p_s}) for $s = 1, \dots, S$, as well as the negative samples v_{n_s} , are randomly sampled. This essentially performs structure augmentation by using a random subset to evaluate the performance of graph experts on the input graph, thereby enhancing the model’s robustness against structural noise.

Complexity Analysis. The training and inference process of our AnyGraph involve only a single expert model, yielding a time complexity of $\mathcal{O}(B \times d^2 \times L')$ per batch. Preprocessing of initial embeddings and expert routing does not add to this batch-wise complexity, making AnyGraph significantly more efficient than typical graph foundation models that use complex GNN models such graph transformers. Additionally, expert routing requires $\mathcal{O}(\sum_{\mathcal{G}_s} |\mathcal{E}_s| \times d \times K + \sum_{\mathcal{G}_s} |\mathcal{V}_s| \times d^2 \times L' \times K)$ computations, with the latter term generally larger and comparable to a simple GCN network. Thus, AnyGraph demonstrates greater efficiency in training and inference compared to existing methods, with the additional routing complexity akin to that of simple GNNs.

A.2 Experimental Datasets

We utilize a total of 38 graph datasets across various domains. The entire experimental data contains

14,437,372 nodes, and 199,265,688 edges. The dataset specifics are detailed below:

E-commerce Datasets. This category includes 15 datasets from various e-commerce contexts such as user rating platforms and online retail services. These datasets vary in terms of the presence and type of node features. For instance, datasets such as Amazon-book, Yelp2018, Gowalla, Yelp-text, Amazon-text, Steam-text, Goodreads, Amazon-Fitness, Amazon-Photo, Movielens-1M, Movielens-10M, Products-home, Products-tech, Home-node, Tech-node are included. Notably, Amazon-text, Steam-text, and Yelp-text utilize the same method for feature generation, while Fitness, Photo, and Goodreads employ a different consistent method.

Academic Network Datasets. We use 13 datasets focused on academic networks, which include citation and collaboration relations among scholars and papers. These datasets represent various research fields and employ diverse feature generation methods, such as NLP embeddings, bag-of-words, and different versions of large language models. The specific datasets are Cora, Pubmed, Arxiv, Cora-link, Pubmed-link, Citeseer, CS, Arxiv-link, Arxiv-t (with features derived using an alternative method), Cite-2019, Cite-20Cent, OGB-Collab.

Biological Information Networks. Our experimental data includes 6 datasets related to biological entities like proteins, drugs, and diseases. This category features networks such as OGB-DDI, OGB-PPA, which record drug-drug and protein-protein relations, respectively, and four other protein relation networks for different species, denoted as Proteins-0, Proteins-1, Proteins-2, Proteins-3.

Other Datasets. In addition to the categories mentioned above, we include 5 datasets from various other fields: an email network Enron, a website network Stanford, a road network dataset Road-PA, a P2P web network dataset Gnutella, and a trust network dataset Epinions.

Dataset Groups. For convenience of performance evaluation, we split the many datasets using different grouping methods. Firstly, two big data groups Link1 and Link2 are made using all the link prediction datasets. Notably, datasets from the same source of collection, such as ML-1M and ML-10M, or uses the same method to generate features, such as Fitness, and Photo, are put into the same group, to avoid information leakage when evaluating zero-shot performance on the other group. Apart from these two datasets, we also conduct

evaluations on domain-specific groups, including E-commerce, Academic, and Others. Specifically, these data groups contain the following datasets:

- **Link1.** This is the first link prediction dataset group, including the following datasets: Products-tech, Yelp2018, Yelp-text, Products-home, Steam-text, Amazon-text, Amazon-book, Cite-2019, Cite-20Cent, Pubmed-link, Citeseer, OGB-PPA, Gnutella, Epinions, Enron.
- **Link2.** This is the other link prediction dataset group, with minimal connection to Link1. It includes the following datasets: Photo, Goodreads, Fitness, Movielens-1M, Movielens10M, Gowalla, Arxiv, Arxiv-t, Cora, CS, OGB-Collab, Proteins-0, Proteins-1, Proteins-2, Proteins-3, OGB-DDI, Stanford, Road-PA.
- **Ecommerce and Academic.** They contain all domain-specific datasets mentioned above.
- **Others:** This group contains all the biological datasets mentioned above, and datasets from other minor domains, including email network data Enron, website network data Stanford, road network data RroadNet-PA, P2P network data Gnutella, and trust network data Epinions.

A.3 Evaluation Protocols

All datasets used in this study are sourced from previous research as referenced (Tang et al., 2024a; Li et al., 2024). We adhere to the original data splits from these sources to delineate our training and testing sets. Given that many baseline methods are not equipped to manage zero-shot prediction across datasets, we instead assess their few-shot capabilities. This allows for a comparative analysis against the zero-shot performance of AnyGraph. We employ specific evaluation settings tailored to each method, detailed as follows:

- **Zero-shot Setting for AnyGraph, GraphGPT, and OpenGraph.** In our study, AnyGraph and two comparative graph foundation models, GraphGPT and OpenGraph, undergo evaluations for zero-shot prediction capabilities. We pre-train two instances of AnyGraph using Link1 and Link2 datasets. The model pre-trained on Link1 is then tested for zero-shot performance on the Link2 group datasets, and vice versa. Results labeled as "zero-shot" for AnyGraph are derived using this cross-evaluation method. Conversely,

results marked as "full-shot" pertain to supervised learning outcomes, where, for example, the model trained on Link1 is tested on the test sets of Link1 group datasets. For GraphGPT and OpenGraph, we utilize the models as released in their respective original studies, which were pre-trained on specified datasets.

- **Zero-shot Node Classification for AnyGraph.** Inspired by prior research (Sun et al., 2022), we approach zero-shot node classification by representing node classes as distinct nodes. We then connect existing nodes that have training labels directly to these new class nodes. This technique eliminates the need for learning specific parameters for each class within the zero-shot learning framework, streamlining the process. We have integrated this innovative approach into baseline methods as well, enhancing their capability to handle unseen node labels effectively.
- **Few-shot Training for GIN and GAT.** The GIN and GAT models, employed as end-to-end training baselines, undergo training from scratch on few-shot subsets of the evaluation datasets. This approach is necessary because these models are not well-suited for cross-dataset transfer, particularly when dealing with datasets that have varying feature dimensionalities.
- **Pre-training and Few-shot Tuning for GraphCL, GPF and GraphPrompt.** These category of baseliens methods follow the pre-training-and-fine-tuning mode. In our evaluations, they are firstly pre-trained using the same pre-training datasets as our AnyGraph. Then, they experience an additional fine-tuning process using the few-shot subsets of the evaluation datasets.

Evaluation Metrics. For link prediction, we follow previous works (He et al., 2020) and utilize Recall@20 and NDCG@20 as the evaluation metrics. Note that we typically use the summary results of the evaluation results across multiple datasets. Results for different datasets are averaged according to their number of test samples. For the node classification task, we employ the widely-used Accuracy and Macro-F1 score as our metrics (Chen et al., 2022; Tang et al., 2024a).

A.4 Hyperparameter Settings

Optimization. Our model, AnyGraph, is implemented using PyTorch. The optimization process

Table 4: Statistics of the experimental datasets.

Dataset	DDI	Collab	ML1m	ML10m	Amazon-book	PPA	Yelp2018	Gowalla	Cora	Pubmed	Citeseer
# Nodes	4,267	235,868	9,746	80,555	144,242	576,289	69,716	70,839	2,708	19,717	3,327
# Edges	1,334,889	1,285,465	920,193	9,200,050	2,984,108	45,495,642	1,561,406	1,027,370	10,556	88,648	9,104
<i>d</i> Feats	0	128	0	0	0	58	0	0	1433	500	3703
Datasets	Proteins-0	Proteins-1	Proteins-2	Proteins-3	Products-home	Products-tech	Yelp-t	Amazon-t	Steam-t	Goodreads	Fitness
# Nodes	25,449	6,568	18,108	13,015	9,790	47,428	22,101	20,332	28,547	676,084	173,055
# Edges	11,660,646	1,845,960	7,418,688	3,962,930	131,843	2,077,241	277,535	200,860	525,922	8,582,306	1,773,500
<i>d</i> Feats	0	0	0	0	100	100	1536	1536	1536	768	768
Datasets	Epinions	Enron	Stanford	Road-PA	Gnutella	Cite-2019	Cite-20Cent	Arxiv	Arxiv-t	Photo	CS
# Nodes	75,879	36,692	281,903	1,088,092	8,717	765,658	1,016,241	169,343	169,343	48,362	18,333
# Edges	508,837	183,831	2,312,497	1,541,898	31,525	1,917,381	5,565,798	1,166,243	1,166,243	500,939	163,788
<i>d</i> Feats	0	0	0	0	128	128	128	128	768	768	6805

employs the Adam optimizer with a learning rate of 1×10^{-4} and a training batch size of 4096. We use cross-entropy loss with a sampled negative set (Wu et al., 2021). The learnable parameters of AnyGraph are initialized using the Xavier uniform initializer. **Network Configurations.** The standard configuration of our AnyGraph includes 512 hidden units and 8 graph expert models. Each expert model comprises 8 fully-connected layers. These layers utilize a ReLU activation function and incorporate a dropout layer with a dropout probability of 0.1. **Algorithm Hyperparameters.** The frequency regularization of our routing mechanism is set with an adjustment range of $\rho = 0.2$. The SVD decomposition is performed using 2 iterations. For structural and feature augmentation, each dataset is reprojected after using 1/10 of its samples for optimization. A minimum of 100 training steps should be executed for each dataset before its initial representations are reprojected. The reassignment of experts occurs after all training datasets have undergone one cycle of re-projection.

The baseline methods are evaluated using their original code or released model. We closely follow the original code to adapt to our experiments. Grid search is conducted to search for the best hyperparameter settings for each baseline method.

A.5 Baseline Methods

This section provides detailed descriptions of the baselines used in our analysis. We employ seven different baseline models across four categories.

Training-from-scratch Graph Neural Networks.

- **GAT** (Veličković et al., 2018). Graph Attention Networks (GAT) leverage an attention mechanism to dynamically weight node-to-node connections, enhancing the model’s ability to adaptively propagate and aggregate information across the graph.
- **GIN** (Xu et al., 2018). The Graph Isomorphism

Network (GIN) significantly boosts the expressive power of Graph Neural Networks by introducing a unique graph encoding technique aimed at effectively distinguishing between non-isomorphic graphs.

Graph Pre-training Models.

- **GraphCL** (Zhu et al., 2021). It enhances the pre-training of graph models via self-discriminative contrastive learning, which is applied to learned node embeddings. The method employs various graph augmentation techniques such as node dropping, edge permutation, random walks, and feature masking to improve robustness.

Graph Prompt Tuning Methods.

- **GraphPrompt** (Liu et al., 2023). It proposes a unified approach that integrates pre-training and prompt tuning for graph models. It features a learnable prompt layer designed to automatically extract crucial information from the pre-trained model to enhance downstream performance.
- **GPF** (Fang et al., 2023). The Graph Prompt Framework (GPF) is a versatile graph prompt tuning framework compatible with various graph pre-training methods. It offers two variants of a learnable graph prompt layer, tailored to different application needs.

Graph Foundation Models.

- **GraphGPT** (Tang et al., 2024a). This approach proposes representation alignment and instruction tuning techniques to align graph representation spaces with text encoding spaces, empowering large language models with the capabilities of zero-shot graph encoding and inference.
- **OpenGraph** (Xia et al., 2024). This method introduces a unified graph tokenizer, a scalable

Table 5: Model and training data configurations.

Model Size Configurations											
Model	M1	M2	M3	M4	M5	M6	M7	M8	M9	M10	M11
Hidden Units	64	128	256	512	512	512	512	512	512	512	1024
# FC Layers	1	1	1	1	2	4	8	8	8	8	8
# Experts	1	1	1	1	1	1	1	2	4	8	8

Training Data Configurations			
ID	Datasets	ID	Datasets
D1	Cora, CS	D2	D1 + Photo
D3	D2 + ML1M	D4	D3 + Gowalla
D5	D4 + Arxiv, Arxiv-t	D6	D5 + collab, ddi, Fitness, proteins-spec1, Yelp2018, web-Stanford, proteins-spec3
D7	D6 + proteins-2, roadNet-PA, Fitness	D8	All Link2 datasets

graph transformer to improve the model’s performance and generalization ability. An LLM-enhanced data augmentation mechanism is proposed to address domain-specific data scarcity.

A.6 Details of the Scaling Law Experiment

For the scaling law experiment (RQ2), we elaborate the configurations of the developed instances of AnyGraph. Table 5 summarizes the configurations. Specifically, for AnyGraph with different model sizes, we begin with the smallest model which has 64 hidden units, 1 fully-connected layer, and 1 expert model. The subsequent 3 model instances increases in their hidden dimensionality, from 64 to 128, 256, and 512. Then 3 larger models with more fully-connected layers are utilized, respectively containing 2, 4, and 8 MLP layers. Then we have MoE versions of AnyGraph, with 2, 4, and 8 experts, respectively. The final largest instance has a larger latent dimensionality of 1024.

For the increase of training data, we begin with a subset of Link2 data including Cora and CS. The next version additionally includes Photo. The thir one includes ML1M. The fourth one includes Gowalla. The fifth one additionally include Arxiv and Arxiv-t. The sixth one adds the following datasets: collab, ddi, Yelp2018, Fitness, proteins-spec1, web-Stanford, proteins-spec3. The seventh one is trained with proteins-2, roadNet-PA, and Fitness additionally. And the final one is trained with all datasets from Link2. In this manner, we gradually increase the amount of training data.

A.7 Supplementary Experimental Results

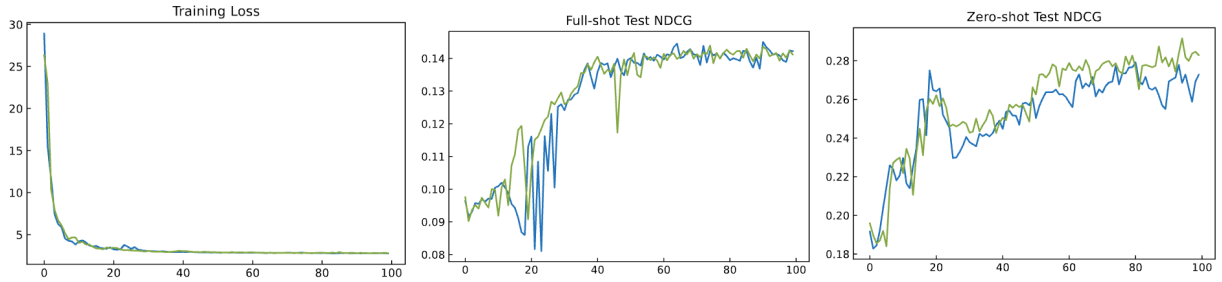
Model Performance Curves. We monitored the training loss and test performance of AnyGraph across each training epoch to understand its training dynamics. This included evaluating AnyGraph’s performance on the test sets of its training datasets

(full-shot performance) as well as its performance on unseen datasets (zero-shot performance), as depicted in Figure 7.

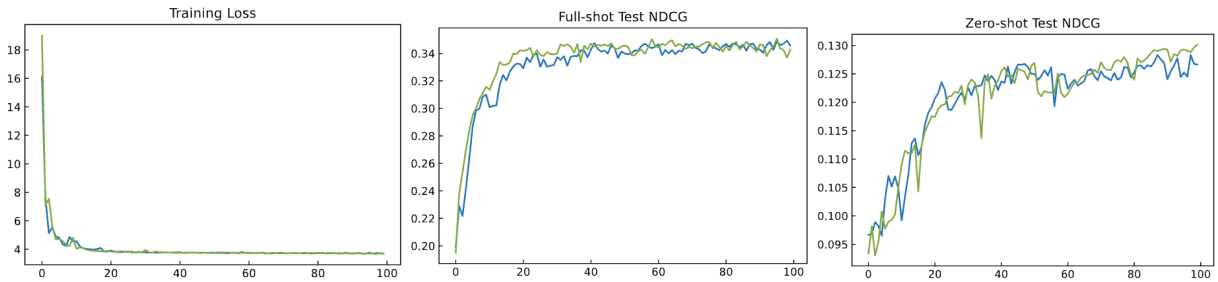
The analysis reveals that training loss and full-shot test performance stop to decrease/increase significantly after approximately 40 epochs. In contrast, zero-shot test performance continues to improve significantly, even up to 100 epochs. This trend underscores a steady enhancement in the model’s generalization abilities, highlighting the potential to further explore and enhance the generalizability of graph models in challenging zero-shot inference tasks.

Performance on Industrial Data. We further assessed the performance of AnyGraph using a real-world dataset from a popular user reading platform, comprising over 1 million user and item nodes. We trained a base graph neural model on historical user behavior data, and evaluated both the base model and AnyGraph using varying amounts of new interaction data to construct the input graph. The results, summarized in Table 6, show that “History” indicates the base model was trained on data from previous days, while “10%”, “20%”, etc. represent the percentages of new data used to construct the input graph. Importantly, the new data was used only as input features, not for tuning, reflecting a real-world scenario where models cannot be promptly fine-tuned on new data. Our key observations are: i) AnyGraph demonstrated superior zero-shot predictive capabilities, outperforming the base model trained on historical data. ii) This underscores the importance of robust zero-shot prediction, as new data may not align with historical patterns in real-world settings.

Recall@20 for Full-shot Performance in Ablation Study. We have expanded our analysis to include full-shot prediction performance, as assessed in our ablation studies. Figure 8 displays the performance of various ablated versions of our AnyGraph alongside the complete model, using Recall@20 as the metric. A notable finding, absent from the original results, is that removing the augmentation actually results in a significant advantage for our AnyGraph in cross-domain evaluations. This phenomenon can be attributed to the fact that data augmentations interfere with the optimization of AnyGraph on the training dataset, thereby impairing the full-shot performance on seen datasets. However, as the zero-shot performance test results indicate, this augmentation technique substantially enhances the generalization capability of AnyGraph. This is



(a) Two instances of AnyGraph independently trained on Link1.



(b) Two instances of AnyGraph independently trained on Link2.

Figure 7: Training loss, test NDCG of full-shot and zero-shot prediction, v.s. the number of training epochs. Two curves in each plot correspond to two independently-trained instances of AnyGraph.

Table 6: Performance on industrial data.

Method	History	10%	20%	30%	40%	50%
Base Method	0.7%	2.0%	5.6%	10.6%	17.3%	19.9%
AnyGraph	6.3%	3.4%	7.5%	14.0%	19.3%	21.7%

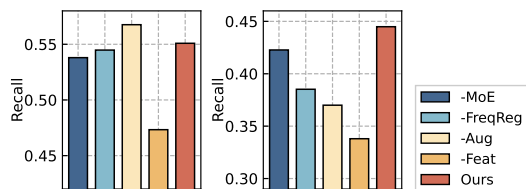


Figure 8: Recall results of ablation study, on cross-domain (left) and academic (right) data.

because the disturbances prevent the model parameters from overfitting to the training data.

A.8 Future Directions

While AnyGraph demonstrates strong generalizabilities in handling unseen graph data and achieving zero-shot generalization, we identify several promising directions for future extensions:

Temporal Graph Dynamics. The current version of AnyGraph focuses on static graph snapshots. For applications involving temporal graphs like social networks and e-commerce interactions, incorporating temporal modeling capabilities could further enhance performance. Future extensions could integrate sequential architectures like temporal attention mechanisms to capture dynamic pat-

terns, enabling AnyGraph to better handle evolving graph structures and temporal dependencies.

Multi-modal Feature Integration. While AnyGraph effectively processes different types of node features through its unified representation approach, there are opportunities to develop more sophisticated multi-modal encoding frameworks. Such extensions could better preserve and integrate complementary signals when nodes have multiple feature types simultaneously (*e.g.*, text, images, and numerical attributes), potentially leading to even richer graph representations.

Relation-aware Modeling. The current AnyGraph architecture successfully handles graphs with varying edge distributions through its MoE framework. This could be extended with relation-aware components to explicitly model multiple edge types (*e.g.*, different interaction types in e-commerce networks). Adding relation-specific attention mechanisms or specialized expert models could help capture more nuanced patterns in graphs with heterogeneous relations.

These future directions highlight the significant potential to build upon AnyGraph’s strong foundation in graph representation learning and zero-shot generalization. Addressing these aspects could further enhance the model’s versatility across an even broader range of real-world applications.