TEST-TIME GRAPH REBIRTH: SERVING GNN GENERALIZATION UNDER DISTRIBUTION SHIFTS

Anonymous authors

Paper under double-blind review

Abstract

011 Distribution shifts between training and test graphs typically lead to the decreased 012 performance of well-trained graph neural network (GNN) models, negatively affecting their ability to generalize in real-world applications. Although there have 013 been advances in addressing graph distribution shifts through various model archi-014 tectures and training strategies, implementing existing solutions in practical GNN 015 deployment and serving at test time can be challenging, as they often necessitate 016 significant modifications or retraining of the GNNs. To address such challenges, 017 in this work, we propose a novel method, *i.e.*, Test-Time Graph **REB**irth, dubbed 018 **TT-GREB**, to effectively generalize the well-trained GNN models to the test-time 019 graphs under distribution shifts by directly manipulating the test graph data. Concretely, we develop an overall framework designed by two principles, corresponding 021 to two sub-modules: (1) prototype extractor for re-extracting the environmentinvariant features of the test-time graph; and (2) environment refiner for re-fining the environment-varying features to explore the potential shifts. Furthermore, we 023 propose a dual test-time graph contrastive learning objective with an effective iterative optimization strategy to obtain optimal prototype components and envi-025 ronmental components of the test graph. By reassembling these two components, 026 we obtain a newly reborn test graph, which is better suited for generalization on 027 the well-trained GNN model with shifts in graph distribution. Extensive experi-028 ments on real-world graphs under diverse test-time distribution shifts verify the 029 effectiveness of the proposed method, showcasing its superior ability to manipulate test-time graphs for better GNN generalization ability.

031 032

033 034

004

010

1 INTRODUCTION

Recent advances in graph neural networks (GNNs) have achieved great success with promising 035 learning abilities for various graph structural data in numerous real-world applications (Zhang et al., 2022; Zheng et al., 2022a;b; 2023c;a; Jin et al., 2022; Zheng et al., 2022c). Well-designed GNN 037 models are ultimately intended for practical deployment and serving on various graph learning tasks (Zheng et al., 2023b; Liu et al., 2023b; Yu et al., 2023). However, these expertly trained GNNs generally experience significant performance degradation due to the graph distribution shift issue 040 between the training and the test graphs (Wu et al., 2022b; Liu et al., 2023a; Chen et al., 2023b; Yu 041 et al., 2023). The main reason behind such distribution mismatch lies in the underlying environmental 042 variations, with time-related attribute changes, agnostic corruptions, and inconsistent graph data 043 collection procedures (Sui et al., 2023; Chen et al., 2023b; Jin et al., 2023). These factors would lead 044 to considerable differences in node contexts, graph structures, and the overall scale of graphs during 045 the test-time stage.

To overcome the model generalization challenge caused by graph distribution shifts, there is a growing focus on research into learning with distribution shifts on graphs, from the *model-centric* perspective (Wu et al., 2022b; You et al., 2023; Chen et al., 2023c; Wu et al., 2020). Typically, these existing methods incorporate invariant representation learning into GNN development through customizing model architectures and training strategies (Xu et al., 2019; Zhu et al., 2021; Liu et al., 2022; Wu et al., 2022a). However, in real-world GNN deployment and serving, it may not be always practical to re-design GNN model architectures or re-train well-trained GNNs are continuously in service online, as accessing and modifying their parameters becomes more difficult. Given 063

064

065 066 067

068

069

071

073

098

099

102

103



Figure 1: Comparison between the (a) model-centric graph learning methods v.s. our proposed (b) data-centric test-time graph rebirth method under graph distribution shifts.

all these circumstances including *complex training-test graph distribution shifts* and *inaccessible online-serving GNN model fine-tuning*, an intriguing problem emerges:

Question: Is it possible to give the rebirth of the test-time graph data from the data-centric perspective, to serve various well-trained GNN models for better generalization under distribution shifts?

In this work, we provide a feasible graph data-centric solution to answer this question by enhancing the test-time graph data quality, without accessing the training graph data or modifying the welltrained GNNs, as shown in Fig. 1. Specifically, the original test graph would undergo a rebirth process to emerge as a new test graph, which is then fed into the well-trained GNN model for direct inference without retraining or fine-tuning. In the test-time graph rebirth process, we identify two essential components that decomposed a graph under distribution shifts: (a) *the environment-invariant component* retains consistent informative features (such as node class label prototypes) across various training and test graph distributions; (b) *the environment-varying component* dictates the extent of shifts in graph data distribution between training and test.

083 In light of this, we propose a novel test-time graph rebirth method for serving good GNN generalization under graph distribution shifts with two principles: (1) re-extracting the environment-invariant 084 features of the test-time graph for identifying the predictive pattern of node class label prototypes with 085 a **prototype extractor**, and (2) re-fining the environment-varying features to explore the potential shifts in the training-test distribution with a **environment refiner**, leading to the alignment of the 087 training-test environment latent space with improved generalization capabilities for online GNN deployment. Furthermore, due to the unknown test-time graph node labels, we propose a dual test-time graph contrastive learning objective with self-supervision signals, along with an effective 090 iterative optimization strategy to obtain expressive prototype features and environmental features 091 of the test graph. In this way, we preserve prototype features to determine the node class label's 092 predictive patterns, meanwhile, we adjust the environmental features of the original test graph to match those of the training graph. Then, these components are reassembled into a newly reborn test graph, which better suits the well-trained GNNs' generalization under distribution shifts. 094

- In summary, the contributions of this work are as follows:
 - **Graph Data-centric Paradigm.** We introduce a novel graph data-centric paradigm, Test-Time Graph REBirth (TT-GREB), designed to enhance the generalization ability of welltrained GNNs to real-world test graphs experiencing distribution shifts at test time.
 - **Innovative Solution.** We develop a comprehensive framework with two core components: a prototype extractor that identifies invariant features within test-time graphs, and an environment refiner that adjusts varying features to align the latent spaces of training and test environments. These components are effectively optimized through a dual test-time graph contrastive learning objective and an iterative optimization strategy.
- Comprehensive Evaluation. We evaluate the proposed method on real-world test graphs under diverse graph distribution shifts. Extensive experimental results reveal consistently strong generalization ability on various well-trained GNN models, providing compelling evidence for the efficacy of the proposed method.

108 2 **RELATED WORK**

109

110 Test-time Adaption. Our work is related to test-time adaption methods, which aim to enable dynamic 111 tuning of the pre-trained model to generalize and adapt well to the test samples (Sun et al., 2020; 112 Wang et al., 2020; Chen et al., 2023a; Liang et al., 2020; Jang et al., 2022). The pioneering work is 113 Test-Time Training (TTT) (Sun et al., 2020), which re-trains model weights at test time on a single test image sample through a self-supervised learning objective with an auxiliary task. Moreover, 114 TENT (Wang et al., 2020) proposes a fully test-time adaption problem that only accesses model 115 parameters and the test data. Considering most existing methods are model-centric and serve the 116 computer vision domain, which might not fit GNN models and graph data, GTRANS (Jin et al., 2023) 117 first proposes to adapt test graph data without accessing the training procedure and GNN architectures. 118 However, GTRANS employs a fully parameterized matrix to represent modified test-time graph node 119 features and engages in graph structure learning within the dynamic node representation learning 120 process. For one thing, the fully parameterized approach to node feature learning merely adds to 121 the original node features, which narrows the scope for learning updated test graphs. For another 122 thing, GTRANS uses a binary-space projected gradient descent method, limiting the flexibility in 123 handling diverse graph structures. In this work, we conduct the parameterized decomposition of the 124 test graph to give it a rebirth by re-extracting the invariant features and re-fining test-time varying 125 features of test graphs. This allows for more adaptable modifications to test graphs with an expended optimization space of our proposed TT-GREB. 126

127

128 Graph Learning Under Distribution Shifts. Our work is also relevant to the research topic of graph learning under distribution shifts, whose goal is to develop a GNN model for better generalization 129 ability on test graphs under graph distribution shifts (Wu et al., 2020; Chen et al., 2023b; Sui et al., 130 2023; Guo et al., 2023; Chen et al., 2022a; Wang et al., 2022). Typically, UDAGCN (Wu et al., 131 2020) conducts unsupervised graph domain adaption with a domain adversarial method to learn 132 domain-invariant embeddings across the source domain and the target domain. AIA (Sui et al., 133 2023) proposes graph data augmentation with effective GNN learning to handle the covariate shift on 134 graphs for the graph classification task. Different from these graph model-centric methods, in this 135 work, we mainly focus on modifying the test graph data with self-supervision signals to deal with 136 the distribution-shifted test graph learning problem. Therefore, the critical distinction between our 137 work and existing methods is that our work is primarily concerned with a graph data-centric learning 138 paradigm to directly manipulate test graph data for serving better the GNN generalization ability 139 under distribution shifts.

140 141

142

TEST-TIME GRAPH REBIRTH (TT-GREB) 3

Preliminary. Given a training graph $G_{tr} = (\mathbf{X}_{tr}, \mathbf{A}_{tr}, \mathbf{Y}_{tr}) \sim P_{tr}$ with N nodes and C-classes of node labels, where $\mathbf{X}_{tr} \in \mathbb{R}^{N \times d}$ is the d-dimension nodes feature matrix indicating node attribute semantics, $\mathbf{A}_{tr} \in \mathbb{R}^{N \times N}$ is the adjacency matrix indicating whether nodes are connected or not by edges with $\mathbf{A}_{tr}^{i,j} = \{0,1\} \in \mathbb{R}$ for *i*-th and *j*-th nodes, $\mathbf{Y}_{tr} \in \mathbb{R}^{N \times C}$ denotes the node labels, and P_{tr} is the training matrix indicating matrix indicating the labels. 143 144 145 146 147 is the training graph distribution.

148 **Training Stage:** A GNN model is trained on G_{tr} according to the following objective function for 149 node classification: $oldsymbol{ heta}_{tr}^{*} = \min_{oldsymbol{ heta}_{tr}} \mathcal{L}_{cls}\left(\hat{\mathbf{Y}}_{tr}, \mathbf{Y}_{tr}
ight),$ where

(1)

150 151 152

 $\mathbf{Z}_{tr}, \hat{\mathbf{Y}}_{tr} = \text{GNN}_{\boldsymbol{\theta}_{tr}}(\mathbf{X}_{tr}, \mathbf{A}_{tr}).$ 153 The parameters of GNN trained on G_{tr} is denoted by θ_{tr} , $\mathbf{Z}_{tr} \in \mathbb{R}^{N \times d_1}$ is the output node embedding 154 of graph G_{tr} from $\text{GNN}_{\theta_{tr}}$, and $\hat{\mathbf{Y}}_{tr} \in \mathbb{R}^{N \times C}$ denotes the output node labels predicted by the trained 155 $\text{GNN}_{\theta_{\text{tr}}}$. By optimizing the node classification loss function \mathcal{L}_{cls} (*e.g.*, cross-entropy loss) between 156 GNN predictions \mathbf{Y}_{tr} and ground-truth node labels \mathbf{Y}_{tr} , the GNN model that is well-trained on G_{tr} 157 can be denotes as $\text{GNN}_{\theta_{\text{tr}}^*}$ with optimal weight parameters θ_{tr}^* . Note that once we obtain the optimal 158 $GNN_{\theta_{tr}}$ that has been well-trained on G_{tr} , the GNN model would be fixed and G_{tr} would not be 159 accessible during test time. 160

Test-time Inference: For practical GNN deployment and serving, given a real-world test graph 161 $G_{\text{te}} = (\mathbf{X}_{\text{te}}, \mathbf{A}_{\text{te}}) \sim P_{\text{te}}$ including M nodes with its feature matrix $\mathbf{X}_{\text{te}} \in \mathbb{R}^{M \times d}$ and its adjacency 162 matrix $\mathbf{A}_{te} \in \mathbb{R}^{M \times M}$, we assume that there are potential distribution shifts between G_{tr} and G_{te} , which mainly lies in node contexts, graph structures, and scales as $P_{tr} \neq P_{te}$, but the label space keeps 163 164 consistent under the covariate shift, i.e., all nodes in Gte are constrained in the same C-classes as 165 G_{tr} as $\{\mathbf{Y}_{\text{tr}}, \mathbf{Y}_{\text{te}}\} \in \mathcal{Y} = \{1, \dots, C\}$. Generally, the well-trained model $\text{GNN}_{\boldsymbol{\theta}_{\text{tr}}^*}$ would be directly 166 used for inferring on the test graph as $\hat{\mathbf{Y}}_{te} = \text{GNN}_{\boldsymbol{\theta}_{tr}^*}(\mathbf{X}_{te}, \mathbf{A}_{te})$. However, due to the latent graph 167 distribution shifts at the test time, the optimal parameters θ_{tr}^* learned on the training graph would not 168 be ideal for inference on the test-time graph. This can result in less accurate node classification on the test graph and does harm to the GNN's generalization ability. 169

170 Compared with existing model-centric methods working on learning optimal GNN parameter $\theta^*_{\{tr,te\}}$ 171 on the joint distribution of the training and test graphs, which requires GNN architecture re-designed 172 and fine-tuned, in this work, we pay attention to a data-centric solution through modifying the test 173 graph G_{te} under distribution shifts by test-time graph rebirth, without re-designing and fine-tuning 174 the well-trained $GNN_{\theta^*_{tr}}$, and without accessing the training graph G_{tr} .

175 176

177

190 191 192

3.1 PROBLEM FORMULATIONS

Through the length of graph structural data generation hypothesis in existing studies (Gui et al., 2022; Ye et al., 2022; Wu et al., 2021; Sui et al., 2023; Chen et al., 2022b; 2023b), a graph can be generated through a mapping $f_{gen} : \{C, S\} \to G$, where $C \subseteq \mathbb{R}^{n_c}$ and $S \subseteq \mathbb{R}^{n_s}$ are the latent variables denotes the environment-invariant part and the environment-varying parts for generating the graph $G \in \mathcal{G} = \bigcup_{N=1}^{\infty} \{0, 1\}^N \times \mathbb{R}^{N \times d}$, where n_c and n_s denote the dimensions of latent variables, respectively. Inspired by such structural causal model (SCM) (Chen et al., 2022b; 2023b) for graph generation progress, we have the following proposition to comprehend test-time graph distribution shifts and the test-time graph discrepancy from the training graph.

Proposition 1 Given the training graph $G_{tr} \sim P_{tr}$, there exist the environment-invariant part G_{tr}^{Inv} and the environment-varying part G_{tr}^{Env} components in this graph, denoting as $G_{tr} = \{G_{tr}^{Inv}, G_{tr}^{Env}\}$, likewise for the the test graph $G_{te} \sim P_{te}$ with $G_{te} = \{G_{te}^{Inv}, G_{te}^{Env}\}$. Then, the GNN model $GNN_{\theta_{tr}^*}$ that has been well trained on the G_{tr} would keep good generalization on the test graph when

$$G_{tr}^{Inv} = G_{te}^{Inv} \sim Q^{Inv}, \quad dist\left(G_{tr}^{Env}, G_{te}^{Env}\right) < \epsilon,$$

where $G_{tr}^{Env} \neq G_{te}^{Env}, G_{tr}^{Env} \sim Q_{tr}^{Env}, G_{te}^{Env} \sim Q_{te}^{Env}.$ (2)

193 In this proposition, the function 194 $dist(\cdot)$ quantifies the discrepancy be-195 tween the components of the train-196 ing graph and the test graph that vary with the environment. Additionally, 197 Q^{Inv} represents the distribution of la-198 tent variables that remain constant 199 across different environments, and is 200 expected to be identical in both the 201 training and test graphs. Furthermore, 202 the environmental variables of G_{tr}^{Env} 203 and $G_{\rm te}^{\rm Env}$ are presumed to adhere to 204 distinct, environment-specific distri-205 butions, denoted as Q_{tr}^{Env} and Q_{te}^{Env} , re-206 spectively.



Figure 2: Illustration of distribution shifts in test-time graphs v.s. training graphs.

As shown in Fig. 2, we elaborate on the distribution shifts in the test-time graph and the training graph, from the view of latent variable decomposition according to Proposition 1. It shows two fundamental insights and principles for test-time graph rebirth:

- Re-extracting the environment-invariant features of the test-time graph, which shares the same informative characteristics with the training graph, to assure that they can reflect the predictive pattern of node class labels, denoting as class-related **prototype** features.
- Re-fining the environment-varying features, which are primarily attributed to possible shifts in the training-test distribution, referred to **environment** features. Essentially, a well-trained GNN model is supposed to perform expressively on the test graph, when the test distribution closely matches



Figure 3: The overall framework of the proposed test-time graph rebirth (TT-GREB) method.

the training graph's distribution. When we push the test-time environment feature distributions more closely with those of the training graph, the well-trained GNN is likely to exhibit improved generalization capabilities under distribution shifts.

According to such two principles, if we (1) keep prototype features and (2) align the environment features on the test graph, then, make a re-composition, we could transform the original test graph to a new test graph, this process can be defined as the problem of test-time graph rebirth:

240 **Definition 3.1 (Test-time Graph Rebirth)** Given the test graph $G_{te} = (\mathbf{X}_{te}, \mathbf{A}_{te})$ and the well-241 trained GNN model GNN $_{\theta_n}$, test-time graph rebirth aims to learn following mapping functions: $f_{Pro}: G_{te} \to G_{te}^{Inv}$ and $f_{Env}: G_{te} \to G_{te}^{Env}$, with the re-composition function $g = f_{Pro} \circ f_{Env}$, the 242 rebirth test graph can be denoted as 243

$$G'_{te} = (\mathbf{X}'_{te}, \mathbf{A}'_{te}) = g(f_{Pro}(G_{te}), f_{Env}(G_{te})).$$

$$(3)$$

In this way, the rebirth test graph would be fed into the well-trained GNN model that has been deployed online in practice for making inference $\hat{\mathbf{Y}}'_{te} = GNN_{\theta_t}(\mathbf{X}'_{te}, \mathbf{A}'_{te})$, where $\hat{\mathbf{Y}}'_{te}$ is expected to be more closely aligned with the actual ground-truth node labels of the test graph compared to the initial predictions $\hat{\mathbf{Y}}_{te}$.

3.2 METHODOLOGY

231 232 233

235 236

237

238

239

244 245 246

247

248 249

250 251

252

266

253 According to the two principles, in this work, we propose a novel method, named TT-GREB, to 254 address the test-time graph rebirth problem for serving good GNN generalization under distribution 255 shifts at the test time. 256

The overall framework is illustrated in Fig. 3. Concretely, the proposed TT-GREB consists of 257 two components: (1) a prototype extractor identifies features that remain unchanged in different 258 environments, mainly determined by node class labels, where these features can be consistent in both 259 training and test graphs and reflect the predictive pattern of GNN models; and (2) an environment 260 refiner adjusts the environment-varying features of the test-time graph, to match the latent distribution 261 of the training graph's environment. This alignment ensures that the GNN, which is well-trained 262 on the training graph, demonstrates strong generalization capability on the rebirth test graph. More 263 details of the modular design of our proposed TT-GREB are presented below. 264

265 3.2.1 MODULAR DESIGN.

Given a test-time graph $G_{\text{te}} = (\mathbf{X}_{\text{te}}, \mathbf{A}_{\text{te}})$ with M nodes with d dimension node attribute features, the prototype extractor $f_{\phi_p}^{\text{Pro}}(\cdot)$ and the environment refiner $f_{\phi_e}^{\text{Env}}(\cdot)$, parameterized by ϕ_p and ϕ_e , 267 268 respectively, take it as the input simultaneously. For ease of reference, we denote them as $f_{\phi_n}(\cdot)$ and 269 $f_{\phi_e}(\cdot)$, by omitting the superscripts in subsequent mentions.

Concretely, these two sub-modules share the same structure, *i.e.*, two full-connected layers, $FC_{node}^{k}(\cdot)$ and $FC_{edge}^{k}(\cdot)$ to generate the soft and dense node attribute reweight matrix $\mathbf{W}_{node}^{k} \in \mathbb{R}^{d \times d}$, and the edge reweight matrix $\mathbf{W}_{edge}^{k} \in \mathbb{R}^{M \times M}$, where $k = \{Pro, Env\}$ for indicating the layers in the prototype extractor and the environment refiner, respectively. In this way, we have:

$$\mathbf{W}_{\text{node}}^{k} = \sigma \left(FC_{\text{node}}^{k}(\mathbf{X}_{\text{te}}) \right),
\mathbf{W}_{\text{edge}}^{k}(i,j) = \sigma \left(FC_{\text{edge}}^{k}\left(\left[\mathbf{x}_{\text{te}}^{i}, \mathbf{x}_{\text{te}}^{j} \right] \right) \right),$$
(4)

where $\sigma(\cdot)$ is the sigmoid function that constrains both the node attribute reweight matrix and the edge reweight matrix to [0, 1], and $\mathbf{x}_{te}^i = \mathbf{X}_{te}[i, :]$ denotes the *i*-th row's feature representation for node *i*, so as for node *j* with $\mathbf{x}_{te}^j = \mathbf{X}_{te}[j, :]$, attributing the value in the location (i, j) for \mathbf{W}_{edge}^k .

Then, we could obtain the graph $G_{\text{Pro}} = (\mathbf{X}_{\text{Pro}}, \mathbf{A}_{\text{Pro}})$ that reflects the node class label prototypes and the graph $G_{\text{Env}} = (\mathbf{X}_{\text{Env}}, \mathbf{A}_{\text{Env}})$ that adjusts the test-time graph characteristics that vary with environments under distribution shifts, as:

$$G_{\text{Pro}} = f_{\phi_p}(G_{\text{te}}) = \left(\mathbf{A}_{\text{te}} \odot \mathbf{W}_{\text{edge}}^{\text{Pro}}, \mathbf{X}_{\text{te}} \odot \mathbf{W}_{\text{node}}^{\text{Pro}}\right),$$

$$G_{\text{Env}} = f_{\phi_e}(G_{\text{te}}) = \left(\mathbf{A}_{\text{te}} \odot \mathbf{W}_{\text{edge}}^{\text{Env}}, \mathbf{X}_{\text{te}} \odot \mathbf{W}_{\text{node}}^{\text{Env}}\right),$$
(5)

where \odot is the broadcasted element-wise product. After these, we make a re-composition with the prototype graph G_{Pro} and the environment graph G_{Env} components to build a new test-time graph $G'_{\text{te}} = (\mathbf{X}'_{\text{te}}, \mathbf{A}'_{\text{te}})$ through:

291 292 293

275 276

277

278

279

281

283

284

285

287

289

290

296 297

298

 $G'_{te} = g(G_{te}) = \left(\mathbf{A}_{te} \odot \mathbf{W}_{edge}^{Comp}, \mathbf{X}_{te} \odot \mathbf{W}_{node}^{Comp}\right), \text{ where}$ $\mathbf{W}_{edge}^{Comp} = \left(\mathbf{1}^{edge} - \mathbf{W}_{edge}^{Pro}\right) \odot \mathbf{W}_{edge}^{Env} + \mathbf{W}_{edge}^{Pro}, \text{ and}$ $\mathbf{W}_{node}^{Comp} = \left(\mathbf{1}^{node} - \mathbf{W}_{node}^{Pro}\right) \odot \mathbf{W}_{node}^{Env} + \mathbf{W}_{node}^{Pro}.$ (6)

$$\mathbf{W}_{node}^{Comp} = \left(\mathbf{1}^{node} - \mathbf{W}_{node}^{Pro}\right) \odot \mathbf{W}_{node}^{Env} + \mathbf{W}_{node}^{Pro}.$$

In this process, $\mathbf{1}^{node}$ and $\mathbf{1}^{edge}$ denote the all-one matrices. Given \mathbf{W}_{node}^{Pro} represents the class-prototype node attribute reweight matrix, $\left(\mathbf{1}^{node} - \mathbf{W}_{node}^{Pro}\right)$ can be viewed as a plain and straightforward

node attri tforward 299 proportion of the environment-sensitive node attributes. Then, $(1^{node} - \mathbf{W}_{node}^{Pro}) \odot \mathbf{W}_{node}^{Env}$ re-composes 300 node attribute reweight matrix by explicitly imposing the environment refinement \mathbf{W}_{node}^{Env} on the 301 environment-sensitive proportion. And then, $+\mathbf{W}_{node}^{Pro}$ makes sure to preserve the environment 302 consistent proportion. The edge reweight matrix composition would follow the same rule. The 303 rationale for such a re-composition schema is based on the understanding that the interplay between 304 node class label prototypes and environment features is typically more complex than a basic additive 305 combination, such as $(\mathbf{W}_{node}^{Env} + \mathbf{W}_{node}^{Pro})$. This complexity becomes particularly evident under the 306 distribution shifts encountered during test time. 307

By this re-composition schema, we jointly keep prototype features and align the environment features
 on the test graph, leading to a transformation from the original test graph to a newly reborn test graph.
 In this way, the new test-time rebirth graph can make effective predictions with good generalization
 ability on the well-trained GNN model with graph data distribution shifts.

312 3.2.2 OPTIMIZATION OBJECTIVE.

314 During the test time, to improve the well-trained GNN model generalization under graph distribution 315 shifts, the significant challenge faced with graph data-centric transformation through test-time graph rebirth is the scarcity of test ground-truth labels. Consequently, this makes it more challenging to 316 conduct supervised learning by minimizing the cross entropy loss, which is the most readily and 317 straightforward solution. Therefore, with (1) the insufficient node class labels of the test graph, (2) the 318 inaccessible training graph for online GNN deployment, and (3) the unknown graph distribution shifts, 319 it is imperative to develop an effective self-supervised learning objective along with an appropriate 320 optimization strategy. 321

In this work, we propose a dual test-time graph contrastive learning objective with an effective iterative optimization strategy. For one thing, we use self-supervise signals from the well-trained GNN's output node representations to guide the learning of prototype extractor with the graph 324 contrastive learning loss \mathcal{L}_{Pro} , following the parameter-free principle (Jin et al., 2023). Through the 325 lens of the general graph contrastive learning scheme, the core idea is to maximize the similarity 326 between two consistent views of the same graph, and to minimize the similarity when the views are 327 not in agreement. For another thing, we perform a decomposition of environment-varying features on 328 the test graph, by ensuring the environmental discrepancy under the graph distribution shifts during the test time. Considering the inaccessible training graph, we encourage the discrepancy between the environmental characteristics of the reborn test-time graph and the original test-time graph, leading 330 to the graph environment refinement loss \mathcal{L}_{Env} to optimize the proposed environment refiner. As 331 illustrated in the lower section of Fig.3, these two optimization objectives are iteratively refined using 332 gradient descent until they reach convergence. More implement details of the dual learning objective 333 of our proposed TT-GREB are presented as follows. 334

Given the obtained prototype graph G_{Pro} , the environment graph G_{Env} , and the new rebirth graph G'_{te} , we fed them into the well-trained GNN model simultaneously, leading to the node representations with $\mathbf{Z}_{Pro} = \text{GNN}_{\theta_{tr}^*}(G_{Pro})$, $\mathbf{Z}_{Env} = \text{GNN}_{\theta_{tr}^*}(G_{Env})$, and $\mathbf{Z}_{te'} = \text{GNN}_{\theta_{tr}^*}(G'_{te})$. Then, the learning objective of structural prototype feature extraction for test-time graph rebirth can be:

$$\min_{\boldsymbol{\phi}_{p}} \mathcal{L}_{\text{Pro}} = \sum_{i=1}^{M} \left(1 - \frac{(\mathbf{z}_{i}^{\text{te'}})^{\top} \mathbf{z}_{i}^{\text{Pro}}}{\||\mathbf{z}_{i}^{\text{ter'}}||} \right) - \sum_{i=1}^{N} \left(1 - \frac{(\mathbf{z}_{i}^{\text{Env}})^{\top} \mathbf{z}_{i}^{\text{Pro}}}{\||\mathbf{z}_{i}^{\text{Pro}}\||} \right) + \alpha \left[Reg \left(\mathbf{W}_{\text{node}}^{\text{Pro}} \right) + Reg \left(\mathbf{W}_{\text{edge}}^{\text{Pro}} \right) \right].$$
(7)

341 342 343

344 345

347 348

349 350 351

352

353

354

355 356

357 358

359

360

361

362

364

339 340

Furthermore, the learning objective of environment refinement can be written as:

$$\max_{\boldsymbol{\phi}_{e}} \mathcal{L}_{\text{Env}} = dist(G_{\text{te}}, G_{\text{Env}}) = ||\mathbf{Z}_{\text{te}} - \mathbf{Z}_{\text{Env}}||_{2}^{2} - \beta \left[Reg\left(\mathbf{W}_{\text{node}}^{\text{Env}}\right) + Reg\left(\mathbf{W}_{\text{edge}}^{\text{Env}}\right)\right],$$
(8)

where $Reg(\mathbf{W}^{\bigstar}) = |\frac{\sum \mathbf{W}^{\bigstar}}{\sum(1-\mathbf{W}^{\bigstar})} - \lambda_s|$ is the regularization term with $s = \{1, 2, 3, 4\}$ and superscript corresponding to \mathbf{W}_{node}^{Pro} , \mathbf{W}_{node}^{Pro} , and \mathbf{W}_{edge}^{Env} , respectively. Here, λ_s acts as a hyper-parameter ranging [0, 1], while α and β balance the loss functions between the primary optimization objectives and the regularization terms. These regularization terms are designed to keep the average ratio of the number of reweighted node features or edges close to λ_s , thereby stabilizing the training process and avoiding trivial solutions.

4 EXPERIMENT

In this section, we verify the effectiveness of the proposed TT-GREB in terms of the GNN generalization ability on test-time graphs under distribution shifts. Concretely, we aim to answer the following questions to demonstrate the effectiveness of the proposed TT-GREB: Q1: How does the proposed TT-GREB perform on the well-trained GNNs for node classification task under various graph distribution shifts at test time? Q2: How does the proposed TT-GREB perform when conducting an ablation study regarding the sub-module components and the learning strategy? Q3: How sensitive are the hyper-parameter λ for the proposed TT-GREB? Q4: How does the proposed TT-GREB perform in terms of running time efficiency and visualization?

365 366 367

368

4.1 EXPERIMENTAL SETTINGS

Datasets. We perform experiments on five real-world graph datasets with diverse graph data distribution shifts containing: node feature shifts: Cora (Yang et al., 2016) and Amazon-Photo (Shchur et al., 2018); domain shifts (Wu et al., 2020): Twitch-E (Rozemberczki et al., 2021); temporal shifts: Elliptic (Pareja et al., 2020) and OGB-arxiv (Pareja et al., 2020). More details of datasets are listed in Appendix A. For all training, validation, and test graphs, we follow the process procedures and splits in previous works (Wu et al., 2022b; Jin et al., 2023; Wu et al., 2020; Zheng et al., 2023b).

Test-time Evaluation Protocol. We test four commonly used GNN models for evaluating GNN generalization under graph distribution shifts following the settings in (Jin et al., 2023), including GCN (Kipf & Welling, 2017), GraphSAGE (Hamilton et al., 2017) (*abbr.* SAGE), GAT (Veličković et al., 2017), and GPR-GNN (Chien et al., 2020) (*abbr.* GPR). For each model, we train it on training

Backbones	Categories	Methods	Amz-Photo	Cora	Elliptic	OGB-Arxiv	Twitch-E	Rank
		ERM	$88.60{\scriptstyle \pm 0.90}$	$87.49{\scriptstyle \pm 7.97}$	$51.09{\scriptstyle\pm5.63}$	$38.39{\scriptstyle\pm2.92}$	$59.80{\scriptstyle \pm 3.77}$	3.8
	Model-centric	EERM	$81.05{\scriptstyle\pm0.95}$	$66.80{\scriptstyle \pm 6.51}$	$45.60{\scriptstyle\pm1.22}$	OOM	$53.28{\scriptstyle\pm1.88}$	6
GCN	woder-centre	DropEdge	$81.73{\scriptstyle\pm1.23}$	$74.05{\scriptstyle\pm8.00}$	$53.83{\scriptstyle \pm 4.52}$	$40.82{\scriptstyle\pm2.18}$	$59.49{\scriptstyle\pm4.14}$	3.8
CON		TENT	$88.60{\scriptstyle \pm 0.90}$	$87.51{\scriptstyle \pm 8.01}$	$47.05{\scriptstyle \pm 2.01}$	$38.45{\scriptstyle\pm2.35}$	$59.79{\scriptstyle \pm 3.77}$	3.8
	Data centric	GTRANS	$89.27{\scriptstyle\pm0.37}$	$\underline{95.20{\scriptstyle\pm0.87}}$	$\underline{56.69{\scriptstyle\pm 6.74}}$	$\underline{40.00{\scriptstyle\pm2.30}}$	$\underline{60.38{\scriptstyle\pm3.86}}$	1.8
	Data-centre	TT-GREB (Ours)	$\underline{89.11{\scriptstyle \pm 0.47}}$	$96.12{\scriptstyle \pm 1.10}$	57.20±8.19	$39.49{\scriptstyle\pm1.72}$	$60.85{\scriptstyle \pm 4.17}$	1.6
		ERM	$84.03{\scriptstyle \pm 7.61}$	$98.48{\scriptstyle \pm 3.68}$	$57.34{\scriptstyle \pm 5.95}$	$39.26{\scriptstyle\pm2.39}$	$62.08{\scriptstyle \pm 4.04}$	4.4
	Model-centric	EERM	$84.97{\scriptstyle\pm7.26}$	$96.73{\scriptstyle\pm6.77}$	$60.94{\scriptstyle\pm5.18}$	OOM	$61.70{\scriptstyle \pm 4.23}$	4.6
SAGE	woder-centre	DropEdge	$80.67{\scriptstyle\pm1.61}$	$92.53{\scriptstyle\pm7.12}$	$52.84{\scriptstyle\pm3.92}$	$37.90{\scriptstyle\pm1.74}$	$\underline{62.19{\scriptstyle\pm4.16}}$	4.8
UNCL		TENT	$84.10{\scriptstyle\pm7.71}$	$98.58{\scriptstyle \pm 3.49}$	$50.16{\scriptstyle \pm 3.89}$	$\underline{39.59{\scriptstyle\pm1.63}}$	$62.04{\scriptstyle\pm4.06}$	3.8
	Data-centric	GTRANS	89.63 ±5.43	$99.89{\scriptstyle \pm 0.03}$	$62.54{\scriptstyle\pm7.94}$	$39.49{\scriptstyle\pm2.34}$	$62.04{\scriptstyle\pm4.06}$	2
		TT-GREB (Ours)	$\underline{88.49{\scriptstyle\pm4.07}}$	$\underline{99.66{\scriptstyle\pm0.48}}$	66.97±8.94	$40.15{\scriptstyle \pm 1.65}$	$62.43{\scriptstyle\pm4.26}$	1.4
		ERM	$91.20{\scriptstyle \pm 2.41}$	$95.53{\scriptstyle\pm4.98}$	65.28 ± 9.59	$40.47{\scriptstyle\pm2.48}$	$58.23{\scriptstyle \pm 3.45}$	3.8
	Model contrie	EERM	$89.13{\scriptstyle \pm 4.06}$	$87.04{\scriptstyle\pm11.07}$	$50.40{\scriptstyle \pm 3.48}$	OOM	$59.51{\scriptstyle \pm 3.26}$	4.6
GAT	Model-centric	DropEdge	$69.52{\scriptstyle\pm6.33}$	$76.71{\scriptstyle \pm 4.60}$	$64.96{\scriptstyle\pm7.12}$	$43.91{\scriptstyle\pm1.93}$	$58.46{\scriptstyle \pm 3.35}$	4
UAI		TENT	$91.40{\scriptstyle \pm 2.36}$	$95.57{\scriptstyle \pm 4.96}$	$56.86{\scriptstyle \pm 5.10}$	$30.36{\scriptstyle\pm1.20}$	$58.23{\scriptstyle \pm 3.45}$	4
	Data-centric	GTRANS	$\underline{94.04}{\scriptstyle \pm 0.73}$	$\underline{97.28}{\scriptstyle \pm 2.92}$	66.85±9.80	$\underline{41.65{\scriptstyle\pm2.26}}$	$58.20{\scriptstyle \pm 3.49}$	2.4
		TT-GREB (Ours)	$94.34{\scriptstyle \pm 0.82}$	$98.05{\scriptstyle\pm1.03}$	$\underline{66.05{\scriptstyle\pm8.92}}$	$41.45{\scriptstyle\pm2.00}$	$\underline{58.53{\scriptstyle\pm3.50}}$	1.8
		ERM	$87.04{\scriptstyle\pm2.86}$	$87.24{\scriptstyle \pm 9.11}$	$64.79{\scriptstyle \pm 7.26}$	$44.38{\scriptstyle \pm 2.97}$	$59.77{\scriptstyle\pm3.73}$	3.4
	Model contrie	EERM	$85.29{\scriptstyle\pm1.48}$	89.50±7.83	$64.41 {\pm} 6.97$	OOM	$61.76{\scriptstyle \pm 4.06}$	3
GPR	woder-centric	DropEdge	$74.20{\scriptstyle\pm6.90}$	$73.29{\scriptstyle\pm10.19}$	$60.62{\scriptstyle\pm6.06}$	$43.96{\scriptstyle \pm 2.37}$	$59.89{\scriptstyle \pm 3.99}$	4.6
Urk		TENT	-	-	-	-	-	-
	Dete centri	GTRANS	$86.94{\scriptstyle\pm2.62}$	87.45 ± 8.91	$67.65{\scriptstyle\pm10.49}$	45.74±2.24	$59.89{\scriptstyle\pm3.61}$	2.4
	Data-Cellulic	TT-GREB (Ours)	$88.55{\scriptstyle \pm 1.68}$	$\underline{88.54{\scriptstyle\pm8.74}}$	$71.34{\scriptstyle \pm 10.01}$	$\underline{45.14}_{\pm 2.41}$	$\underline{60.00{\scriptstyle\pm3.86}}$	1.6

Table 1: Average classification results (%) over the test graphs under various graph distribution shifts
 on different backbone GNN models. The best results are in bold, and the second-bests are with
 underlines. 'Rank' indicates the average rank of each algorithm for each backbone; 'OOM' indicates an out-of-memory error on 32 GB
 GPU memory; TENT with '-' means it cannot be applied to GNNs without batch normalization layers.

407 408

sets, until the model achieves the optimal node classification on its validation sets following the
standard training process, so that we can obtain the 'well-trained' GNN model that keeps fixed in the
whole test-time graph rebirth process. We report the average classification performance, and for all
experiments, we report the average results of 10 repeated times with different random seeds.

Baseline Methods. We compare the proposed TT-GREB with the following baselines that fall in two groups: graph model-centric methods: empirical risk minimization (ERM) for standard training (Wu et al., 2022b), data augmentation technique DropEdge (Rong et al., 2019), Explore-to-Extrapolate Risk Minimization (EERM) (Wu et al., 2022b) customized for node-level graph OOD generalization, and test-time training method <u>TENT</u> (Wang et al., 2020); And the recent SOTA graph data-centric method: test-time graph transformation method <u>GTRANS</u> (Jin et al., 2023). More demonstrations of the differences among these baselines are presented in Appendix A.

420 421

422

4.2 EXPERIMENTAL RESULTS

In Table 1, we report the average node-level classification results over the test graphs under various
 graph distribution shifts on different backbone GNN models, along with the average rank of each
 comparison method for each backbone.

As can be observed, our proposed TT-GREB generally delivers great performance across various graph datasets and models, achieving the highest ranks overall: 1.6, 1.4, 1.8, and 1.6 for GCN, SAGE, GAT, and GPR, respectively. These results could verify the outstanding effectiveness of the proposed TT-GREB for modifying graph data at test time to serve better GNN generalization ability.

431 Moreover, compared with the recent SOTA graph data-centric method GTRANS, our proposed TT-GREB achieves significant improvements in some cases: for example, our method has 5.5%

average improvement of the classification performance from GTRANS's 67.65% to 71.34% on
 Elliptic dataset with GPR model, and 7.1% average improvement with SAGE model, respectively.

Besides, we can also observe that some model-centric comparison methods, achieve great performance in some cases, for instance, DropEdge and EERM could deliver excellent performance on OGB-arxiv with GCN and Twitch-e with GAT models, respectively. Nevertheless, DropEdge and EERM can not be applied to the test-time application scenario, since it has to modify the training process of GNNs to achieve better generalization ability. Although TENT is suited for the test-time adaption and training scenario, it can not be directly used for models without batch normalization layers, significantly limiting its usage on GNN models.

In summary, our proposed TT-GREB significantly improves GNN generalization for different graph distribution shifts and GNN models during test time, achieving superior average rankings compared to existing approaches. This success is due to the collaboration between the prototype extractor and environment refiner, which enhances the representations of test graph nodes and edges. Additionally, the incorporation of a dual graph contrastive learning objective, coupled with an effective iterative optimization strategy, further contributes to the method's outstanding performance.

447 448

449

462

463

464

465

466

467

468

477

478 479

4.3 ABLATION STUDY OF TT-GREB

450 In Table 2, we evaluate the effec-451 tiveness of the overall framework of 452 the proposed TT-GREB, from the perspectives of sub-module compo-453 nents and learning strategies, respec-454 tively. We observe the effectiveness 455 of the prototype extractor (ProExtrac-456 tor), and the environment refiner (En-457 vRefiner), respectively. For learning 458 strategies, we test the effectiveness 459 of with and without the dual graph 460 contrastive learning objective (Con-461 trastive Obj) as well as the iterative

Table 2: Ablation study components of the proposed method. For Idx02, \checkmark^* denotes enabling the complete framework in the test-time graph rebirth process but only using the output G_{Pro} of the prototype extractor for final inference.

Ablation Index	Sub-module	Components	Learning Strategies		
	ProExtractor	EnvRefiner	Contrastive_Obj	Iterative_Opt	
Baseline	×	×	×	×	
Idx00	\checkmark	×	×	×	
Idx01	\checkmark	\checkmark	✓	×	
Idx02	√*	\checkmark	✓	\checkmark	
Idx03 (TT-GREB)	\checkmark	\checkmark	✓	\checkmark	

optimization method (Iterative_Opt). Baseline denotes the original test graph classification performance directly inferring on the well-trained GNNs without any test-time modification. For Idx00 without the contrastive learning objective, the optimization objective would be degraded to the close distance constraint between the output prototype extractor and the original test graph, which means with the weakest supervision signals to instruct the learning process. For Idx02 with \checkmark^* , it denotes that we use the overall proposed framework to give a test-time graph rebirth, but only access the partial output, *i.e.*, the output G_{Pro} of the prototype extractor for final test-time inference.



Figure 4: Ablation study results (%) on Cora with GCN, SAGE, and GAT models demonstrated with Box-plot on all test graphs.

The results on all test graphs of Cora on GCN, SAGE, and GAT with a fixed seed run are presented
 in Fig. 4. As can be observed, our proposed TT-GREB achieves consistently good classification
 performance on all test graphs on average, also with the smallest standard deviations across all
 models (shown in Appendix B). Besides, it can also be observed that, generally, each component
 of sub-modules and learning strategies could contribute to performance improvement in different
 degrees when these components are coupled together to achieve the best performance of the proposed
 TT-GREB.



Figure 5: Hyper-parameter sensitivity study results (%) on Cora and Elliptic with GCN model, from left to right: (1) α on Elliptic with GCN, (2) β on Elliptic with GCN, (3) λ_{Pro} on Cora with GCN.

4.4 HYPER-PARAMETER SENSITIVITY ANALYSIS

In Fig. 5, we evaluate the sensitivity of three hyper-parameters in terms of the regularization in our proposed TT-GREB in Eq. (7) and Eq. (8). Concretely, λ_{Pro} indicates $s = \{1, 2\}$ in Eq. (7) and Eq. (8) for $\lambda_1 = \lambda_2$, corresponding to \mathbf{W}_{node}^{Pro} , \mathbf{W}_{edge}^{Pro} , and we empirically set the $\lambda_3 = 1$ for \mathbf{W}_{node}^{Env} . More observation on the sensitivity of hyper-parameters, *i.e.*, λ_{Env} indicates λ_4 , corresponding to \mathbf{W}_{edge}^{Pro} is presented in Appendix B. For regularization weights in balancing the optimization objective, we observe α and β in a set of $\{0, 0.1, 0.5, 1, 10, 100\}$. For λ_{Pro} and λ_{Env} , we observe the parameter range of [0, 1] with an interval of 0.1. Observations indicate that the hyper-parameters demonstrate a moderate level of sensitivity within specific ranges, underscoring the robustness of our proposed method to hyper-parameter tuning. More results on visualization are listed in Appendix C.

4.5 RUNNING TIME COMPARISON

512 In Table 3, we compare the running 513 time of our proposed TT-GREB with 514 existing baseline methods, *i.e.*, EERM 515 and GTRANS, which are specifically designed for the graph distribution 516 shift issue. The results are obtained in 517 a single NVIDIA A100 GPU across 518 all datasets with GCN model in 20 519 epochs. It can be observed that our 520 proposed method achieves a compara-521

Table 3: Running time (seconds) comparison in 20 epochs with a single NVIDIA A100 GPU on all graph distribution shift datasets with GCN model.

Methods	Amz-Photo	Cora	Elliptic	OGB-Arxiv	Twitch-E
EERM GTRANS	14.66 0.24	2.67 0.13	230.50 0.32	191.48 0.89	22.14 0.20
TT-GREB (Ours)	2.06	1.08	1.53	4.29	1.76

ble running time with GTRANS, and significantly exceeds the EERM method, demonstrating its
 great time efficiency. This efficiency stems primarily from the fact that EERM, a graph data-centric
 method, necessitates retraining the GNN, which is inherently time-consuming. In contrast, both
 GTRANS and our proposed method employ test-time graph modifications to enhance performance.
 However, our method incurs a slight increase in time consumption due to the implementation of a
 dual iterative optimization strategy.

527 528

495

496 497 498

499 500

501

502

503

504

505

506

507

508

509 510

511

5 CONCLUSION

529 530

In this work, we proposed a new graph data-centric method, test-time graph rebirth (TT-GREB), 531 aimed at enhancing the generalization ability of GNN models to test-time graphs affected by distri-532 bution shifts through direct manipulation of the test graph data. The overall framework includes a 533 prototype extractor for learning environment-invariant features and an environment refiner for adjust-534 ing environment-sensitive features, followed by a dual test-time graph contrastive learning objective and an efficient iterative optimization strategy, facilitating the extraction of optimal prototype and 536 environmental components of the reborn test graph. Our extensive experiments on real-world graph 537 datasets under various test-time distribution shifts confirm the superiority of our method, underscoring its innovative capacity to modify test-time graphs for enhanced GNN generalization. A potential 538 limitation of this work is its current focus on node-level tasks, but future extensions are expected to adapt it for broader applications in graph-level and edge-level tasks.

540 REFERENCES

- Guanzi Chen, Jiying Zhang, Xi Xiao, and Yang Li. Graphtta: Test time adaptation on graph neural
 networks. *arXiv preprint arXiv:2208.09126*, 2022a.
- Liang Chen, Yong Zhang, Yibing Song, Ying Shan, and Lingqiao Liu. Improved test-time adaptation for domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 24172–24182, 2023a.
- Yongqiang Chen, Yonggang Zhang, Yatao Bian, Han Yang, MA Kaili, Binghui Xie, Tongliang Liu, Bo Han, and James Cheng. Learning causally invariant representations for out-of-distribution generalization on graphs. *Advances in Neural Information Processing Systems (NeurIPS)*, 35: 22131–22148, 2022b.
- Yongqiang Chen, Yatao Bian, Kaiwen Zhou, Binghui Xie, Bo Han, and James Cheng. Does invariant graph learning via environment augmentation learn invariance? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023b.
- Yongqiang Chen, Wei Huang, Kaiwen Zhou, Yatao Bian, Bo Han, and James Cheng. Towards understanding feature learning in out-of-distribution generalization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023c.
- Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized pagerank
 graph neural network. In *International Conference on Learning Representations (ICLR)*, 2020.
- Shurui Gui, Xiner Li, Limei Wang, and Shuiwang Ji. Good: A graph out-of-distribution benchmark.
 Advances in Neural Information Processing Systems (NeurIPS), 35:2059–2073, 2022.
- Yuxin Guo, Cheng Yang, Yuluo Chen, Jixi Liu, Chuan Shi, and Junping Du. A data-centric framework to endow graph neural networks with out-of-distribution detection ability. In Ambuj K. Singh, Yizhou Sun, Leman Akoglu, Dimitrios Gunopulos, Xifeng Yan, Ravi Kumar, Fatma Ozcan, and Jieping Ye (eds.), *Proceedings of ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 638–648. ACM, 2023.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.
 Advances in Neural Information Processing Systems (NeurIPS), 2017.
- 572 Minguk Jang, Sae-Young Chung, and Hye Won Chung. Test-time adaptation via self-training with
 573 nearest neighbor information. In *International Conference on Learning Representations (ICLR)*,
 574 2022.
- 575
 576
 576
 577
 578
 578
 578
 579
 579
 570
 570
 571
 572
 573
 574
 575
 575
 576
 577
 578
 577
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
- Wei Jin, Tong Zhao, Jiayuan Ding, Yozen Liu, Jiliang Tang, and Neil Shah. Empowering graph
 representation learning with test-time graph transformation. In *International Conference on Learning Representations (ICLR)*, 2023.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source
 hypothesis transfer for unsupervised domain adaptation. In *International Conference on Machine Learning (ICML)*, pp. 6028–6039. PMLR, 2020.
- Hongrui Liu, Binbin Hu, Xiao Wang, Chuan Shi, Zhiqiang Zhang, and Jun Zhou. Confidence may cheat: Self-training on graph neural networks under distribution shift. In *Proceedings of the ACM Web Conference 2022*, pp. 1248–1258, 2022.
- Yixin Liu, Kaize Ding, Huan Liu, and Shirui Pan. Good-d: On unsupervised graph out-of-distribution
 detection. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining (WSDM)*, pp. 339–347, 2023a.

- Yixin Liu, Yizhen Zheng, Daokun Zhang, Vincent Lee, and Shirui Pan. Beyond smoothing: Unsupervised graph representation learning with edge heterophily discriminating. In *Proceedings of the Association for the Advanced of Artificial Intelligence (AAAI)*, 2023b.
- Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kanezashi, Tim Kaler, Tao Schardl, and Charles Leiserson. Evolvegcn: Evolving graph convolutional networks for dynamic graphs. In *Proceedings of AAAI Conference on Artificial Intelligence*, volume 34, pp. 5363–5370, 2020.
- Yu Rong, Wenbing Huang, Tingyang Xu, and Junzhou Huang. Dropedge: Towards deep graph convolutional networks on node classification. In *International Conference on Learning Representations* (*ICLR*), 2019.
- Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding. *Journal* of Complex Networks, 9(2):cnab014, 2021.
- Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls
 of graph neural network evaluation. *arXiv preprint arXiv:1811.05868*, 2018.
- Yongduo Sui, Qitian Wu, Jiancan Wu, Qing Cui, Longfei Li, Jun Zhou, Xiang Wang, and Xiangnan
 He. Unleashing the power of graph data augmentation on covariate distribution shift. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *International Conference on Machine Learning*, pp. 9229–9248. PMLR, 2020.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
 Bengio. Graph attention networks. In *International Conference on Learning Representations* (*ICLR*), 2017.
- Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test time adaptation by entropy minimization. In *International Conference on Learning Representations* (*ICLR*), 2020.
- Yiqi Wang, Chaozhuo Li, Wei Jin, Rui Li, Jianan Zhao, Jiliang Tang, and Xing Xie. Test-time training for graph neural networks. *arXiv preprint arXiv:2210.08813*, 2022.
- Man Wu, Shirui Pan, Chuan Zhou, Xiaojun Chang, and Xingquan Zhu. Unsupervised domain adaptive graph convolutional networks. In *Proceedings of The Web Conference (WWW)*, pp. 1457–1467, 2020.
- Man Wu, Shirui Pan, and Xingquan Zhu. Attraction and repulsion: Unsupervised domain adaptive graph contrastive learning network. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(5):1079–1091, 2022a.
- Qitian Wu, Hengrui Zhang, Junchi Yan, and David Wipf. Handling distribution shifts on graphs: An invariance perspective. In *International Conference on Learning Representations (ICLR)*, 2022b.
- Yingxin Wu, Xiang Wang, An Zhang, Xiangnan He, and Tat-Seng Chua. Discovering invariant
 rationales for graph neural networks. In *International Conference on Learning Representations* (*ICLR*), 2021.
- Kaidi Xu, Hongge Chen, Sijia Liu, Pin-Yu Chen, Tsui-Wei Weng, Mingyi Hong, and Xue Lin. Topology attack and defense for graph neural networks: An optimization perspective. In *Proceedings of International Joint Conference on Artificial Intelligence (IJCAI)*, 2019.
- Karal Karal
- Nanyang Ye, Kaican Li, Haoyue Bai, Runpeng Yu, Lanqing Hong, Fengwei Zhou, Zhenguo Li, and Jun Zhu. Ood-bench: Quantifying and understanding two dimensions of out-of-distribution generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7947–7958, 2022.

648	Yuning You Tianlong Chen Zhazengyang Wang and Yang Shen Graph domain adaptation via
649	theory-grounded spectral regularization. In International Conference on Learning Representations
650	(<i>ICLR</i>), 2023.

- Junchi Yu, Jian Liang, and Ran He. Mind the label shift of augmentation-based graph ood generaliza-tion. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11620-11630, 2023.
- He Zhang, Bang Wu, Xingliang Yuan, Shirui Pan, Hanghang Tong, and Jian Pei. Trustworthy graph neural networks: Aspects, methods and trends. arXiv preprint arXiv:2205.07424, 2022.
- Xin Zheng, Yixin Liu, Shirui Pan, Miao Zhang, Di Jin, and Philip S Yu. Graph neural networks for graphs with heterophily: A survey. arXiv preprint arXiv:2202.07082, 2022a.
- Xin Zheng, Miao Zhang, Chunyang Chen, Chaojie Li, Chuan Zhou, and Shirui Pan. Multi-relational graph neural architecture search with fine-grained message passing. In 2022 IEEE International Conference on Data Mining (ICDM), pp. 783–792. IEEE, 2022b.
- Xin Zheng, Yixin Liu, Zhifeng Bao, Meng Fang, Xia Hu, Alan Wee-Chung Liew, and Shirui Pan. To-wards data-centric graph machine learning: Review and outlook. arXiv preprint arXiv:2309.10979, 2023a.
- Xin Zheng, Miao Zhang, Chunyang Chen, Soheila Molaei, Chuan Zhou, and Shirui Pan. Gnnevaluator: Evaluating gnn performance on unseen graphs without labels. In Advances in Neural Information Processing Systems (NeurIPS), 2023b.
- Xin Zheng, Miao Zhang, Chunyang Chen, Quoc Viet Hung Nguyen, Xingquan Zhu, and Shirui Pan. Structure-free graph condensation: From large-scale graphs to condensed graph-free data. In Advances in Neural Information Processing Systems (NeurIPS), 2023c.
- Yizhen Zheng, Shirui Pan, Vincent Cs Lee, Yu Zheng, and Philip S Yu. Rethinking and scaling up graph contrastive learning: An extremely efficient approach with group discrimination. In Advances in Neural Information Processing Systems (NeurIPS), 2022c.
 - Qi Zhu, Natalia Ponomareva, Jiawei Han, and Bryan Perozzi. Shift-robust gnns: Overcoming the limitations of localized graph training data. In Advances in Neural Information Processing Systems (NeurIPS), volume 34, pp. 27965–27977, 2021.

	10	

702 APPENDIX 703

708 709

710 711

712

713 714

715

716

717

718

719

720

721

722 723 724

725

726

727

728

729

730

731

732

733 734

735

736

737 738 739

740 741

742

743 744

745

746

747 748

704 This is the appendix of our work: Test-Time Graph Rebirth: Serving GNN Generalization 705 **Under Distribution Shifts.** In this appendix, we provide more details of the proposed TT-GREB in 706 terms of more experiments, covering dataset statistics, baseline method comparison, and additional experimental results.

BASELINE METHOD COMPARISON А

The statistics of datasets are presented in Table A1. In the following, we demonstrate the differences among these baselines:

- Except for TENT, GTRANS, and our proposed TT-GREB, other baseline methods do NOT perform test-time adaption only with a single-stage training process.
- TENT, GTRANS, and our proposed TT-GREB use two-stage training and test-time adaption, where all the GNN backbones with fixed optimal parameters are trained on common cross-entropy loss under the standard training.
- TENT falls into the model-centric method group by fine-tuning and adapting well-trained GNN models' parameters at the test time, while GTRANS, and our proposed TT-GREB do NOT fine-tune the model parameters but only modify graph data at the test time.



Figure A1: Ablation study results (%) on Cora on GPR model demonstrated with Boxplot on all test graphs.



Figure A2: Hyper-parameter λ_{Env} sensitivity study results (ACC%) on Cora with GCN model.

В ADDITIONAL EXPERIMENT RESULTS

We provide more ablation study results covering GPR-GNN model in Fig. A1. Additional hyperparameter sensitivity analysis results are presented in Fig. A2, where λ_{Env} indicates λ_4 , corresponding to \mathbf{W}_{edge}^{Pro} in Eq. (7) and Eq. (8) of the main manuscript.

Note that the outcomes for each hyper-parameter are presented under the condition that the remaining parameters are set to their optimal values. Thus, the optimal set of hyper-parameters is achieved by combining the best values from these analyses.

С VISUALIZATION COMPARISON

749 750

For a comprehensive understanding of the reborn test graph by our proposed method, in Fig. A3, we 751 present the t-SNE visualization of the original test graph and our reborn test graph on Cora's first test 752 graph, in terms of the output node representations of the well-trained GCN model. 753

It can be seen that the clusters corresponding to different node class labels are more distinctly 754 separated in the t-SNE latent space after experiencing our proposed test-graph rebirth process. This 755 could effectively verify that the proposed TT-GREB can be beneficial to improve node representation

Distribution shifts	Datasets	#Nodes	#Edges	#Classes	Metrics	Splits
Node feature shifts	Cora (Yang et al., 2016)	2,703	5,278	10	Accuracy	1/1/8
	Amazon-Photo (Shchur et al., 2018)	7,650	119,081	10	Accuracy	1/1/8
Domain shifts	Twitch-E (Rozemberczki et al., 2021)	1,9129,498	31,299 - 153,138	2	ROC-AUC	1/1/5
Temporal shifts	Elliptic (Pareja et al., 2020)	203,769	234,355	2	F1 Score	5/5/33
	OGB-arxiv (Pareja et al., 2020)	169,343	1,166,243	40	Accuracy	1/1/3

Table A1: Dataset statistics with various test-time graph data distribution shifts. 'Splits' denotes the number of training/validation/test graphs.



Figure A3: Visualization comparison of t-SNE on the embeddings of the original test graph (1-st test graph) and our reborn test graph with the well-trained GCN model on Cora.

learning, and better separation in the latent space test data demonstrates good generalization ability of our method under graph distribution shifts.