

ZERO-SHOT GENERALIST GRAPH ANOMALY DETECTION WITH UNIFIED NEIGHBORHOOD PROMPTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Graph anomaly detection (GAD), which aims to identify nodes in a graph that significantly deviate from normal patterns, plays a crucial role in broad application domains. Existing GAD methods, whether supervised or unsupervised, are one-model-for-one-dataset approaches, *i.e.*, training a separate model for each graph dataset. This limits their applicability in real-world scenarios where training on the target graph data is not possible due to issues like data privacy. To overcome this limitation, we propose a novel zero-shot generalist GAD approach **UNPrompt** that trains a one-for-all detection model, requiring the training of one GAD model on a single graph dataset and then effectively generalizing to detect anomalies in other graph datasets without any retraining or fine-tuning. The key insight in UNPrompt is that i) the predictability of latent node attributes can serve as a generalized anomaly measure and ii) highly generalized normal and abnormal graph patterns can be learned via latent node attribute prediction in a properly normalized node attribute space. UNPrompt achieves generalist GAD through two main modules: one module aligns the dimensionality and semantics of node attributes across different graphs via coordinate-wise normalization in a projected space, while another module learns generalized neighborhood prompts that support the use of latent node attribute predictability as an anomaly score across different datasets. Extensive experiments on real-world GAD datasets show that UNPrompt significantly outperforms diverse competing methods under the generalist GAD setting, and it also has strong superiority under the one-model-for-one-dataset setting.

1 INTRODUCTION

Graph anomaly detection (GAD) aims to identify anomalous nodes that exhibit significant deviations from the majority of nodes in a graph. GAD has attracted extensive research attention in recent years (Ma et al., 2021; Pang et al., 2021; Qiao et al., 2024) due to the board applications in various domains such as spam review detection in online shopping networks (McAuley & Leskovec, 2013; Rayana & Akoglu, 2015) and malicious user detection in social networks (Yang et al., 2019). To handle high-dimensional node attributes and complex structural relations between nodes, graph neural networks (GNNs) (Kipf & Welling, 2016; Wu et al., 2020) have been widely exploited for GAD due to their strong ability to integrate the node attributes and graph structures. These methods can be roughly divided into two categories, *i.e.*, supervised and unsupervised methods. One category formulates GAD as a binary classification problem and aims to capture anomaly patterns under the guidance of labels (Tang et al., 2022; Peng et al., 2018; Gao et al., 2023). By contrast, due to the difficulty of obtaining these class labels, another category of methods takes the unsupervised approach that aims to learn normal graph patterns, *e.g.*, via data reconstruction or other proxy learning tasks that are related to GAD (Qiao & Pang, 2023; Liu et al., 2021b; Ding et al., 2019; Huang et al., 2022).

Despite their remarkable detection performance, these methods need to train a dataset-specific model for each graph dataset for GAD. This one-model-for-one-dataset paradigm limits their applicability in real-world scenarios since training a model from scratch incurs significant computation costs and requires even a large amount of labeled data for supervised GAD methods (Liu et al., 2024; Qiao et al., 2024). Training on a target graph may even not be possible due to data privacy protection and regulation. To address this limitation, a new one-for-all anomaly detection (AD) paradigm, called generalist anomaly detection (Zhu & Pang, 2024; Zhou et al., 2024), has been proposed for image AD with the emergence of foundation models such as CLIP (Radford et al., 2021). This new

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

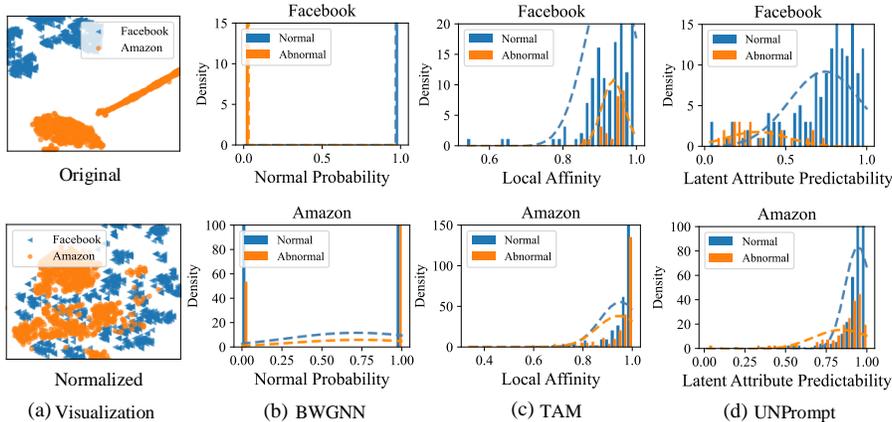


Figure 1: (a) Visualization of two popular GAD datasets: Facebook and Amazon, where the node attributes are unified into a common semantic space via our proposed normalization compared to the original heterogeneous raw attributes. (b)-(d) The anomaly scores of BWGNN (normal probability) (Tang et al., 2022), TAM (local affinity) (Qiao & Pang, 2023) and UNPrompt (latent attribute predictability) on the two datasets, where the methods are all trained on Facebook and tested on Amazon under the zero-shot setting. It is clear that BWGNN and TAM struggle to generalize from Facebook to Amazon, while UNPrompt can learn well to generalize across the datasets.

direction aims to learn a generalist detection model on auxiliary datasets so that it can generalize to detect anomalies effectively in diverse target datasets without any re-training or fine-tuning. This paper explores this direction in the area of GAD.

Compared to image AD, there are some unique challenges for learning generalist models for GAD. First, unlike image data where raw features are in the same RGB space, the node attributes in graphs from different applications and domains can differ significantly in node attribute dimensionality and semantics. For example, as a shopping network dataset, Amazon contains the relationships between users and reviews, and the node attribute dimensionality is 25. Differently, Facebook, a social network dataset, describes relationships between users with 576-dimensional attributes. Second, generalist AD models on image data rely on the superior generalizability learned in large visual-language models (VLMs) through pre-training on web-scale image-text-aligned data (Zhu & Pang, 2024; Zhou et al., 2024), whereas there are no such foundation models for graph data (Liu et al., 2023a). Therefore, the key question here is: *can we learn generalist models for GAD on graph data with heterogeneous node attributes and structure without the support of foundation models?*

To address these challenges, we propose **UNPrompt**, a novel generalist GAD approach that learns *Unified Neighborhood Prompts* on a single auxiliary graph dataset and then effectively generalizes to directly detect anomalies in other graph datasets under a **zero-shot** setting. The key insight in UNPrompt is that i) the predictability of latent node attributes can serve as a generalized anomaly measure and ii) highly generalized normal and abnormal graph patterns can be learned via latent node attribute prediction in a properly normalized node attribute space. UNPrompt achieves this through two main modules including *coordinate-wise normalization-based node attribute unification* and *neighborhood prompt learning*. The former module aligns the dimensionality of node attributes across graphs and transforms the semantics into a common space via coordinate-wise normalization, as shown in Figure 1(a). In this way, the diverse distributions of node attributes are calibrated into the same semantic space. On the other hand, the latter module learns graph-agnostic normal and abnormal patterns via a neighborhood-based latent attribute prediction task. Specifically, we incorporate learnable prompts into the normalized attributes of the neighbors of a target node to predict the latent attributes of the target node. Despite being trained on a small pre-trained GNN using a single graph, UNPrompt can effectively generalize to detect anomalous nodes in different unseen graphs without any re-training at the inference stage, as shown in Figure 1(b)-(d).

Overall, the main contributions of this paper are summarised as follows. (1) We propose a novel zero-shot generalist GAD approach, UNPrompt. To the best of our knowledge, this is the first method that exhibits effective zero-shot GAD performance across various graph datasets. There is a

concurrent work on generalist GAD (Liu et al., 2024), but it can only work under a few-shot setting. (2) We reveal that a simple yet effective coordinate-wise normalization can be utilized to unify the heterogeneous distributions in the node attributes across different graphs. (3) We further introduce a novel neighborhood prompt learning module that utilizes a neighborhood-based latent node attribute prediction task to learn generalized prompts in the normalized attribute space, enabling the zero-shot GAD across different graphs. (4) Extensive experiments on real-world GAD datasets show that UNPrompt significantly outperforms state-of-the-art competing methods under the zero-shot generalist GAD. (5) We show that UNPrompt can also work in the conventional one-model-for-one-dataset setting, outperforming state-of-the-art models in this popular GAD setting.

2 RELATED WORK

Graph Anomaly Detection. Existing GAD methods can be roughly categorized into unsupervised and supervised approaches (Ma et al., 2021; Qiao et al., 2024). The unsupervised methods are typically built using data reconstruction, self-supervised learning, and learnable graph anomaly measures (Qiao et al., 2024; Liu et al., 2022). The reconstruction-based approaches like DOMINANT (Ding et al., 2019) and AnomalyDAE (Fan et al., 2020) aim to capture the normal patterns in the graph, where the reconstruction error in both graph structure and attributes is utilized as the anomaly score. CoLA (Liu et al., 2021b) and SL-GAD (Zheng et al., 2021) are representative self-supervised learning methods assuming that normality is reflected in the relationship between the target node and its contextual nodes. The graph anomaly measure methods typically leverage the graph structure-aware anomaly measures to learn intrinsic normal patterns for GAD, such as node affinity in TAM (Qiao & Pang, 2023). In contrast to the unsupervised approaches, the supervised anomaly detection approaches have shown substantially better detection performance in recent years due to the incorporation of labeled anomaly data (Liu et al., 2021a; Chai et al., 2022). Most supervised methods concentrate on the design of propagation mechanisms and spectral feature transformations to address the notorious over-smoothing feature representation issues (Tang et al., 2022; Gao et al., 2023; Chai et al., 2022). Although both approaches can be adapted for zero-shot GAD by directly applying the trained GAD models to the target datasets, they struggle to capture generalized normal and abnormal patterns for GAD across different graph datasets. There are some studies working on cross-domain GAD (Ding et al., 2021b; Wang et al., 2023) that aim to transfer knowledge from a labeled graph dataset for GAD on a target dataset, but it is a fundamentally different problem from generalist GAD since cross-domain GAD requires the training on both source and target graph datasets.

Graph Prompt Learning. Prompt learning, initially developed in natural language processing, seeks to adapt large-scale pre-trained models to different downstream tasks by incorporating learnable prompts while keeping the pre-trained models frozen (Liu et al., 2023b). Specifically, it designs task-specific prompts capturing the knowledge of the corresponding tasks and enhances the compatibility between inputs and pre-trained models to enhance the pre-trained models in downstream tasks. Recently, prompt learning has been explored in graphs to unify multiple graph tasks (Sun et al., 2023; Liu et al., 2023c) or improve the transferability of graph models on the datasets across the different domains (Li et al., 2024; Zhao et al., 2024), *e.g.*, by optimizing the prompts with labeled data of various downstream tasks (Fang et al., 2024; Liu et al., 2023c). Although being effective in popular graph learning tasks like node classification and link prediction, they are inapplicable to generalist GAD due to the unsupervised nature and/or irregular distributions of anomalies.

Generalist Anomaly Detection. Generalist AD has been very recently emerging as a promising solution to tackle sample efficiency and model generalization problems in AD. There have been a few studies on non-graph data that have large pre-trained models to support the generalized pattern learning, such as image generalist AD (Zhou et al., 2023; Zhu & Pang, 2024). However, it is a very challenging task for data like graph data due to the lack of such pre-trained models. Recently a concurrent approach, ARC (Liu et al., 2024), introduces an effective framework that leverages in-context learning to achieve generalist GAD without relying on large pre-trained GNNs. Unlike ARC which focuses on a few-shot GAD setting, *i.e.*, requiring the availability of some labeled nodes in the target testing graph dataset, we tackle a zero-shot GAD setting assuming no access to any labeled data during inference stages.

Inductive Graph Learning. Similar to generalist setting, inductive graph learning (Hamilton et al., 2017; Ding et al., 2021a; Li et al., 2023b; Huang et al., 2023; Fang et al., 2023) also focuses

on inference on unseen graph data. However, these methods are not applicable to the generalist setting. Specifically, inductive graph learning trains the model on partial data of the whole graph dataset Hamilton et al. (2017); Ding et al. (2019); Li et al. (2023b) or the previously observed data of dynamic graphs (Fang et al., 2023). Then, the learned model is evaluated on the unseen data of the whole dataset or the future graph. These unseen testing data are from the same source of the training data with the same dimensionality and semantics. In contrast, the unseen data in our method are from different distributions/domains with significantly different dimensionality and semantics. This cross-dataset nature, specifically referred to as a zero-shot problem (Jeong et al., 2023; Zhou et al., 2024), makes our setting significantly different from the current inductive graph learning setting.

3 METHODOLOGY

3.1 PRELIMINARIES

Notations. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be an attributed graph with N nodes, where $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$ represents the node set and \mathcal{E} is the edge set. The attributes of nodes can be denoted as $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times d}$ and the edges between nodes can be presented by an adjacency matrix $A \in \{0, 1\}^{N \times N}$ with $A_{ij} = 1$ if there is an edge between v_i and v_j and $A_{ij} = 0$ otherwise. For simplicity, the graph can be represented as $\mathcal{G} = (A, X)$. In GAD, the node set can be divided into a set of the normal nodes \mathcal{V}_n and a set of anomalous nodes \mathcal{V}_a . Typically, the number of normal nodes is significantly larger than the anomalous nodes, *i.e.*, $|\mathcal{V}_n| \gg |\mathcal{V}_a|$. Moreover, the anomaly labels can be denoted as $\mathbf{y} \in \{0, 1\}^N$ with $\mathbf{y}_i = 1$ if $v_i \in \mathcal{V}_a$ and $\mathbf{y}_i = 0$ otherwise.

Conventional GAD. Conventional GAD typically focuses on model training and anomaly detection on the same graph. Specifically, given a graph \mathcal{G} , an anomaly scoring model $f : \mathcal{G} \rightarrow \mathbb{R}$ is optimized on \mathcal{G} in a supervised or unsupervised manner. Then, the model is used to detect anomalies within the same graph. The model is expected to generate higher anomaly scores for abnormal nodes than normal nodes, *i.e.*, $f(v_i) < f(v_j)$ if $v_i \in \mathcal{V}_n$ and $v_j \in \mathcal{V}_a$.

Generalist GAD. Generalist GAD aims to learn a generalist model f on a single training graph so that f can be directly adapted to different target graphs across diverse domains without any fine-tuning or re-training. More specifically, the model is optimized on $\mathcal{G}_{\text{train}}$ with the corresponding anomaly labels $\mathbf{y}_{\text{train}}$. After model optimization, the learned f is utilized to detect anomalies within different unseen target graphs $\mathcal{T}_{\text{test}} = \{\mathcal{G}_{\text{test}}^{(1)}, \dots, \mathcal{G}_{\text{test}}^{(n)}\}$ which has heterogeneous attributes and/or structure to $\mathcal{G}_{\text{train}}$, *i.e.*, $\mathcal{G}_{\text{train}} \cap \mathcal{T}_{\text{test}} = \emptyset$. Depending on whether labeled nodes of the target graph are provided during inference, the generalist GAD problem can be further divided into two categories, *i.e.*, **few-shot** and **zero-shot** settings. We focus on the zero-shot setting where the generalist models cannot get access to any labeled data of the testing graphs during both training and inference.

3.2 OVERVIEW OF THE PROPOSED APPROACH – UNPROMPT

The framework is illustrated in Figure 2, which consists of two main modules, coordinate-wise normalization-based node attribute unification and neighborhood prompt learning. For all graphs, the node attribute unification aligns the dimensionality of node attributes and transforms the semantics into a common space via coordinate-wise normalization in a projected space. Then, in the normalized space, the generalized latent attribute prediction task is performed with the neighborhood prompts to learn generalized GAD patterns at the training stage. In this prompt learning module, UNPrompt aims to maximize the predictability of the latent attributes of normal nodes while minimizing those of abnormal nodes. In this paper, we evaluate the predictability via the similarity. In doing so, the graph-agnostic normal and abnormal patterns are incorporated into the prompts. During inference, the target graph is directly fed into the learned models after node attribute unification without any re-training or labeled nodes of the graph. For each node, the predictability of latent node attributes is directly used as the normal score for final anomaly detection.

3.3 NODE ATTRIBUTE UNIFICATION

Graphs from different distributions and domains significantly differ in the dimensionality and semantics of node attributes. Therefore, the premise of developing a generalist GAD model is to unify

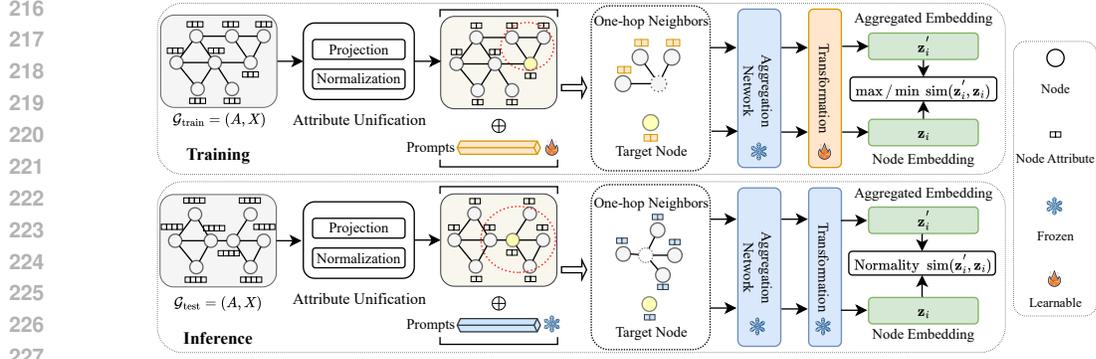


Figure 2: Overview of UNPrompt. Node attribute unification is used to align the attribute dimensionality and semantics. During training, the neighborhood prompts are optimized to capture generalized patterns by maximizing the predictability of the latent attributes (*i.e.*, the embedding z_i) of normal nodes while minimizing that of abnormal nodes. During inference, the learned prompts are directly applied to the testing nodes, and the latent attribute predictability of each node is used for GAD.

the dimensionality and semantics of node attributes into the same space. In this paper, we propose a simple yet effective node attribute unification module to address this issue, which consists of feature projection and coordinate-wise normalization. Different from ARC (Liu et al., 2024) which aligns the attributes based on feature reordering using feature smoothness, we calibrate the feature distributions of diverse graphs into the same frame, resulting in a simpler yet effective alignment.

Feature Projection. To address the inconsistent attribute dimensions across graphs, various feature projection methods can be utilized, such as singular value decomposition (Stewart, 1993) (SVD) and principal component analysis (Abdi & Williams, 2010) (PCA). Formally, given the attribute matrix $X^{(i)} \in \mathbb{R}^{N^{(i)} \times d^{(i)}}$ of any graph $\mathcal{G}^{(i)}$ from $\mathcal{G}_{\text{train}} \cup \mathcal{T}_{\text{test}}$, we transform it into $\tilde{X}^{(i)} \in \mathbb{R}^{N^{(i)} \times d'}$ with the common dimensionality of d' ,

$$X^{(i)} \in \mathbb{R}^{N^{(i)} \times d^{(i)}} \xrightarrow[\text{Projection}]{\text{Feature}} \tilde{X}^{(i)} \in \mathbb{R}^{N^{(i)} \times d'}. \quad (1)$$

Coordinate-wise normalization. Despite the attribute dimensionality being unified, the semantics and distributions of each attribute dimension are still divergent across graphs, posing significant challenges to learning a generalist GAD model. A recent study (Li et al., 2023a) has demonstrated that semantic differences across datasets are mainly reflected in the distribution shifts and calibrating the distributions into a common frame helps learn more generalized AD models. Inspired by this, we propose to use coordinate-wise normalization to align the semantics and unify the distributions across graphs. Specifically, the transformed attribute matrix $\tilde{X}^{(i)}$ is shifted and re-scaled to have mean zeros and variance ones via the following equation:

$$\bar{X}^{(i)} = \frac{\tilde{X}^{(i)} - \boldsymbol{\mu}^{(i)}}{\boldsymbol{\sigma}^{(i)}}, \quad (2)$$

where $\boldsymbol{\mu}^{(i)} = [\mu_1^{(i)}, \dots, \mu_{d'}^{(i)}]$ and $\boldsymbol{\sigma}^{(i)} = [\sigma_1^{(i)}, \dots, \sigma_{d'}^{(i)}]$ are the coordinate-wise mean and variance of $\tilde{X}^{(i)}$ of the graph $\mathcal{G}^{(i)}$. In this way, the distributions of normalized attributes along each dimension are the same within and across graphs, as shown in Figure 1(a). This helps to capture the generalized normal and abnormal patterns for generalist GAD (see Table 2).

3.4 NEIGHBORHOOD PROMPT LEARNING VIA LATENT NODE ATTRIBUTE PREDICTION

Latent Node Attribute Predictability as Anomaly Score. To build a generalist GAD model, one must capture the generalized normal and abnormal patterns across graphs. Otherwise, the model would overfit the dataset-specific knowledge of the training graph which can be very different from that in target graphs. In this paper, we reveal that the predictability of latent node attributes can serve as a generalized anomaly measure, and thus, highly generalized normal and abnormal graph

270 patterns can be learned via latent node attribute prediction in the normalized node attribute space
 271 with the neighborhood prompts. The key intuition of this anomaly measure is that normal nodes
 272 tend to have more connections with normal nodes of similar attributes due to prevalent graph ho-
 273 mophily relations, resulting in a more homogeneous neighborhood in the normal nodes (Qiao &
 274 Pang, 2023); by contrast, the presence of anomalous connections and/or attributes makes abnormal
 275 nodes deviate significantly from their neighbors. Therefore, for a target node, its latent attributes
 276 (*i.e.*, node embedding) is more predictable based on the latent attributes of its neighbors if the node
 277 is normal node, compared to abnormal nodes. The neighborhood-based latent attribute prediction is
 278 thus used to measure the normality for GAD. As shown in our experiments (see Figures 1(b)-(d) and
 279 Tables 1 and 3), it is a generalized anomaly scoring method that works effectively across graphs.
 280 However, due to the existence of irrelevant and noisy attribute information in the original attribute
 281 space, the attribute prediction is not as effective as expected in the simply projected space after at-
 282 tribute unification. To address this issue, we propose to learn discriminative prompts via the latent
 283 attribute prediction task to enhance the effectiveness of this anomaly measure.

284 To achieve this, we first design a simple graph neural network g , a neighborhood aggregation net-
 285 work, to generate the aggregated neighborhood embedding of each target node. Specifically, given a
 286 graph $\mathcal{G} = (A, \bar{X})$, the aggregated neighborhood embeddings for each node are obtained as follows:

$$287 \quad \tilde{Z} = g(\mathcal{G}) = \tilde{A}\bar{X}W, \quad (3)$$

288 where \tilde{Z} is the aggregated representation of neighbors, $\tilde{A} = (D)^{-1}A$ is the normalized adjacency
 289 matrix with D being a diagonal matrix and its elements $D_{kk} = \sum_j A_{kj}$, and W is the learnable
 290 parameters. Compared to conventional GNNs such as GCN (Kipf & Welling, 2016) and SGC (Wu
 291 et al., 2019), we do not require \tilde{A} to be self-looped and symmetrically normalized as we aim to
 292 obtain the aggregated representation of all the neighbors for each node. To design the latent node
 293 attribute prediction task, we further obtain the latent attributes of each node as follows:

$$294 \quad Z = \bar{X}W, \quad (4)$$

295 where Z serves as the prediction ground truth for the latent attribute prediction task. The adjacency
 296 matrix A is discarded to avoid carrying neighborhood-based attribute information into Z which
 297 would lead to ground truth leakage in this prediction task. We further propose to utilize the cosine
 298 similarity to measure this neighborhood-based latent attribute predictability for each node:

$$299 \quad s_i = \text{sim}(\mathbf{z}_i, \tilde{\mathbf{z}}_i) = \frac{\mathbf{z}_i(\tilde{\mathbf{z}}_i)^T}{\|\mathbf{z}_i\| \|\tilde{\mathbf{z}}_i\|}, \quad (5)$$

300 where \mathbf{z} and $\tilde{\mathbf{z}}_i$ are the i -th node embeddings in Z and \tilde{Z} respectively. A higher similarity denotes
 301 the target node can be well predicted by its neighbors and indicates the target is normal with a higher
 302 probability. Therefore, we directly utilize the similarity to measure the normal score of the nodes.
 303

304 **GNN Pre-training.** To build generalist models, pre-training is required. Here we pre-train the
 305 above neighborhood aggregation network via graph contrastive learning due to the ability to obtain
 306 robust and transferable models (You et al., 2020; Zhu et al., 2020) across graphs (see Appendix B for
 307 the details). Without pre-training, the dataset-specific knowledge would be captured by the model
 308 if it is directly optimized based on the neighborhood-based latent attribute prediction of normal and
 309 abnormal nodes, limiting the generalizability of the model to other graphs (see Table 2).
 310

311 **Neighborhood Prompting via Latent Attribute Prediction.** After the pre-training, we aim to
 312 further learn more generalized normal and abnormal patterns via prompt tuning in the normalized
 313 space. Thus, we devise learnable prompts appending to the attributes of the neighboring nodes of
 314 the target nodes, namely *neighborhood prompts*, for learning robust and discriminative patterns that
 315 can detect anomalous nodes in different unseen graphs without any re-training during inference.
 316

317 Specifically, neighborhood prompting aims to learn some prompt tokens that help maximize the
 318 neighborhood-based latent prediction of normal nodes while minimizing that of abnormal nodes si-
 319 multaneously. To this end, the prompt is designed as a set of shared and learnable tokens that can
 320 be incorporated into the normalized node attributes. Formally, the neighborhood prompts are repre-
 321 sented as $P = [\mathbf{p}_1, \dots, \mathbf{p}_k]^T \in \mathbb{R}^{K \times d'}$ where K is the number of vector-based tokens \mathbf{p}_i . For each
 322
 323

node in $\mathcal{G} = (A, \bar{X})$, the node attributes in the unified feature space are augmented by the weighted combination of these tokens, with the weights obtained from K learnable linear projections:

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}}_i + \sum_j^K \alpha_j \mathbf{p}_j, \quad \alpha_j = \frac{e^{(\mathbf{w}_j)^T \mathbf{x}_i^t}}{\sum_l^K e^{(\mathbf{w}_l)^T \mathbf{x}_i^t}}, \quad (6)$$

where α_j denotes the importance score of the token \mathbf{p}_j in the prompt and \mathbf{w}_j is a learnable projection. For convenience, we denote the graph modified by the graph prompt as $\tilde{\mathcal{G}} = (A, \bar{X} + P)$. Then, $\tilde{\mathcal{G}}$ is fed into the frozen pre-trained model g to obtain the corresponding aggregated embeddings \tilde{Z} and node latent attributes Z via Eq.(3) and Eq.(4) respectively to measure the attribute predictability. To further enhance the representation discrimination, a transformation layer h is applied on the learned \tilde{Z} and Z to transform them into a more anomaly-discriminative feature space,

$$\tilde{Z} = h(\tilde{Z}), \quad Z = h(Z). \quad (7)$$

The transformed representations are then used to measure the latent node attribute predictability with Eq.(5). To optimize P and h , we employ the following training objective,

$$\min_{P, h} \sum \ell(\mathbf{z}_i, \tilde{\mathbf{z}}_i), \quad (8)$$

where $\ell(\mathbf{z}_i, \tilde{\mathbf{z}}_i) = -\text{sim}(\mathbf{z}_i, \tilde{\mathbf{z}}_i)$ if $\mathbf{y}_i = 0$, and $\ell(\mathbf{z}_i, \tilde{\mathbf{z}}_i) = \text{sim}(\mathbf{z}_i, \tilde{\mathbf{z}}_i)$ if $\mathbf{y}_i = 1$.

3.5 TRAINING AND INFERENCE OF UNPROMPT

Training. The training process of UNPrompt can be divided into two parts. First, given $\mathcal{G}_{\text{train}}$, a neighborhood aggregation network g is optimized via graph contrastive learning. Then, the neighborhood prompts P and the transformation layer h are optimized to capture the graph-agnostic normal and abnormal patterns while keeping the pre-trained model g frozen. In this way, the transferable knowledge of the pre-trained g is maintained, while the neighborhood prompt learning helps learn the generalized normal and abnormal patterns.

Inference. During inference, given $\mathcal{G}_{\text{test}}^{(i)} \in \mathcal{T}_{\text{test}}$, the node attributes are first aligned. Then, the test graph $\mathcal{G}_{\text{test}}^{(i)}$ is augmented with the learned neighborhood prompt P and fed into the model g and the transformation layer h to obtain the neighborhood aggregated representations and the latent node attributes. Finally, the similarity (Eq.(5)) is used as the normal score for the test nodes for anomaly detection. Note that the inference does not require any further re-training and labeled nodes of $\mathcal{G}_{\text{test}}^{(i)}$. The algorithms of the training and inference of UNPrompt are provided in Appendix C.

4 EXPERIMENTS

4.1 PERFORMANCE ON ZERO-SHOT GENERALIST GAD

Datasets. We evaluate the proposed UNPrompt on seven real-world GAD datasets from diverse social networks and online shopping co-review networks. Specifically, the social networks include Facebook (Xu et al., 2022), Reddit (Kumar et al., 2019) and Weibo (Kumar et al., 2019). The co-review networks consist of Amazon (McAuley & Leskovec, 2013), YelpChi (Rayana & Akoglu, 2015), Amazon-all (McAuley & Leskovec, 2013) and YelpChi-all (Rayana & Akoglu, 2015).

Competing Methods. Since there is no zero-shot generalist GAD method, a set of eight state-of-the-art (SotA) unsupervised and supervised competing methods are employed for comparison in our experiments. The unsupervised methods comprise reconstruction-based AnomalyDAE (Fan et al., 2020), contrastive learning-based CoLA (Liu et al., 2021b), hop prediction-based HCM-A (Huang et al., 2022), local affinity-based TAM (Qiao & Pang, 2023) and GADAM (Chen et al., 2024). Supervised methods include two conventional GNNs – GCN (Kipf & Welling, 2016) and GAT (Veličković et al., 2017) – and three SotA GAD GNNs – BWGNN (Tang et al., 2022), GHRN (Gao et al., 2023) and XGBGraph (Tang et al., 2023).

Following (Liu et al., 2024; Qiao & Pang, 2023; Qiao et al., 2024), two widely-used metrics, AU-ROC and AUPRC, are used to evaluate the performance of all methods. For both metrics, the higher value denotes the better performance. Moreover, for each method, we report the average performance with standard deviations after 5 independent runs with different random seeds.

Table 1: AUROC and AUPRC results on six real-world GAD datasets with the models trained on Facebook only. For each dataset, the best performance per column within each metric is boldfaced, with the second-best underlined. ‘‘Avg’’ denotes the averaged performance of each method.

Metric	Method	Dataset						Avg.
		Amazon	Reddit	Weibo	YelpChi	Aamazon-all	YelpChi-all	
AUROC	Unsupervised Methods							
	AnomalyDAE	0.5818 \pm 0.039	0.5016 \pm 0.032	<u>0.7785</u> \pm 0.058	0.4837 \pm 0.094	0.7228 \pm 0.023	0.5002 \pm 0.018	<u>0.5948</u>
	CoLA	0.4580 \pm 0.054	0.4623 \pm 0.005	0.3924 \pm 0.041	0.4907 \pm 0.017	0.4091 \pm 0.052	0.4879 \pm 0.010	0.4501
	HCM-A	0.4784 \pm 0.005	0.5387 \pm 0.041	0.5782 \pm 0.048	0.5000 \pm 0.000	0.5056 \pm 0.059	0.5023 \pm 0.005	0.5172
	GADAM	<u>0.6646</u> \pm 0.063	<u>0.4532</u> \pm 0.024	<u>0.3652</u> \pm 0.052	<u>0.3376</u> \pm 0.012	<u>0.5959</u> \pm 0.080	<u>0.4829</u> \pm 0.016	<u>0.4832</u>
	TAM	0.4720 \pm 0.005	0.5725 \pm 0.004	0.4867 \pm 0.028	0.5035 \pm 0.014	0.7543 \pm 0.002	0.4216 \pm 0.002	0.5351
	Supervised Methods							
	GCN	0.5988 \pm 0.016	0.5645 \pm 0.000	0.2232 \pm 0.074	0.5366 \pm 0.019	0.7195 \pm 0.002	<u>0.5486</u> \pm 0.001	0.5319
	GAT	0.4981 \pm 0.008	0.5000 \pm 0.025	0.4521 \pm 0.101	<u>0.5871</u> \pm 0.016	0.5005 \pm 0.012	0.4802 \pm 0.004	0.5030
	BWGN	0.4769 \pm 0.020	0.5208 \pm 0.016	0.4815 \pm 0.108	0.5538 \pm 0.027	0.3648 \pm 0.050	0.5282 \pm 0.015	0.4877
	GHRN	0.4560 \pm 0.033	0.5253 \pm 0.006	0.5318 \pm 0.038	0.5524 \pm 0.020	0.3382 \pm 0.085	0.5125 \pm 0.016	0.4860
	XGBGraph	<u>0.4179</u> \pm 0.000	<u>0.4601</u> \pm 0.000	<u>0.5373</u> \pm 0.000	<u>0.5722</u> \pm 0.000	<u>0.7950</u> \pm 0.000	<u>0.4945</u> \pm 0.000	<u>0.5462</u>
	UNPrompt (Ours)	0.7525 \pm 0.016	0.5337 \pm 0.002	0.8860 \pm 0.007	0.5875 \pm 0.016	0.7962 \pm 0.022	0.5558 \pm 0.012	0.6853
	AUPRC	Unsupervised Methods						
AnomalyDAE		0.0833 \pm 0.015	0.0327 \pm 0.004	<u>0.6064</u> \pm 0.031	0.0624 \pm 0.017	0.1921 \pm 0.026	0.1484 \pm 0.009	<u>0.1876</u>
CoLA		0.0669 \pm 0.002	0.0391 \pm 0.004	0.1189 \pm 0.014	0.0511 \pm 0.000	0.0861 \pm 0.019	0.1466 \pm 0.003	0.0848
HCM-A		0.0669 \pm 0.002	0.0391 \pm 0.004	0.1189 \pm 0.014	0.0511 \pm 0.000	0.0861 \pm 0.019	0.1466 \pm 0.003	0.0848
GADAM		<u>0.1562</u> \pm 0.103	<u>0.0293</u> \pm 0.001	<u>0.0830</u> \pm 0.005	<u>0.0352</u> \pm 0.001	<u>0.1595</u> \pm 0.121	<u>0.1371</u> \pm 0.006	<u>0.1001</u>
TAM		0.0666 \pm 0.001	<u>0.0413</u> \pm 0.001	0.1240 \pm 0.014	0.0524 \pm 0.002	0.1736 \pm 0.004	0.1240 \pm 0.001	0.0970
Supervised Methods								
GCN		<u>0.0891</u> \pm 0.007	0.0439 \pm 0.000	0.1109 \pm 0.020	0.0648 \pm 0.009	0.1536 \pm 0.002	<u>0.1735</u> \pm 0.000	0.1060
GAT		0.0688 \pm 0.002	0.0329 \pm 0.002	0.1009 \pm 0.017	0.0810 \pm 0.005	0.0696 \pm 0.001	0.1400 \pm 0.002	0.0822
BWGN		0.0652 \pm 0.002	0.0389 \pm 0.003	0.2241 \pm 0.046	0.0708 \pm 0.018	0.0586 \pm 0.003	0.1605 \pm 0.005	0.1030
GHRN		0.0633 \pm 0.003	0.0407 \pm 0.002	0.1965 \pm 0.059	0.0661 \pm 0.010	0.0569 \pm 0.006	0.1505 \pm 0.005	0.0957
XGBGraph		<u>0.0536</u> \pm 0.000	<u>0.0330</u> \pm 0.000	<u>0.2256</u> \pm 0.000	<u>0.0655</u> \pm 0.000	<u>0.2307</u> \pm 0.000	<u>0.1449</u> \pm 0.000	<u>0.1256</u>
UNPrompt (Ours)		0.1602 \pm 0.013	0.0351 \pm 0.000	0.6406 \pm 0.026	<u>0.0712</u> \pm 0.008	0.2430 \pm 0.028	0.1810 \pm 0.012	0.2219

Implementation Details. To ensure a fair comparison, the common dimensionality is set to eight to unify the node attribute across graphs for all methods, and SVD is used for feature projection. The number of layers in GNNs is set to one and the number of hidden units is 128. The transformation layer is also implemented via a one-layer MLP with the same number of hidden units. The size of the neighborhood prompt is set to one by default. For all baselines, we adopt their official code and follow the recommended optimization and hyperparameter settings to conduct the experiments. UNPrompt and all its competing methods are trained on Facebook and then directly tested on the other six GAD datasets without any further training or additional knowledge of the target graphs.

Main Results. The AUROC and AUPRC results of all methods are presented in Table 1. From the table, we can have the following observations. (1) Under the proposed generalist GAD scenario where a model is trained on a single dataset and evaluated on six other datasets, all the competing baselines fail to work well, demonstrating that it is very challenging to build a generalist GAD model that generalizes across different datasets under zero-shot setting. (2) For supervised methods, the simple GCN achieves better performance than the specially designed GAD GNNs. This can be attributed to more dataset-specific knowledge being captured in these specialized GAD models, limiting their generalization capacity to the unseen testing graphs. (3) Unsupervised methods perform more stable than supervised methods across the target graphs and generally outperform supervised ones. This is because the unsupervised objectives are closer to the shared anomaly patterns across graphs compared to the supervised ones, especially for TAM which employs a fairly generalized local affinity-based objective to train the model. (4) The proposed method UNPrompt demonstrates strong and stable generalist GAD capacity across graphs from different distributions and domains. Specifically, UNPrompt achieves the best AUROC performance on 5 out of 6 datasets and the average performance outperforms the best-competing method by over 9%. In terms of AUPRC, UNPrompt outperforms all baselines on 4 out of 6 datasets and also achieves the best average performance. The superiority of UNPrompt is attributed to the fact that i) the proposed coordinate-wise normalization effectively aligns the features across graphs, and ii) the shared generalized normal and abnormal patterns are well captured in the neighborhood prompts.

Ablation Study. To evaluate the importance of each component in UNPrompt, we design four variants, *i.e.*, w/o coordinate-wise normalization, w/o graph contrastive learning-based pre-training, without neighborhood prompts, and w/o transformation layer. The results of these variants are re-

Table 2: AUROC results of the proposed method UNPrompt and its four variants.

Method	Amazon	Reddit	Weibo	YelpChi	Aamazon-all	YelpChi-all	Avg.
UNPrompt	0.7525	0.5337	0.8860	0.5875	0.7962	0.5558	0.6853
w/o Normalization	0.4684	0.5006	0.1889	0.5620	0.3993	0.5466	0.4443
w/o Pre-training	0.5400	0.5233	0.5658	0.4672	0.3902	0.4943	0.4968
w/o Prompt	0.5328	0.5500	0.4000	0.4520	0.4096	0.4894	0.4723
w/o Transformation	0.7331	0.5556	0.7406	0.5712	0.7691	0.5545	0.6540

ported in the Table 2. From the table, we can see that all four components contribute to the overall superior performance of UNPrompt. More specifically, (1) without the coordinate-wise normalization, the method fails to calibrate the distributions of diverse node attributes into a common space, leading to large performance drop across all datasets. (2) Besides the semantics alignment, the graph contrastive learning-based pre-training ensures our GNN network is transferable to other graphs instead of overfitting to the training graph. As expected, the performance of the variant without pre-training also drops significantly. (3) If the neighborhood prompts are removed, the learning of latent node attribute prediction is ineffective for capturing generalized normal and abnormal patterns. (4) The variant without the transformation layer achieves inferior performance on nearly all the datasets, demonstrating the importance of mapping the features into a more anomaly-discriminative space.

Sensitivity w.r.t the Neighborhood Prompt Size. We evaluate the sensitivity of UNPrompt w.r.t the size of the neighborhood prompts, *i.e.*, the number of tokens K . We vary K in the range of $[1, 9]$ and report the results in Figure 3(a). It is clear that the performances on Reddit, Weibo and YelpChi-all remain stable with varying sizes of neighborhood prompts while the other datasets show slight fluctuation, demonstrating that the generalized normal and abnormal patterns can be effectively captured in our neighborhood prompts even with a small size.

Prompt learning using latent attribute prediction vs. alternative graph anomaly measures. To further justify the effectiveness of latent attribute predictability on learning generalized GAD patterns in our prompt learning module, we compare this proposed learnable anomaly measure to the recently proposed anomaly measure, local node affinity in TAM (Qiao & Pang, 2023). All modules of UNPrompt are fixed with only the latent attribute prediction task replaced as the maximization of local affinity as in TAM. The results are presented in Figure 3(b). We can see that the latent attribute predictability consistently and significantly outperforms the local affinity-based measure across all graphs, demonstrating its superiority in learning generalized patterns for generalist GAD.

4.2 PERFORMANCE ON CONVENTIONAL UNSUPERVISED GAD

We also evaluate the applicability of UNPrompt unsupervised GAD setting to further verify the effectiveness of the latent node attribute prediction-based anomaly scores using our proposed neighborhood prompt learning. Specifically, we adopt the same pipeline as in the generalist GAD setting, *i.e.*, graph contrastive-based pre-training and neighborhood prompt learning. Different from the training process in the generalist setting, there is no label information available in unsupervised GAD since models are trained and evaluated on the same graph data. To address this issue, we employ the pseudo-labeling technique to provide supervision for neighborhood prompt learning. In a nutshell, we enforce the neighborhood prompts to maximize the latent attribute predictability of high-score nodes. More details on unsupervised GAD are provided in Appendix D.

Experimental Setup. Six datasets from different distributions and domains are used, *i.e.*, Amazon, Facebook, Reddit, YelpCHI, Amazon-all, and YelpChi-all. Following (Qiao & Pang, 2023), eight SotA unsupervised baselines are used for comparison, *i.e.*, iForest (Liu et al., 2012), ANOMALOUS (Peng et al., 2018), CoLA (Liu et al., 2021b), SL-GAD (Zheng et al., 2021), HCM-A (Huang et al., 2022), DOMINANT (Ding et al., 2019), ComGA (Luo et al., 2022) and TAM (Qiao & Pang, 2023). For each method, we report the average performance with standard deviations after 5 independent runs with different random seeds. The implementation details of UNPrompt remain the same as in the generalist GAD setting. More experimental details on unsupervised GAD are in Appendix F.2.

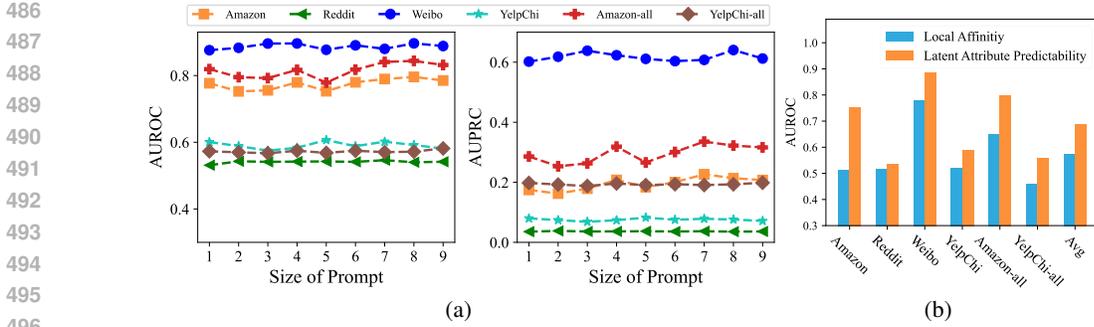


Figure 3: (a) AUROC and AUPRC results of UNPrompt w.r.t. varying neighborhood prompt size. (b). The AUROC performance of generalist GAD with different prompt learning objectives.

Table 3: AUROC and AUPRC results of unsupervised GAD methods on six real-world GAD datasets. The best performance per column within each metric is boldfaced, with the second-best underlined. ‘‘Avg’’ denotes the averaged performance of each method.

Metric	Method	Dataset						Avg.
		Amazon	Facebook	Reddit	YelpChi	Amazon-all	YelpChi-all	
AUROC	iForest	0.5621 \pm 0.008	0.5382 \pm 0.015	0.4363 \pm 0.020	0.4120 \pm 0.040	0.1914 \pm 0.002	0.3617 \pm 0.001	0.4169
	ANOMALOUS	0.4457 \pm 0.003	0.9021 \pm 0.005	0.5387 \pm 0.012	0.4956 \pm 0.003	0.3230 \pm 0.021	0.3474 \pm 0.018	0.5087
	DOMINANT	0.5996 \pm 0.004	0.5677 \pm 0.002	0.5555 \pm 0.011	0.4133 \pm 0.010	0.6937 \pm 0.028	0.5390 \pm 0.014	0.5615
	CoLA	0.5898 \pm 0.008	0.8434 \pm 0.011	<u>0.6028</u> \pm 0.007	0.4636 \pm 0.001	0.2614 \pm 0.021	0.4801 \pm 0.016	0.5402
	SL-GAD	0.5937 \pm 0.011	0.7936 \pm 0.005	0.5677 \pm 0.005	0.3312 \pm 0.035	0.2728 \pm 0.012	0.5551 \pm 0.015	0.5190
	HCM-A	0.3959 \pm 0.014	0.7387 \pm 0.032	0.4593 \pm 0.011	0.4593 \pm 0.005	0.4191 \pm 0.011	0.5691 \pm 0.018	0.5069
	ComGA	0.5895 \pm 0.008	0.6055 \pm 0.000	0.5453 \pm 0.003	0.4391 \pm 0.000	0.7154 \pm 0.014	0.5352 \pm 0.006	0.5716
	TAM	<u>0.7064</u> \pm 0.010	<u>0.9144</u> \pm 0.008	0.6023 \pm 0.004	<u>0.5643</u> \pm 0.007	<u>0.8476</u> \pm 0.028	<u>0.5818</u> \pm 0.033	<u>0.7028</u>
	UNPrompt (Ours)	0.7335 \pm 0.020	0.9379 \pm 0.006	0.6067 \pm 0.006	0.6223 \pm 0.007	0.8516 \pm 0.004	0.6084 \pm 0.001	0.7267
AUPRC	iForest	0.1371 \pm 0.002	0.0316 \pm 0.003	0.0269 \pm 0.001	0.0409 \pm 0.000	0.0399 \pm 0.001	0.1092 \pm 0.001	0.0643
	ANOMALOUS	0.0558 \pm 0.001	0.1898 \pm 0.004	0.0375 \pm 0.004	0.0519 \pm 0.002	0.0321 \pm 0.001	0.0361 \pm 0.005	0.0672
	DOMINANT	0.1424 \pm 0.002	0.0314 \pm 0.041	0.0356 \pm 0.002	0.0395 \pm 0.020	0.1015 \pm 0.018	0.1638 \pm 0.007	0.0857
	CoLA	0.0677 \pm 0.001	0.2106 \pm 0.017	<u>0.0449</u> \pm 0.002	0.0448 \pm 0.002	0.0516 \pm 0.001	0.1361 \pm 0.015	0.0926
	SL-GAD	0.0634 \pm 0.005	0.1316 \pm 0.020	0.0406 \pm 0.004	0.0350 \pm 0.000	0.0444 \pm 0.001	0.1711 \pm 0.011	0.0810
	HCM-A	0.0527 \pm 0.015	0.0713 \pm 0.004	0.0287 \pm 0.005	0.0287 \pm 0.012	0.0565 \pm 0.003	0.1154 \pm 0.004	0.0589
	ComGA	0.1153 \pm 0.005	0.0354 \pm 0.001	0.0374 \pm 0.001	0.0423 \pm 0.000	0.1854 \pm 0.003	0.1658 \pm 0.003	0.0969
	TAM	0.2634 \pm 0.008	0.2233 \pm 0.016	0.0446 \pm 0.001	0.0778 \pm 0.009	0.4346 \pm 0.021	0.1886 \pm 0.017	0.2054
	UNPrompt (Ours)	0.2688 \pm 0.060	0.2622 \pm 0.028	0.0450 \pm 0.001	0.0895 \pm 0.004	0.6094 \pm 0.014	0.2068 \pm 0.004	0.2470

Main Results. The AUROC and AUPRC results of all methods are presented in Table 3. Despite being a generalist GAD method, UNPrompt works very well as a specialized GAD model too. UNPrompt substantially outperforms all the competing methods on all datasets in terms of both AUROC and AUPRC. Particularly, the average performance of UNPrompt surpasses the best-competing method TAM by over 2% in both metrics. Moreover, UNPrompt outperforms the best-competing method by 2%-6% in AUROC on most of the datasets. The superior performance shows that the latent node attribute predictability can be a generalized GAD measure that holds for different graphs, and this property can be effectively learned by the proposed neighborhood prompting method.

5 CONCLUSION

In this paper, we propose a novel zero-shot generalist GAD method, UNPrompt, that trains one detector on a single dataset and can effectively generalize to other unseen target graphs without any further re-training or labeled nodes of target graphs during inference. The attribute inconsistency and the absence of generalized anomaly patterns are the main obstacles for generalist GAD. To address these issues, two main modules are proposed, *i.e.*, coordinate-wise normalization-based attribute unification and neighborhood prompt learning. The first module aligns node attribute dimensionality and semantics, while the second module captures generalized normal and abnormal patterns via the neighborhood-based latent node attribute prediction. Extensive experiments on various real-world GAD datasets from different distributions and domains demonstrate the effectiveness of UNPrompt for generalist GAD. Besides, the experiments conducted on the unsupervised GAD with UNPrompt further support the rationality of the learned anomaly patterns in the generalist model.

REFERENCES

- 540
541
542 Hervé Abdi and Lynne J Williams. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics*, 2(4):433–459, 2010.
- 543
544
545 Ziwei Chai, Siqi You, Yang Yang, Shiliang Pu, Jiarong Xu, Haoyang Cai, and Weihao Jiang. Can abnormality be detected by graph neural networks? In *IJCAI*, pp. 1945–1951, 2022.
- 546
547
548 Bo Chen, Jing Zhang, Xiaokang Zhang, Yuxiao Dong, Jian Song, Peng Zhang, Kaibo Xu, Evgeny Kharlamov, and Jie Tang. Gccad: Graph contrastive coding for anomaly detection. *IEEE Transactions on Knowledge and Data Engineering*, 35(8):8037–8051, 2022.
- 549
550
551
552 Jingyan Chen, Guanghui Zhu, Chunfeng Yuan, and Yihua Huang. Boosting graph anomaly detection with adaptive message passing. In *The Twelfth International Conference on Learning Representations*, 2024.
- 553
554
555
556 Kaize Ding, Jundong Li, Rohit Bhanushali, and Huan Liu. Deep anomaly detection on attributed networks. In *Proceedings of the 2019 SIAM international conference on data mining*, pp. 594–602. SIAM, 2019.
- 557
558
559
560 Kaize Ding, Jundong Li, Nitin Agarwal, and Huan Liu. Inductive anomaly detection on attributed networks. In *Proceedings of the twenty-ninth international conference on international joint conferences on artificial intelligence*, pp. 1288–1294, 2021a.
- 561
562
563
564
565
566 Kaize Ding, Kai Shu, Xuan Shan, Jundong Li, and Huan Liu. Cross-domain graph anomaly detection. *IEEE Transactions on Neural Networks and Learning Systems*, 33(6):2406–2415, 2021b.
- 567
568
569
570
571
572 Haoyi Fan, Fengbin Zhang, and Zuoyong Li. Anomalydae: Dual autoencoder for anomaly detection on attributed networks. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5685–5689. IEEE, 2020.
- 573
574
575
576
577
578 Lanting Fang, Kaiyu Feng, Jie Gui, Shanshan Feng, and Aiqun Hu. Anonymous edge representation for inductive anomaly detection in dynamic bipartite graph. *Proceedings of the VLDB Endowment*, 16(5):1154–1167, 2023.
- 579
580
581
582
583
584
585 Taoran Fang, Yunchao Zhang, Yang Yang, Chunping Wang, and Lei Chen. Universal prompt tuning for graph neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- 586
587
588
589
590
591 Yuan Gao, Xiang Wang, Xiangnan He, Zhenguang Liu, Huamin Feng, and Yongdong Zhang. Addressing heterophily in graph anomaly detection: A perspective of graph spectrum. In *Proceedings of the ACM Web Conference 2023*, pp. 1528–1538, 2023.
- 592
593
594
595 Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.
- 596
597
598
599
600
601 Tianjin Huang, Yulong Pei, Vlado Menkovski, and Mykola Pechenizkiy. Hop-count based self-supervised anomaly detection on attributed networks. In *Joint European conference on machine learning and knowledge discovery in databases*, pp. 225–241. Springer, 2022.
- 602
603
604
605
606
607 Yihong Huang, Liping Wang, Fan Zhang, and Xuemin Lin. Unsupervised graph outlier detection: Problem revisit, new insight, and superior method. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*, pp. 2565–2578. IEEE, 2023.
- 608
609
610
611
612
613 Jongheon Jeong, Yang Zou, Taewan Kim, Dongqing Zhang, Avinash Ravichandran, and Onkar Dabeer. Winclip: Zero-/few-shot anomaly classification and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19606–19616, 2023.
- 614
615
616
617
618
619 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 620
621
622
623
624
625 Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory in temporal interaction networks. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 1269–1278, 2019.

- 594 Aodong Li, Chen Qiu, Marius Kloft, Padhraic Smyth, Maja Rudolph, and Stephan Mandt. Zero-shot
595 anomaly detection via batch normalization. In *Thirty-seventh Conference on Neural Information*
596 *Processing Systems*, 2023a.
- 597 Xujia Li, Yuan Li, Xueying Mo, Hebing Xiao, Yanyan Shen, and Lei Chen. Diga: guided diffusion
598 model for graph recovery in anti-money laundering. In *Proceedings of the 29th ACM SIGKDD*
599 *Conference on Knowledge Discovery and Data Mining*, pp. 4404–4413, 2023b.
- 600 Yuhan Li, Peisong Wang, Zhixun Li, Jeffrey Xu Yu, and Jia Li. Zerog: Investigating cross-dataset
601 zero-shot transferability in graphs. In *Proceedings of the 30th ACM SIGKDD Conference on*
602 *Knowledge Discovery and Data Mining*, pp. 1725–1735, 2024.
- 603 Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation-based anomaly detection. *ACM Trans-*
604 *actions on Knowledge Discovery from Data (TKDD)*, 6(1):1–39, 2012.
- 605 Jiawei Liu, Cheng Yang, Zhiyuan Lu, Junze Chen, Yibo Li, Mengmei Zhang, Ting Bai, Yuan Fang,
606 Lichao Sun, Philip S Yu, et al. Towards graph foundation models: A survey and beyond. *arXiv*
607 *preprint arXiv:2310.11829*, 2023a.
- 608 Kay Liu, Yingdong Dou, Yue Zhao, Xueying Ding, Xiyang Hu, Ruitong Zhang, Kaize Ding, Canyu
609 Chen, Hao Peng, Kai Shu, et al. Bond: Benchmarking unsupervised outlier node detection on
610 static attributed graphs. *Advances in Neural Information Processing Systems*, 35:27021–27035,
611 2022.
- 612 Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-
613 train, prompt, and predict: A systematic survey of prompting methods in natural language pro-
614 cessing. *ACM Computing Surveys*, 55(9):1–35, 2023b.
- 615 Yang Liu, Xiang Ao, Zidi Qin, Jianfeng Chi, Jinghua Feng, Hao Yang, and Qing He. Pick and
616 choose: a gnn-based imbalanced learning approach for fraud detection. In *Proceedings of the*
617 *web conference 2021*, pp. 3168–3177, 2021a.
- 618 Yixin Liu, Zhao Li, Shirui Pan, Chen Gong, Chuan Zhou, and George Karypis. Anomaly detection
619 on attributed networks via contrastive self-supervised learning. *IEEE transactions on neural*
620 *networks and learning systems*, 33(6):2378–2392, 2021b.
- 621 Yixin Liu, Shiyuan Li, Yu Zheng, Qingfeng Chen, Chengqi Zhang, and Shirui Pan. Arc: A generalist
622 graph anomaly detector with in-context learning. *arXiv preprint arXiv:2405.16771*, 2024.
- 623 Zemin Liu, Xingtong Yu, Yuan Fang, and Xinming Zhang. Graphprompt: Unifying pre-training and
624 downstream tasks for graph neural networks. In *Proceedings of the ACM Web Conference 2023*,
625 pp. 417–428, 2023c.
- 626 Xuexiong Luo, Jia Wu, Amin Beheshti, Jian Yang, Xiankun Zhang, Yuan Wang, and Shan Xue.
627 Comga: Community-aware attributed graph anomaly detection. In *Proceedings of the Fifteenth*
628 *ACM International Conference on Web Search and Data Mining*, pp. 657–665, 2022.
- 629 Xiaoxiao Ma, Jia Wu, Shan Xue, Jian Yang, Chuan Zhou, Quan Z Sheng, Hui Xiong, and Leman
630 Akoglu. A comprehensive survey on graph anomaly detection with deep learning. *IEEE Trans-*
631 *actions on Knowledge and Data Engineering*, 35(12):12012–12038, 2021.
- 632 Julian John McAuley and Jure Leskovec. From amateurs to connoisseurs: modeling the evolution
633 of user expertise through online reviews. In *Proceedings of the 22nd international conference on*
634 *World Wide Web*, pp. 897–908, 2013.
- 635 Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. Deep learning for
636 anomaly detection: A review. *ACM computing surveys (CSUR)*, 54(2):1–38, 2021.
- 637 Zhen Peng, Minnan Luo, Jundong Li, Huan Liu, Qinghua Zheng, et al. Anomalous: A joint model-
638 ing approach for anomaly detection on attributed networks. In *IJCAI*, volume 18, pp. 3513–3519,
639 2018.
- 640 Hezhe Qiao and Guansong Pang. Truncated affinity maximization: One-class homophily modeling
641 for graph anomaly detection. *Advances in Neural Information Processing Systems*, 36, 2023.

- 648 Hezhe Qiao, Hanghang Tong, Bo An, Irwin King, Charu Aggarwal, and Guansong Pang. Deep graph
649 anomaly detection: A survey and new perspectives. *arXiv preprint arXiv:2409.09957*, 2024.
- 650
- 651 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
652 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
653 models from natural language supervision. In *International conference on machine learning*, pp.
654 8748–8763. PMLR, 2021.
- 655 Shebuti Rayana and Leman Akoglu. Collective opinion spam detection: Bridging review networks
656 and metadata. In *Proceedings of the 21th acm sigkdd international conference on knowledge
657 discovery and data mining*, pp. 985–994, 2015.
- 658 Gilbert W Stewart. On the early history of the singular value decomposition. *SIAM review*, 35(4):
659 551–566, 1993.
- 660
- 661 Xiangguo Sun, Hong Cheng, Jia Li, Bo Liu, and Jihong Guan. All in one: Multi-task prompting
662 for graph neural networks. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge
663 Discovery and Data Mining*, pp. 2120–2131, 2023.
- 664 Jianheng Tang, Jiajin Li, Ziqi Gao, and Jia Li. Rethinking graph neural networks for anomaly
665 detection. In *International Conference on Machine Learning*, pp. 21076–21089. PMLR, 2022.
- 666 Jianheng Tang, Fengrui Hua, Ziqi Gao, Peilin Zhao, and Jia Li. Gadbench: Revisiting and bench-
667 marking supervised graph anomaly detection. *Advances in Neural Information Processing Sys-
668 tems*, 36:29628–29653, 2023.
- 669
- 670 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
671 Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- 672 Qizhou Wang, Guansong Pang, Mahsa Salehi, Wray Buntine, and Christopher Leckie. Cross-
673 domain graph anomaly detection via anomaly-aware contrastive alignment. In *Proceedings of
674 the AAAI Conference on Artificial Intelligence*, volume 37, pp. 4676–4684, 2023.
- 675
- 676 Mark Weber, Giacomo Domeniconi, Jie Chen, Daniel Karl I Weidele, Claudio Bellei, Tom Robin-
677 son, and Charles E Leiserson. Anti-money laundering in bitcoin: Experimenting with graph
678 convolutional networks for financial forensics. *arXiv preprint arXiv:1908.02591*, 2019.
- 679 Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Sim-
680 plifying graph convolutional networks. In *International conference on machine learning*, pp.
681 6861–6871. PMLR, 2019.
- 682 Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A
683 comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and
684 Learning Systems*, 32(1):4–24, 2020.
- 685
- 686 Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie
687 Jegelka. Representation learning on graphs with jumping knowledge networks. In *International
688 conference on machine learning*, pp. 5453–5462. PMLR, 2018.
- 689 Zhiming Xu, Xiao Huang, Yue Zhao, Yushun Dong, and Jundong Li. Contrastive attributed network
690 anomaly detection with data augmentation. In *Pacific-Asia conference on knowledge discovery
691 and data mining*, pp. 444–457. Springer, 2022.
- 692
- 693 Yang Yang, Yuhong Xu, Yizhou Sun, Yuxiao Dong, Fei Wu, and Yueting Zhuang. Mining fraudsters
694 and fraudulent strategies in large-scale mobile social networks. *IEEE Transactions on Knowledge
695 and Data Engineering*, 33(1):169–179, 2019.
- 696
- 697 Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph
698 contrastive learning with augmentations. *Advances in neural information processing systems*, 33:
5812–5823, 2020.
- 699 Shijie Zhang, Hongzhi Yin, Tong Chen, Quoc Viet Nguyen Hung, Zi Huang, and Lizhen Cui. Gcn-
700 based user representation learning for unifying robust recommendation and fraudster detection.
701 In *Proceedings of the 43rd international ACM SIGIR conference on research and development in
information retrieval*, pp. 689–698, 2020.

702 Haihong Zhao, Aochuan Chen, Xiangguo Sun, Hong Cheng, and Jia Li. All in one and one for all: A
703 simple yet effective method towards cross-domain graph pretraining. In *Proceedings of the 30th*
704 *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 4443–4454, 2024.
705

706 Tong Zhao, Chuchen Deng, Kaifeng Yu, Tianwen Jiang, Daheng Wang, and Meng Jiang. Error-
707 bounded graph anomaly loss for gnns. In *Proceedings of the 29th ACM International Conference*
708 *on Information & Knowledge Management*, pp. 1873–1882, 2020.

709 Yu Zheng, Ming Jin, Yixin Liu, Lianhua Chi, Khoa T Phan, and Yi-Ping Phoebe Chen. Genera-
710 tive and contrastive self-supervised learning for graph anomaly detection. *IEEE Transactions on*
711 *Knowledge and Data Engineering*, 35(12):12220–12233, 2021.

712 Qihang Zhou, Guansong Pang, Yu Tian, Shibo He, and Jiming Chen. Anomalyclip: Object-agnostic
713 prompt learning for zero-shot anomaly detection. *arXiv preprint arXiv:2310.18961*, 2023.
714

715 Qihang Zhou, Guansong Pang, Yu Tian, Shibo He, and Jiming Chen. AnomalyCLIP: Object-
716 agnostic prompt learning for zero-shot anomaly detection. In *The Twelfth International Con-*
717 *ference on Learning Representations*, 2024.

718 Jiawen Zhu and Guansong Pang. Toward generalist anomaly detection via in-context residual learn-
719 ing with few-shot sample prompts. In *Proceedings of the IEEE/CVF Conference on Computer*
720 *Vision and Pattern Recognition*, pp. 17826–17836, 2024.
721

722 Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. Deep graph contrastive
723 representation learning. *arXiv preprint arXiv:2006.04131*, 2020.
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
749
750
751
752
753
754
755

756 A GRAPH SIMILARITY
757

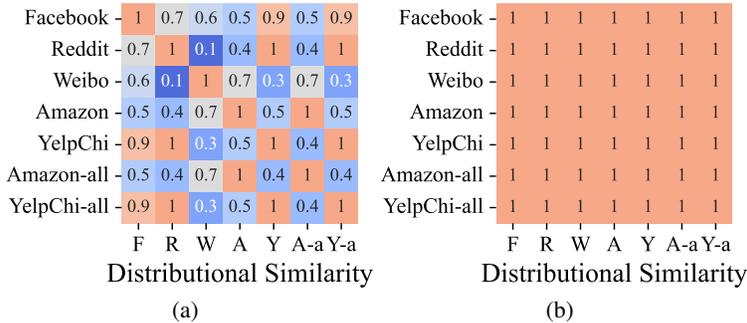
758 In addition to the visualization results presented in Figure 1, we further provide the distributional
759 similarity of various graphs in this section. Specifically, for dimension-aligned graphs across dif-
760 ferent distributions and domains, we measure their distributional similarity to analyze their diverse
761 semantics.

762 Given a graph $\mathcal{G}^{(i)} = (A^{(i)}, \tilde{X}^{(i)})$, the coordinate-wise mean $\boldsymbol{\mu}^{(i)} = [\mu_1^{(i)}, \dots, \mu_{d'}^{(i)}]$ and variance
763 $\boldsymbol{\sigma}^{(i)} = [\sigma_1^{(i)}, \dots, \sigma_{d'}^{(i)}]$ of $\tilde{X}^{(i)}$ are calculated and concatenated to form the distributional vector of
764 $\mathcal{G}^{(i)}$, *i.e.*, $\mathbf{d}_i = [\boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{(i)}]$. Then, the distribution similarity between $\mathcal{G}^{(i)}$ and $\mathcal{G}^{(j)}$ is measured via
765 the cosine similarity,
766

$$767 s_{ij} = \text{sim}(\mathbf{d}_i, \mathbf{d}_j). \tag{9}$$

768 The distributional similarities between graphs from different domains or distributions are shown
769 in Figure 4(a). From the figure, we can see that the distributional similarities are typically small,
770 demonstrating the diverse semantics of node features across graphs. Noth that, for Amazon &
771 Amazon-all and YelpChi & YelpChi-all, their distribution similarity is one, which can be attributed
772 to the fact that they are from the same distributions respectively but with different numbers of nodes
773 and structures.

774 To reduce the semantic gap among graphs for generalist GAD, we propose to calibrate the distri-
775 butions of all graphs into the same frame with coordinate-wise normalization. The distributional
776 similarity with normalization is illustrated in Figure 4(b). It is clear that the node attributes share
777 the same distribution after the normalization. In this way, the generalist model can better capture the
778 shared GAD patterns and generalize to different target graphs, as demonstrated in our experimental
779 results.



781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
Figure 4: (a) Distributional similarity between different graphs without coordinate-wise normaliza-
tion. (b) Distributional similarity between different graphs with the coordinate-wise normalization.

800 B DETAILS ON PRE-TRAINING OF NEIGHBORHOOD AGGREGATION
801 NETWORKS

802 We pre-train the neighborhood aggregation network g via graph contrastive learning (Zhu et al.,
803 2020) for subsequent graph prompt learning so that the generic normality and abnormality can be
804 captured in the prompts.

805 Specifically, given the training graph $\mathcal{G} = (A, X)$, to construct contrastive views for graph con-
806 trastive learning, two widely used graph augmentations are employed, *i.e.*, edge removal and at-
807 tribute masking (Zhu et al., 2020). The edge removal randomly drops a certain portion of existing
808 edges in \mathcal{G} and the attribute masking randomly masks a fraction of dimensions with zeros in node
809 attributes, *i.e.*,

$$\hat{A} = A \circ R, \quad \hat{X} = [\mathbf{x}_1 \circ \mathbf{m}, \dots, \mathbf{x}_N \circ \mathbf{m}]^T, \tag{10}$$

where $R \in \{0, 1\}^{N \times N}$ is the edge masking matrix whose entry is drawn from a Bernoulli distribu-
tion controlled by the edge removal probability, $\mathbf{m} \in \{0, 1\}^d$ is the attribute masking vector whose

entry is independently drawn from a Bernoulli distribution with the attribute masking ratio, and \circ denotes the Hadamard product.

By applying the graph augmentations to the original graph, the corrupted graph $\hat{\mathcal{G}} = (\hat{A}, \hat{X})$ forms the contrastive view for the original graph $\mathcal{G} = (A, X)$. Then, $\hat{\mathcal{G}}$ and \mathcal{G} are fed to the shared model g followed by the non-linear projection to obtain the corresponding node embeddings, *i.e.*, \hat{Z}' and Z' . For graph contrastive learning, the embeddings of the same node in different views are pulled closer while the embeddings of other nodes are pushed apart. The pairwise objective for each node pair (\hat{z}'_i, z'_i) can be formulated as:

$$\ell(\hat{z}'_i, z'_i) = -\log \frac{e^{\text{sim}(\hat{z}'_i, z'_i)/\tau}}{e^{\text{sim}(\hat{z}'_i, z'_i)/\tau} + \sum_{j \neq i}^N e^{\text{sim}(\hat{z}'_i, z'_j)/\tau} + \sum_{j \neq i}^N e^{\text{sim}(\hat{z}'_j, z'_i)/\tau}}, \quad (11)$$

where $\text{sim}(\cdot)$ represents the cosine similarity and τ is a temperature hyperparameter. Therefore, the overall objective can be defined as follows:

$$\mathcal{L}_{\text{contrast}} = \frac{1}{2N} \sum_{i=1}^N (\ell(\hat{z}'_i, z'_i) + \ell(z'_i, \hat{z}'_i)). \quad (12)$$

With the objective Eq.(12), the model g is optimized to learn transferable discriminative representations of nodes.

C ALGORITHMS

The training and inference processes of UNPrompt are summarized in Algorithms 1 and Algorithm 2, respectively.

Algorithm 1: Training of UNPrompt

- 1: **Input:** Training graph $\mathcal{G}_{\text{train}} = (A, X)$; training epoch E
 - 2: **Output:** Neighborhood aggregation network g , graph prompts $P = [\mathbf{p}_1, \dots, \mathbf{p}_K]$, and transformation h .
 - 3: Perform feature unification of X .
 - 4: Pre-train g on $\mathcal{G}_{\text{train}}$ with graph contrastive learning in Eq.(12).
 - 5: Keep model g frozen.
 - 6: **for** $e = 1, \dots, E$ **do**
 - 7: Obtain modified node attribute with prompts via Eq.(6).
 - 8: Obtain the neighborhood aggregated representation \tilde{Z} via Eq.(3).
 - 9: Obtain the node representations Z via Eq.(4).
 - 10: Transform \tilde{Z} and Z with h via Eq.(7).
 - 11: Optimize P and h by minimizing Eq.(8).
 - 12: **end for**
-

Algorithm 2: Inference of UNPrompt

- 1: **Input:** Testing graphs $\mathcal{T}_{\text{test}} = \{\mathcal{G}_{\text{test}}^{(1)}, \dots, \mathcal{G}_{\text{test}}^{(n)}\}$, neighborhood aggregation network g , graph prompts $P = [\mathbf{p}_1, \dots, \mathbf{p}_K]$, and transformation h .
 - 2: **Output:** Normal score of testing nodes.
 - 3: **for** $\mathcal{G}_{\text{test}}^{(i)} = (A^{(i)}, X^{(i)}) \in \mathcal{T}_{\text{test}}$ **do**
 - 4: Perform feature unification of $X^{(i)}$.
 - 5: Obtain modified node attribute with prompts via Eq.(6).
 - 6: Obtain the neighborhood aggregated representation $\tilde{Z}^{(i)}$ via Eq.(3).
 - 7: Obtain the node representations $Z^{(i)}$ via Eq.(4).
 - 8: Transform $\tilde{Z}^{(i)}$ and $Z^{(i)}$ with h via Eq.(7).
 - 9: Obtain the normal score of nodes via Eq.(5).
 - 10: **end for**
-

D UNSUPERVISED GAD WITH UNPROMPT

To demonstrate the wide applicability of the proposed method UNPrompt, we further perform unsupervised GAD with UNPrompt which focuses on detecting anomalous nodes within one graph and does not have access to any node labels during training. Specifically, we adopt the same pipeline in the generalist GAD setting, *i.e.*, graph contrastive pertaining and neighborhood prompt learning. Since we focus on anomaly detection on each graph separately, the node attribute unification is discarded for unsupervised GAD. However, the absence of node labels poses a challenge to learning meaningful neighborhood prompts for anomaly detection. To overcome this issue, we propose to utilize the pseudo-labeling technique to guide the prompt learning. Specifically, the normal score of each node is calculated by the neighborhood-based latent attribute predictability after the graph contrastive learning process:

$$s_i = \text{sim}(\mathbf{z}_i, \tilde{\mathbf{z}}_i), \quad (13)$$

where \mathbf{z}_i is the node representation learned by graph contrastive learning and $\tilde{\mathbf{z}}_i$ is the corresponding aggregated neighborhood representation. Higher s_i of node v_i typically indicates a higher probability of v_i being normal nodes. Therefore, more emphasis should be put on high-score nodes when learning neighborhood prompts. To achieve this, the normal score s_i is transformed into the loss weight $w_i = \text{Sigmoid}(\alpha(s_i - t))$ where t is a threshold and α is the scaling parameter. In this way, w_i would approach 1 if $s_i > t$ and 0 otherwise. Overall, the objective for unsupervised GAD using UNPrompt can be formulated as follows:

$$\mathcal{L} = \sum_i^N (-w_i \text{sim}(\mathbf{z}_i, \tilde{\mathbf{z}}_i)) + \lambda \sum_{j, j \neq i}^N \text{sim}(\mathbf{z}_i, \tilde{\mathbf{z}}_j), \quad (14)$$

where the second term is a regularization term employed to prevent the node embeddings from being collapsed into the same and λ is a trade-off hyperparameter.

Note that, we only focus on maximizing the latent attribute predictability of high-score nodes without minimizing the predictability of low-score nodes in the above objective. These low-score nodes could also be normal nodes with high probability as the score from Eq.(13) is only obtained from the pre-trained model, resulting in the score not being fully reliable. If the predictability is also minimized for these nodes, conflicts would be induced for neighborhood prompt learning, limiting the performance of unsupervised GAD. After optimization, the latent attribute predictability is also directly used as the normal score for the final unsupervised GAD.

E TIME COMPLEXITY ANALYSIS

Theoretical Analysis. In this section, we analyze the time complexity of training UNPrompt. As discussed in the main paper, UNPrompt first pre-trains the aggregation network with graph contrastive learning. Then, the model remains frozen when optimizing neighborhood graph prompts and the transformation layer to capture the generalized normal and abnormal graph patterns. In the experimental section, we employ a one-layer aggregation network, denoting the number of hidden units as d_h . The time complexity of the graph contrastive learning is $\mathcal{O}(4E_1(|A|d_h + Nd_hd' + 6Nd_h^2))$, where $|A|$ returns of the number of edges of the $\mathcal{G}_{\text{train}}$, N is the number of nodes, d' represents the predefined dimensionality of node attributes, and E_1 is the number of training epoch. After that, we freeze the learned model and learn the learnable neighborhood prompt tokens and the transformation layer to capture the shared anomaly patterns. In our experiments, we set the size of each graph prompt to K and implement the classification head as a single-layer MLP with the same hidden units d_h . Given the number of the training epoch E_2 , the time complexity of optimizing the graph prompt and the transformation layer is $\mathcal{O}((4KNd' + 2Nd_h^2)E_2)$, which includes both the forward and backward propagation. Note that, despite the neighborhood aggregation model being frozen, the forward and backward propagations of the model are still needed to optimize the task-specific graph prompts and the transformation layer. Therefore, the overall time complexity of UNPrompt is $\mathcal{O}(4E_1(|A|d_h + Nd_hd' + 6Nd_h^2) + 2E_2(|A|d_h + Nd_hd' + 2KNd' + Nd_h^2))$, which is linear to the number of nodes, the number of edges, and the number of node attributes of the training graph. Note that, after the training, the learned generalist model is directly utilized to perform anomaly detection on various target graphs without any further training.

Table 4: Training time and inference time (seconds) for different methods.

Methods	AnomalyDAE	TAM	GAT	BWGNN	UNPrompt (Ours)
Training Time	86.04	479.70	2.43	4.86	2.08
Inference Time	264.29	521.92	300.90	330.99	58.95

Empirical Computational Complexity Analysis. In Table 4, we report the training time and inference time of different methods, where two representative unsupervised methods (AnomalyDAE and TAM) and two supervised methods (GAT and BWGNN) are used for comparison to our method UNPrompt. The results show that the proposed method requires much less training and inference time compared to other baselines, demonstrating the efficiency of the proposed UNPrompt. Note that, TAM has the highest time consumption, which can be attributed to that it performs multiple graph truncation and learns multiple local affinity maximization networks.

F EXPERIMENTAL SETUP

F.1 DETAILS ON DATASETS

We conduct the experiments using seven real-world with genuine anomalies in diverse online shopping services and social networks, including Facebook (Xu et al., 2022), Reddit (Kumar et al., 2019), Weibo (Zhao et al., 2020), Amazon (McAuley & Leskovec, 2013) and YelpChi (Rayana & Akoglu, 2015) as well as two large-scale graph datasets including Amazon-all (McAuley & Leskovec, 2013) and YelpChi-all (Rayana & Akoglu, 2015). The statistical information including the number of nodes, edge, the dimension of the feature, and the anomaly rate of the datasets can be found in Table 5. The more detailed description of each dataset is given as follows

- **Facebook** (Xu et al., 2022). Facebook is a social network where each node represents a user, and edges signify relationships between users. Ground truth anomalies are nodes that either connect to randomly selected circles or exhibit abnormal attributes, as described in (Ding et al., 2019; Liu et al., 2021b).
- **Reddit** (Kumar et al., 2019). Reddit is a forum-based network derived from the social media platform Reddit, where nodes represent users, and the embeddings of post texts serve as attributes. Users who have been banned from the platform are labeled as anomalies.
- **Weibo** (Kumar et al., 2019). Weibo is a social network and their associated hashtags are obtained from the Tencent Weibo platform. Users who engaged in at least five of these activities are labeled as anomalies while the others are classified as normal samples. Suspicious activities are defined as two posts made within a specific timeframe, such as 60 seconds. The attributes of each node include the location of a micro-blog post and bag-of-words features.
- **Amazon** (McAuley & Leskovec, 2013). Amazon is a graph dataset that captures the relations between users and product reviews. There are 25 handcrafted features used as the node attribute (Zhang et al., 2020). The users with more than 80% helpful votes are labeled as normal entities and users with less than 20% helpful votes as anomalies. Amazon is constructed by extracting the Amazon-UPU dataset that connects the users who give reviews to at least one common product.
- **YelpChi** (Rayana & Akoglu, 2015). YelpChi includes hotel and restaurant reviews filtered (spam) and recommended (legitimate) by Yelp. There are 32 handcrafted features used as node attributes (Rayana & Akoglu, 2015). The users with more than 80% helpful votes are labeled as benign entities and users with less than 20% helpful votes as fraudulent entities. The YelpChi is constructed by extracting YelpChi-RUR which connects reviews posted by the same user.
- **Amazon-all** (McAuley & Leskovec, 2013). Amazon-all includes three types of relations: U-P-U (users reviewing at least one same product), U-S-U (users giving at least one same star rating within one week), and U-V-U (users with top-5% mutual review similarities). Amazon-all is formed by treating the different relations as a single relation following Chen et al. (2022); Qiao & Pang (2023).

Table 5: Key statistics of the real-world GAD datasets with real anomalies.

Data set	Type	Nodes	Edges	Attributes	Anomalies(Rate)
Facebook	Social Networks	1,081	55,104	576	27(2.49%)
Reddit	Social Networks	10,984	168,016	64	366(3.33%)
Weibo	Social Networks	8,405	407,963	400	868(10.30%)
Amazon	Co-review	10,244	175,608	25	693(6.66%)
YelpChi	Co-review	24,741	49,315	32	1,217(4.91%)
Amazon-all	Co-review	11,944	4,398,392	25	821(6.87%)
YelpChi-all	Co-review	45,941	3,846,979	32	6,674(14.52%)
Disney	Co-purchase	124	335	28	6(4.84%)
Elliptic	Payment Flow	203,769	234,355	166	4,545(9.76%)

- **YelpChi-all** (Rayana & Akoglu, 2015). Similar to Amazon-all, YelpChi-all includes three types of edges: R-U-R (reviews posted by the same user), R-S-R (reviews for the same product with the same star rating), and R-T-R (reviews for the same product posted in the same month). YelpChi-all is formed by treating the different relations as a single relation following Chen et al. (2022); Qiao & Pang (2023).

F.2 MORE IMPLEMENTATION DETAILS

Generalist GAD. For the graph contrastive learning-based pre-training, the probabilities of edge removal and attribute masking are by default set to 0.2 and 0.3 respectively. Besides, the learning rate is set to 0.001 with the Adam optimizer, the training epoch is set to 200 and the temperature τ is 0.5.

For the neighborhood prompt learning, the learning rate is also set to 0.001 with the Adam optimizer, and the training epoch is set to 900. Note that, since we focus on generalist GAD, we do not perform any hyperparameter search for specific target graphs. Instead, the results of all target graphs are obtained with the same hyperparameter settings.

Unsupervised GAD. Similar to the generalist GAD setting, the hidden units of the neighborhood aggregation network and the transformation layer are set to 128 for all graphs. The threshold t is determined by the 40th percentile of the normal scores obtained by the pre-trained model g , and the scaling parameter α is set to 10 for all graphs. Besides, we utilize random search to find the optimal hyperparameters of the size of neighborhood prompts K and the trade-off parameter λ .

For both generalist and unsupervised GAD, the code is implemented with Pytorch (version: 1.13.1), DGL (version: 1.0.1), OGB (version: 1.3.6), and Python 3.8.19. All experiments are conducted on a Linux server with an Intel CPU (Intel Xeon Gold 6346 3.1GHz) and a Nvidia A40 GPU.

G MORE EXPERIMENTAL RESULTS

G.1 GENERALIST PERFORMANCE WITH DIFFERENT COMMON DIMENSIONALITIES

For the results reported in the main paper, the common dimensionality is set to eight. In this subsection, we further evaluate the generalist anomaly detection with different common dimensionalities. Specifically, the dimensionality varies in $[2, 4, 6, 8, 10, 12]$ and the results are reported in Figure 5.

From the figure, we can see that small dimensionality leads to poor generalist anomaly detection performance. This is attributed to the fact that much attribute information would be discarded with a small dimensionality. By increasing the common dimensionality, more attribute information is retained, generally resulting in much better detection performance.

G.2 RESULTS ON TWO OTHER GRAPHS FROM DIFFERENT DOMAINS

Besides the social networks and co-review graphs, we further evaluate the performance of UN-Prompt on Disney (Liu et al., 2022) and Elliptic (Weber et al., 2019). These two datasets consist of

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

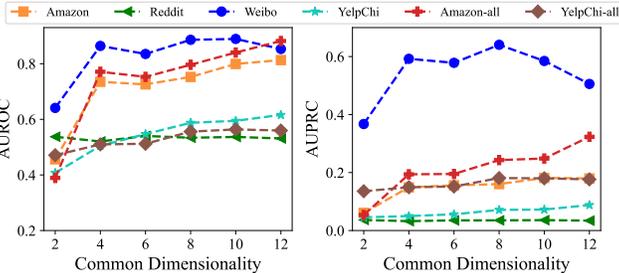


Figure 5: AUROC and AUPRC results of UNPrompt w.r.t. varying common dimensionality.

Table 6: AUROC and AUPRC results on two additional real-world GAD datasets with the models trained on Facebook only. For each dataset, the best performance per column within each metric is boldfaced, with the second-best underlined. ‘‘Avg’’ denotes the averaged performance of each method.

Metric	Method	Dataset		Avg.
		Disney	Elliptic	
AUROC	Unsupervised Methods			
	AnomalyDAE	0.4853 \pm 0.003	0.4197 \pm 0.109	0.4525
	CoLA	0.4696 \pm 0.065	0.5572 \pm 0.019	0.5134
	HCM-A	0.2014 \pm 0.015	0.2975 \pm 0.004	0.2495
	GADAM	0.4288 \pm 0.023	0.3922 \pm 0.012	0.4105
	TAM	0.4773 \pm 0.003	0.3282 \pm 0.003	0.4028
	Supervised Methods			
	GCN	0.5000 \pm 0.000	0.7640 \pm 0.002	0.6320
	GAT	0.5175 \pm 0.054	0.6588 \pm 0.019	0.5882
	BWGNN	0.6073 \pm 0.026	0.5843 \pm 0.101	0.5958
	GHRN	0.5336 \pm 0.030	0.5400 \pm 0.103	0.5368
	XGBGraph	0.6692 \pm 0.000	0.4274 \pm 0.000	0.5483
UNPrompt (Ours)	<u>0.6412</u> \pm 0.030	0.5901 \pm 0.026	<u>0.6157</u>	
AUPRC	Unsupervised Methods			
	AnomalyDAE	0.0566 \pm 0.000	0.0798 \pm 0.014	0.0682
	CoLA	0.0701 \pm 0.023	0.0998 \pm 0.005	0.0850
	HCM-A	0.0355 \pm 0.001	0.0776 \pm 0.000	0.0566
	GADAM	0.0651 \pm 0.012	0.0733 \pm 0.001	0.0692
	TAM	0.0628 \pm 0.001	0.0697 \pm 0.001	0.0663
	Supervised Methods			
	GCN	0.0484 \pm 0.000	0.1963 \pm 0.002	0.1224
	GAT	0.0530 \pm 0.004	0.1366 \pm 0.010	0.0948
	BWGNN	0.0624 \pm 0.003	0.1158 \pm 0.026	0.0891
	GHRN	0.0519 \pm 0.003	0.1148 \pm 0.041	0.0834
	XGBGraph	0.1215 \pm 0.000	0.0816 \pm 0.000	0.1016
UNPrompt (Ours)	0.1236 \pm 0.031	0.1278 \pm 0.004	0.1257	

co-purchase network and financial network respectively. The statistics of them are also summarized in Table 5. We follow exactly the same experimental settings in the main paper.

The results of all competing methods are reported in Table 6. It is clear that UNPrompt can still achieve promising performance, demonstrating the generality of UNPrompt across different graphs. Although GCN and XGBGraph obtain the best AUROC performance on Disney and Elliptic respectively, they perform poorly on most of the other datasets. UNPrompt ranks in second in the average AUROC performance and the best AUPRC performance here. This is consistent with the superior performance of UNPrompt in Table 1.

G.3 INCORPORATING COORDINATE-WISE NORMALIZATION INTO BASELINES

We further conduct experiments by incorporating the proposed coordinate-wise normalization into the baselines to evaluate whether the normalization could facilitate the baselines. Without loss of the generality, three unsupervised methods (AnomalyDAE, CoLA and TAM) and three supervised methods (GCN, BWGNN and GHRN) are used and the results are reported in Table 7.

From the table, we can see that the proposed coordinate-wise normalization does not improve the baselines consistently but downgrades most of the baselines. This can be attributed to two reasons. First, while the proposed coordinate-wise normalization unifies the semantics of different graphs into the common space, the discrimination between normal and abnormal patterns would also be

Table 7: AUROC and AUPRC results of several baselines with coordinate-wise normalization (CN).

Metric	Method	Dataset						Avg.
		Amazon	Reddit	Weibo	YelpChi	Aamazon-all	YelpChi-all	
AUROC	Unsupervised Methods							
	AnomalyDAE	0.5818 \pm 0.039	0.5016 \pm 0.032	0.7785 \pm 0.058	0.4837 \pm 0.094	0.7228 \pm 0.023	0.5002 \pm 0.018	0.5948
	+ CN	0.4359 \pm 0.053	0.4858 \pm 0.063	0.4526 \pm 0.074	0.5992 \pm 0.028	0.2833 \pm 0.039	0.5080 \pm 0.013	0.4608
	CoLA	0.4580 \pm 0.054	0.4623 \pm 0.005	0.3924 \pm 0.041	0.4907 \pm 0.017	0.4091 \pm 0.052	0.4879 \pm 0.010	0.4501
	+ CN	0.4729 \pm 0.019	0.5299 \pm 0.008	0.3401 \pm 0.026	0.3640 \pm 0.006	0.5424 \pm 0.019	0.4882 \pm 0.008	0.4563
	TAM	0.4720 \pm 0.005	0.5725 \pm 0.004	0.4867 \pm 0.028	0.5035 \pm 0.014	0.7543 \pm 0.002	0.4216 \pm 0.002	0.5351
	+ CN	0.4509 \pm 0.015	0.5526 \pm 0.006	0.4723 \pm 0.007	0.5189 \pm 0.006	0.7580 \pm 0.004	0.4057 \pm 0.002	0.5264
	Supervised Methods							
	GCN	0.5988 \pm 0.016	0.5645 \pm 0.000	0.2232 \pm 0.074	0.5366 \pm 0.019	0.7195 \pm 0.002	0.5486 \pm 0.001	0.5319
	+ CN	0.5694 \pm 0.014	0.5349 \pm 0.008	0.0632 \pm 0.005	0.3954 \pm 0.002	0.6798 \pm 0.009	0.5550 \pm 0.005	0.4663
	BWGNN	0.4769 \pm 0.020	0.5208 \pm 0.016	0.4815 \pm 0.108	0.5538 \pm 0.027	0.3648 \pm 0.050	0.5282 \pm 0.015	0.4877
	+ CN	0.4745 \pm 0.048	0.4942 \pm 0.011	0.2538 \pm 0.038	0.4727 \pm 0.016	0.6307 \pm 0.077	0.5221 \pm 0.025	0.4747
	GHRN	0.4560 \pm 0.033	0.5253 \pm 0.006	0.5318 \pm 0.038	0.5524 \pm 0.020	0.3382 \pm 0.085	0.5125 \pm 0.016	0.4860
	+ CN	0.4308 \pm 0.024	0.5061 \pm 0.026	0.2621 \pm 0.043	0.4781 \pm 0.018	0.5712 \pm 0.046	0.5200 \pm 0.009	0.4614
	UNPrompt (Ours)	0.7525 \pm 0.016	0.5337 \pm 0.002	0.8860 \pm 0.007	0.5875 \pm 0.016	0.7962 \pm 0.022	0.5558 \pm 0.012	0.6853
	AUPRC	Unsupervised Methods						
AnomalyDAE		0.0833 \pm 0.015	0.0327 \pm 0.004	0.6064 \pm 0.031	0.0624 \pm 0.017	0.1921 \pm 0.026	0.1484 \pm 0.009	0.1876
+ CN		0.0596 \pm 0.009	0.0333 \pm 0.007	0.1910 \pm 0.049	0.0874 \pm 0.011	0.0495 \pm 0.006	0.1527 \pm 0.007	0.0956
CoLA		0.0669 \pm 0.002	0.0391 \pm 0.004	0.1189 \pm 0.014	0.0511 \pm 0.000	0.0861 \pm 0.019	0.1466 \pm 0.003	0.0848
+ CN		0.0669 \pm 0.002	0.0360 \pm 0.002	0.1618 \pm 0.027	0.0370 \pm 0.000	0.0934 \pm 0.017	0.1446 \pm 0.005	0.0899
TAM		0.0666 \pm 0.001	0.0413 \pm 0.001	0.1240 \pm 0.014	0.0524 \pm 0.002	0.1736 \pm 0.004	0.1240 \pm 0.001	0.0970
+ CN		0.0606 \pm 0.003	0.0394 \pm 0.001	0.1044 \pm 0.005	0.0542 \pm 0.001	0.2482 \pm 0.013	0.1213 \pm 0.001	0.1047
Supervised Methods								
GCN		0.0891 \pm 0.007	0.0439 \pm 0.000	0.1109 \pm 0.020	0.0648 \pm 0.009	0.1536 \pm 0.002	0.1735 \pm 0.000	0.1060
+ CN		0.0770 \pm 0.003	0.0355 \pm 0.001	0.0548 \pm 0.000	0.0401 \pm 0.000	0.1383 \pm 0.006	0.1789 \pm 0.002	0.0874
BWGNN		0.0652 \pm 0.002	0.0389 \pm 0.003	0.2241 \pm 0.046	0.0708 \pm 0.018	0.0586 \pm 0.003	0.1605 \pm 0.005	0.1030
+ CN		0.0684 \pm 0.014	0.0320 \pm 0.001	0.2576 \pm 0.031	0.0516 \pm 0.004	0.1557 \pm 0.115	0.1585 \pm 0.010	0.1206
GHRN		0.0633 \pm 0.003	0.0407 \pm 0.002	0.1965 \pm 0.059	0.0661 \pm 0.010	0.0569 \pm 0.006	0.1505 \pm 0.005	0.0957
+ CN		0.0586 \pm 0.004	0.0330 \pm 0.002	0.2663 \pm 0.038	0.0525 \pm 0.004	0.0898 \pm 0.015	0.1570 \pm 0.007	0.1095
UNPrompt (Ours)		0.1602 \pm 0.013	0.0351 \pm 0.000	0.6406 \pm 0.026	0.0712 \pm 0.008	0.2430 \pm 0.028	0.1810 \pm 0.012	0.2219

compressed. This requires the generalist anomaly detector to capture the fine-grained differences between normal and abnormal patterns. Second, these baselines are not designed to capture generalized abnormality and normality across graphs, failing to capture and discriminate the generalized nuance. By contrast, we reveal that the predictability of latent node attributes can serve as a generalized anomaly measure and learn highly generalized normal and abnormal patterns via latent node attribute prediction. In this way, the graph-agnostic anomaly measure can be well generalized across graphs.

G.4 GENERALIST PERFORMANCE WITH DIFFERENT TRAINING GRAPH

In the main paper, we report the generalist performance of UNPrompt by using Facebook as the training graph. To further demonstrate the generalizability of UNPrompt, we conduct additional experiments by using Amazon as the training graph and testing the learned generalist model on the rest graphs. Note that, Facebook and Amazon are from different domains, which are the social network and co-review network respectively.

The AUROC and AUPRC results of all methods are reported in Table 8. Similar to the observations when taking Facebook as the training graph, UNPrompt achieves the best average performance in terms of both AUROC and AUPRC when training on Amazon, demonstrating the generalizability and effectiveness of UNPrompt with different training graphs. Note that, the training graph Amazon and the target graph Amazon-all come from the same distribution but have different numbers of nodes and graph structures. Intuitively, all the methods should achieve promising performance on Amazon-all. However, only a few methods achieve this goal, including BWGNN, GHRN, and our method. The failures of other baselines can be attributed to the more complex graph structure of Amazon-all hinders the generalizability of these methods. Moreover, compared to BWGNN and GHRN, our method performs more stably across different datasets. This demonstrates the importance of capturing intrinsic normal and abnormal patterns for graph anomaly detection.

Table 8: AUROC and AUPRC results on six real-world GAD datasets with the generalist model trained on Amazon. For each dataset and metric, the best performance per column is boldfaced, with the second-best underlined. “Avg” denotes the averaged performance of each method.

Metric	Method	Dataset						Avg.
		Facebook	Reddit	Weibo	YelpChi	Aamazon-all	YelpChi-all	
AUROC	Unsupervised Methods							
	AnomalyDAE	0.6123 \pm 0.141	0.5799 \pm 0.035	<u>0.7884</u> \pm 0.031	0.4788 \pm 0.046	0.6233 \pm 0.070	0.4912 \pm 0.009	0.5957
	CoLA	0.5427 \pm 0.109	0.4962 \pm 0.025	0.3987 \pm 0.017	0.3358 \pm 0.012	0.4751 \pm 0.014	0.4937 \pm 0.003	0.4570
	HCM-A	0.5044 \pm 0.047	0.4993 \pm 0.057	0.4937 \pm 0.056	0.5000 \pm 0.000	0.4785 \pm 0.016	0.4958 \pm 0.003	0.4953
	GADAM	0.6024 \pm 0.033	0.4720 \pm 0.062	0.4324 \pm 0.047	0.4299 \pm 0.023	0.5199 \pm 0.072	0.5289 \pm 0.017	0.4976
	TAM	0.5496 \pm 0.038	<u>0.5764</u> \pm 0.003	0.4876 \pm 0.029	0.5091 \pm 0.014	0.7525 \pm 0.002	0.4268 \pm 0.002	0.5503
	Supervised Methods							
	GCN	0.6892 \pm 0.004	0.5658 \pm 0.000	0.2355 \pm 0.019	<u>0.5277</u> \pm 0.002	0.7503 \pm 0.002	0.5565 \pm 0.000	0.5542
	GAT	0.3886 \pm 0.118	0.4997 \pm 0.012	0.3897 \pm 0.134	0.5051 \pm 0.019	0.5007 \pm 0.006	0.4977 \pm 0.006	0.4636
	BWGNN	0.5441 \pm 0.020	0.4026 \pm 0.028	0.4214 \pm 0.039	0.4908 \pm 0.013	<u>0.9684</u> \pm 0.005	0.5841 \pm 0.002	0.5686
	GHRN	0.5242 \pm 0.013	0.4096 \pm 0.021	0.4783 \pm 0.021	0.5036 \pm 0.016	0.9601 \pm 0.018	0.6045 \pm 0.022	0.5800
	XGBGraph	0.4869 \pm 0.069	0.4869 \pm 0.069	0.7843 \pm 0.090	0.4773 \pm 0.022	0.9815 \pm 0.000	0.5869 \pm 0.014	0.6340
	Our	0.7917 \pm 0.021	0.5356 \pm 0.005	0.8192 \pm 0.015	0.5362 \pm 0.007	0.9289 \pm 0.007	0.5448 \pm 0.009	0.6927
	AUPRC	Unsupervised Methods						
AnomalyDAE		0.0675 \pm 0.028	0.0413 \pm 0.005	0.6172 \pm 0.015	0.0647 \pm 0.016	0.1025 \pm 0.026	0.1479 \pm 0.006	0.1735
CoLA		0.0468 \pm 0.026	0.0327 \pm 0.002	0.0956 \pm 0.005	0.0361 \pm 0.001	0.0678 \pm 0.005	0.1474 \pm 0.001	0.0711
HCM-A		0.0249 \pm 0.003	0.0374 \pm 0.008	0.0979 \pm 0.011	0.0511 \pm 0.000	0.0727 \pm 0.006	0.1453 \pm 0.000	0.0716
GADAM		0.0461 \pm 0.014	0.0299 \pm 0.004	0.0917 \pm 0.007	0.0428 \pm 0.002	0.0773 \pm 0.024	0.1602 \pm 0.010	0.0747
TAM		0.0243 \pm 0.002	<u>0.0417</u> \pm 0.001	0.1266 \pm 0.015	0.0532 \pm 0.002	0.1771 \pm 0.002	0.1271 \pm 0.001	0.0917
Supervised Methods								
GCN		0.0437 \pm 0.001	0.0449 \pm 0.000	0.2527 \pm 0.026	0.0763 \pm 0.001	0.1738 \pm 0.002	0.1759 \pm 0.000	0.1279
GAT		0.0445 \pm 0.039	0.0327 \pm 0.001	0.0892 \pm 0.016	0.0595 \pm 0.003	0.0697 \pm 0.001	0.1478 \pm 0.003	0.0739
BWGNN		0.0289 \pm 0.003	0.0263 \pm 0.002	0.2735 \pm 0.026	0.0543 \pm 0.004	<u>0.8406</u> \pm 0.012	0.1975 \pm 0.031	0.2369
GHRN		0.0254 \pm 0.001	0.0265 \pm 0.002	0.3103 \pm 0.013	0.0541 \pm 0.005	0.8142 \pm 0.045	0.2015 \pm 0.015	0.2387
XGBGraph		0.0268 \pm 0.006	0.0315 \pm 0.000	0.4116 \pm 0.040	0.0500 \pm 0.003	0.8673 \pm 0.000	0.1994 \pm 0.012	0.2644
Our		0.2291 \pm 0.023	0.0340 \pm 0.001	<u>0.4746</u> \pm 0.033	0.0610 \pm 0.003	0.7329 \pm 0.042	0.1767 \pm 0.004	0.2847

H INDUCTIVE LEARNING VS. ZERO-SHOT GENERALIST LEARNING

Despite our method and inductive graph learning (Hamilton et al., 2017; Xu et al., 2018) are both focused on evaluating the learned models on unseen graph data during inference, there are fundamental differences between our zero-shot learning and the inductive graph learning. We clarify the differences as follows:

- For inductive graph learning, the training dataset and testing dataset come from the same graph source. For example, for the graph classification task in Xu et al. (2018), 20 graphs of protein-protein interaction (PPI) datasets are used for training, and 2 other graphs are used for testing. These graphs both belong to the same protein-protein interaction graph dataset with the same attribute distribution and semantics. Therefore, the learned model can be easily generalized to the test graphs.
- In our zero-shot setting, the training dataset and testing dataset are from different domains and distributions. They differ in the dimensionality of node attributes and graph semantics. For example, as a shopping network dataset, Amazon contains the relationships between users and reviews, and the node attribute dimensionality is 25. Differently, Facebook, a social network dataset, describes relationships between users with 576-dimensional attributes. This is one fundamental difference between the inductive setting and our zero-shot setting. Moreover, our zero-shot setting requires the learned models to be directly applied to other graphs from different domains without any further tuning/training or labeled nodes of the target graphs. This requires the learned model to capture the more generalized patterns for anomaly detection based on only one training graph, resulting in a task being more challenging than the mentioned inductive learning.
- There are also studies (Ding et al., 2021a) on inductive graph anomaly detection, but the problem setting is also fundamentally different from our setting. In particular, to allow the evaluation of inductive detection, Ding et al. (2021a) samples nodes from the same graph to construct two graph datasets, with one graph used for training and another used for testing, leading to the fact that the training and test datasets are essentially from the same distribution. This is fundamentally different from our settings, where training and testing datasets are separately from highly heterogeneous distributions and domains.

Table 9: AUROC and AUPRC results on six real-world GAD datasets with the generalist model trained on Facebook. For each dataset and metric, the best performance per column is boldfaced, with the second-best underlined. “Avg” denotes the averaged performance of each method.

Metric	Method	Dataset						Avg.
		Amazon	Reddit	Weibo	YelpChi	Aamazon-all	YelpChi-all	
AUROC	GraphSAGE	0.4276 \pm 0.156	<u>0.5275</u> \pm 0.011	0.0975 \pm 0.002	0.4593 \pm 0.000	0.3276 \pm 0.074	0.4720 \pm 0.006	0.3853
	AEGIS	0.4664 \pm 0.030	0.3530 \pm 0.016	0.4979 \pm 0.048	<u>0.5267</u> \pm 0.016	0.4375 \pm 0.149	0.5116 \pm 0.022	0.4655
	Ours	0.7525 \pm 0.016	0.5337 \pm 0.002	0.8860 \pm 0.007	0.5875 \pm 0.016	0.7962 \pm 0.022	0.5558 \pm 0.012	0.6853
AUPRC	GraphSAGE	0.0601 \pm 0.015	0.0361 \pm 0.001	0.0858 \pm 0.000	0.0468 \pm 0.001	0.0493 \pm 0.006	0.1358 \pm 0.001	0.0690
	AEGIS	0.0600 \pm 0.003	0.0233 \pm 0.001	<u>0.2158</u> \pm 0.028	<u>0.0677</u> \pm 0.007	<u>0.0928</u> \pm 0.057	<u>0.1628</u> \pm 0.010	<u>0.1037</u>
	Ours	0.1602 \pm 0.013	<u>0.0351</u> \pm 0.000	0.6406 \pm 0.026	0.0712 \pm 0.008	0.2430 \pm 0.028	0.1810 \pm 0.012	0.2219

Moreover, our problem setting is the same as existing work on zero-shot image anomaly detection (Jeong et al., 2023; Zhou et al., 2024), with the only difference in the data type used. Considering all these factors, we think “zero-shot” is more suitable for characterizing the nature of the problem complexity and more consistent with the terms/concepts used in the anomaly detection community.

To further demonstrate the difference between inductive learning and zero-shot generalist learning, we adopt the inductive learning methods to the zero-shot setting. Specifically, two representative inductive methods, GraphSAGE (Hamilton et al., 2017) and AEGIS (Ding et al., 2021a) are used and we follow the experimental setup in the main paper to unify the node attribute dimensionality with SVD. The results of GraphSAGE and AEGIS are reported in Table 9.

From the table, we can see whether the general inductive learning method or the inductive anomaly detection method does not achieve promising performance for the zero-shot generalist anomaly detection. This highlights the difference and incompatibility between inductive learning and the problem studied in this paper.