# Raising the Bar in Graph-level Anomaly Detection

**Chen Qiu**[1,2] , **Marius Kloft**[2] , **Stephan Mandt**[3] and **Maja Rudolph**[1]

[1]Bosch Center for Artificial Intelligence
[2]TU Kaiserslautern, Germany
[3]University of California, Irvine, USA

chen.qiu@de.bosch.com, kloft@cs.uni-kl.de, mandt@uci.edu, maja.rudolph@us.bosch.com

## Abstract

Graph-level anomaly detection has become a critical topic in diverse areas, such as financial fraud detection and detecting anomalous activities in social networks. While most research has focused on anomaly detection for visual data such as images, where high detection accuracies have been obtained, existing deep learning approaches for graphs currently show considerably worse performance. This paper raises the bar on graph-level anomaly detection, i.e., the task of detecting abnormal graphs in a set of graphs. By drawing on ideas from self-supervised learning and transformation learning, we present a new deep learning approach that significantly improves existing deep one-class approaches by fixing some of their known problems, including hypersphere collapse and performance flip. Experiments on nine real-world data sets involving nine techniques reveal that our method achieves an average performance improvement of 11.8% AUC compared to the best existing approach.

## 1 Introduction

Anomaly detection (AD) is an important tool for scanning systems for unknown threats. Many web-based systems are best represented by graphs and there has been work on detecting anomalous nodes and edges within a graph (Akoglu *et al.*, 2015). However, in many applications, it is much more relevant to ask whether an entire graph is abnormal.

For example, in a financial network with nodes representing individuals, businesses, and banks and with edges representing transactions, it might be difficult to detect certain criminal activity by looking at individual nodes and edges (Jullum *et al.*, 2020). Clever criminals can hide their intentions behind innocent-looking transactions. However, the entire network associated with a money-laundering scheme is harder to obfuscate and will still exhibit properties of criminal activity. By using tools for graph-level AD, we might be able to detect an entire criminal network rather than flag individual entities.

Unfortunately, there has been limited success in adapting advances in deep anomaly detection to graph-level AD (Zhao
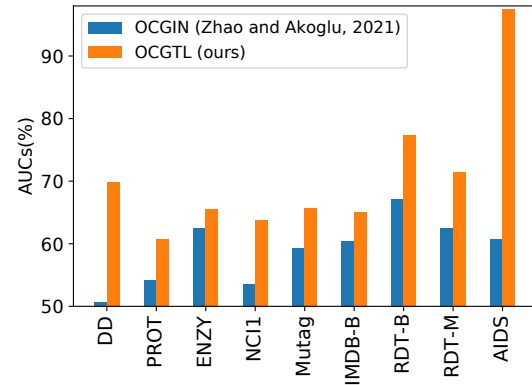


Figure 1: AUC comparison between one-class graph transformation learning (OCGTL) (ours) and One-class GIN (OCGIN) (Zhao and Akoglu, 2021) on several datasets from Sec. 4. OCGTL improves anomaly detection accuracy.

and Akoglu, 2021). Our work addresses this shortcoming. For graph-level AD, we assume to have access to a large dataset of typical graphs, such as a dataset of communities in a social network or a dataset of snapshots of a financial network. All graphs in the training data are considered "normal". The goal is to use the data to learn an anomaly scoring function which can then be used to score how likely it is that a new graph is either normal or abnormal. Importantly, the term graph-level AD refers to detecting *entire* abnormal graphs, rather than localizing anomalies within graphs.

Recently there has been a trend of using deep learning in AD on images (Golan and El-Yaniv, 2018) or tabular and sequential data (Qiu *et al.*, 2021). However, there has been limited research on deep AD for graphs. This may seem surprising since it appears straightforward to adopt a deep AD method for tabular data into one for graphs by defining an appropriate feature map. Yet, Zhao and Akoglu (2021) found that the resulting methods often perform close to random, and so far, attempts to adopt modern AD methods (based on deep learning) to graph-level AD have not been successful.

We develop one-class graph transformation learning (OCGTL), a new model for graph level AD that combines deep one-class classification (OCC) and self-supervision. Fig. 2 provides a sketch of the approach. The OCGTL archi-
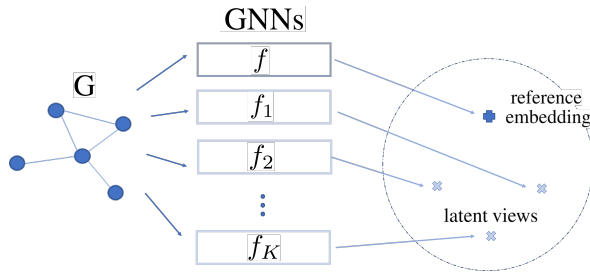
Figure 2: Sketch of the OCGTL procedure. Given a graph (left), we use a set of GNNs to embed the latter into a latent space. The different GNNs embeddings are trained to be both diverse while also being close to a so-called 'reference embedding'. See Sec. 3 for more details.

tecture consists of $K + 1$ graph neural networks (GNNs) that are jointly trained on two complementary deep AD losses. In Sec. 3.2 we prove that this new combined loss mitigates known issues of previous deep AD approaches (Ruff *et al.*, 2018; Zhao and Akoglu, 2021).

Fig. 1 shows that the approach significantly raises the bar in graph-level AD performance. In an extensive empirical study, we evaluate nine methods on nine real-world datasets. Our work brings deep AD on graphs up to speed with other domains, contributes a completely new method (OCGTL), and paves the way for future progress.

In summary, our main contributions are as follows:

- We develop OCGTL, a novel end-to-end method for graph-level AD, which combines the advantages of OCC and neural transformation learning.[1]

- In Sec. 3.2, we prove theoretically that OCGTL is not susceptible to the useless trivial solutions that are optimal under the objective of Zhao and Akoglu (2021).

- In Sec. 4, we study nine methods (four newly developed) on nine real-world graph datasets, ranging from social networks to bioinformatics datasets. We improve the architectures of existing deep approaches to graph-level AD. Yet, OCGTL significantly raises the anomaly detection accuracy over previous work.

## 2   Related Work

**Deep Anomaly Detection.**   Deep AD has received massive attention in various domains (Ruff *et al.*, 2021). Related work on deep AD can be summarized in the following classes. Deep autoencoder variants (Zong *et al.*, 2018) detect anomalies based on the reconstruction error. Deep one-class networks (Ruff *et al.*, 2018) are trained on an OCC objective to map normal data close to a center in an embedding space. They score anomalies based on their distance to the center. Deep generative models detect anomalies based on density estimation (Zhang *et al.*, 2021) or using the discriminator of generative adversarial networks (Deecke *et al.*, 2018). Self-supervised AD methods have achieved great success on im-

---

[1]Code is available at
https://github.com/boschresearch/GraphLevel-AnomalyDetection

ages (Ruff *et al.*, 2021). They rely on a self-supervision task for training and anomaly scoring. The self-supervision tasks often require data augmentation, e.g., for transformation prediction (Golan and El-Yaniv, 2018). Another self-supervised paradigm using data augmentations is contrastive learning (Chen *et al.*, 2020). In contrastive learning for AD, transformed (e.g. rotated) images are treated as negative samples (Sohn *et al.*, 2020). For data types beyond images, handcrafting transformations is challenging. Data-driven neural transformations (Qiu *et al.*, 2021) have shown success in AD.

**Graph Anomaly Detection.**   Finding abnormal nodes or edges within a large graph is widely studied (Akoglu *et al.*, 2015). Deep learning based methods for node- and edge-level AD have been successfully applied to both static (Ding *et al.*, 2019, 2020) and dynamic graphs (Yoon *et al.*, 2019; Zheng *et al.*, 2019). In contrast, deep learning for graph-level AD, the problem we study in this paper, has received less attention. Zhao and Akoglu (2021) first explore how to extend deep OCC for graph-level AD and develop one-class GIN (OC-GIN). They also introduce two-stage graph-level AD frameworks with either graph embedding models or graph kernels. However, all these attempts have not yet produced a solid baseline for graph-level AD. Zhao and Akoglu (2021) report that these methods suffer from "performance flip", which is an issue where the model performs worse than random on at least one experimental variant. We study how to overcome it.

**Deep Learning for Graphs.**   GNNs automate feature extraction from graphs and play an important role in graph classification and representation learning. Graph Convolution Networks (Kipf and Welling, 2016) learn a node representation by aggregating the representations of the node's neighbors. Errica *et al.* (2020) provide a fair comparison of various GNNs for graph classification. An outstanding example is graph isomorphism network (GIN) (Xu *et al.*, 2018a), which is proven to be as powerful as the Weisfeiler-Lehman graph isomorphism test and is the architecture we use for all GNN-based AD methods we study in this paper.

## 3   One-Class Graph Transformation Learning for Anomaly Detection

We propose a new method for graph-level AD. One-class graph transformation learning (OCGTL) combines the complementary strengths of deep OCC (Ruff *et al.*, 2018; Zhao and Akoglu, 2021) and self-supervised AD with learnable transformations (Qiu *et al.*, 2021).

Our new self-supervised approach is designed to overcome known issues of deep OCC. The deep OCC objective is prone to a trivial solution called *hypersphere collapse*. The deep one-class objective encourages all the graph embeddings of the graphs in the training data to concentrate within a hypersphere. This task can be solved perfectly when the feature extractor learns to map all inputs to the center of the hypersphere. Our model provably overcomes hypersphere collapse by regularizing the one-class terms in the objective with a transformation learning term. The resulting model is more flexible (for example, the hypersphere centers can be treated as trainable parameters) and training is more robust, despite the added flexibility.

Zhao and Akoglu (2021) have developed the first deep one-class approach to graph-level AD. In their paper, they report an additional, practical difficulty in graph-level AD which they call the *performance flip issue*. In many of their experiments, their trained model (OCGIN) systematically confuses anomalies with normal samples. The goal of this work is to overcome both hypersphere collapse and performance flip.

Our model consists of an ensemble of GNNs. One of them – the reference feature extractor – produces a reference embedding of its input graph. The other GNN feature extractors produce alternative "latent views" of the graph. The objective of our approach has a one-class term and a transformation learning term. The one-class term aims at concentrating all the latent views within a hyper-sphere in the embedding space. Transformation learning has the competing objective to make each view predictive of the reference embedding. It encourages the latent views to be diverse yet semantically meaningful. By counteracting the one-class term in this manner, hypersphere collapse can be provably avoided.

The tension that arises from satisfying both aspects of the objective has further advantages. In particular, it leads to a harder self-supervision task, which in turn leads to better anomaly detection performance. When the training objective is difficult to satisfy, the trained model has to be more sensitive to typical salient features of normal data. New graphs which do not exhibit these features incur a higher loss and are then more easily detected as anomalies. Also, the two loss contributions focus on different notions of distance between the graph embeddings. The one-class term is based on Euclidean distances, while the transformation learning loss is based on angles between embeddings. With the combined loss as the anomaly score, our method is sensitive to abnormal embedding configurations both in terms of angles between the latent views and in terms of Euclidean distances.

In this section, we first introduce OCGTL and then detail its main ingredients, including self-supervised AD with learnable transformations, deep OCC, and feature extraction with GNNs. We then present the theory behind OCGTL.

### 3.1 Proposed Method - OCGTL

OCGTL combines the best of OCC and neural transformation learning. The OCGTL architecture consists of a reference feature extractor $f$ and $K$ additional feature extractors $f_k$ ($k = 1, \cdots, K$), which are trained jointly as illustrated in Fig. 2. Each of the feature extractors is a parameterized function (e.g. GNN) which takes as input an attributed graph $G = \{\mathcal{V}, \mathcal{E}, \mathcal{X}\}$, with vertex set $\mathcal{V}$, edges $\mathcal{E}$, and node features (attributes) $\mathcal{X} = \{x_v | v \in \mathcal{V}\}$ and maps it into an embedding space $\mathcal{Z}$. These $K + 1$ feature extractors are trained jointly on the OCGTL loss, $\mathcal{L}_{\text{OCGTL}} = \mathbb{E}_G [\mathcal{L}_{\text{OCGTL}}(G)]$. Each graph in the training data contributes two terms to the loss,

$$\mathcal{L}_{\text{OCGTL}}(G) = \mathcal{L}_{\text{OCC}}(G) + \mathcal{L}_{\text{GTL}}(G). \quad (1)$$

The first term, $\mathcal{L}_{\text{OCC}}(G)$, is a one-class term; it encourages all the embeddings to be as close as possible to the same point $\theta \in \mathcal{Z}$. The second term, $\mathcal{L}_{\text{GTL}}$, enforces each GNN's embeddings to be diverse and semantically meaningful representations of the input graph G.

The two terms are presented in detail below.

### The Graph Transformation Learning Term

Neural transformation learning (Qiu *et al.*, 2021) is a self-supervised training objective for deep AD which has seen success on time series and tabular data. Here we generalize the training objective of Qiu *et al.* (2021) (by dropping their parameter sharing constraint) and adapt it to graphs.

For a graph G, the loss of graph transformation learning encourages the embeddings of each GNN, $f_k(G)$, to be similar to the embedding of the reference GNN, $f(G)$, while being dissimilar from each other. Consequently, each GNN $f_k$ is able to extract graph-level features to produce a different view of G. The contribution of each graph to the objective is

$$\mathcal{L}_{\text{GTL}}(G) = -\sum_{k=1}^{K} \log \frac{c_k}{C_k} \quad (2)$$

with

$$c_k = \exp\left(\frac{1}{\tau}\text{sim}(f_k(G), f(G))\right),$$

$$C_k = c_k + \sum_{l \neq k}^{K} \exp\left(\frac{1}{\tau}\text{sim}(f_k(G), f_l(G))\right),$$

where $\tau$ denotes a temperature parameter. The similarity here is defined as the cosine similarity $\text{sim}(z, z') := z^T z' / \|z\|\|z'\|$. Note that the above loss is more general than the one proposed in Qiu *et al.* (2021) as it omits a parameter sharing constraint between transformations. This choice is inspired by the observation in You *et al.* (2020) that different graph categories prefer different types of transformations.

### The One-Class Term

One-class classification (OCC) is a popular paradigm for AD (Noumir *et al.*, 2012). The idea is to map data into a minimal hypersphere encompassing all normal training data. Data points outside the boundary are considered anomalous. The contribution of each graph G to our OCC objective is

$$\mathcal{L}_{\text{OCC}}(G) = \sum_{k=1}^{K} \|(f_k(G) - \theta)\|_2 \quad (3)$$

The loss function penalizes the distance of the graph G to the center $\theta$ which we treat as a trainable parameter. In previous deep OCC approaches, the center $\theta$ has to be a fixed hyperparameter to avoid trivial solutions to Eqn. (3).

### Feature Extraction with GNNs

For graph data, parametrizing the feature extractors $f$ and $f_1, \cdots, f_K$ by GNNs is advantageous. At each layer $l$, a GNN maintains node representation vectors $h_v^{(l)}$ for each node $v$. The representation is computed based on the previous layer's representations of $v$ and its neighbors $\mathcal{N}(v)$,

$$h_v^{(l)} = \text{GNN}^{(l)}\left(h_v^{(l-1)}, h_u^{(l-1)} \mid u \in \mathcal{N}(v)\right). \quad (4)$$

Each layer's node representations are then combined into layer-specific graph representations,

$$h_G^{(l)} = \text{READOUT}^{(l)}\left(h_v^{(l)} \mid v \in G\right), \quad (5)$$

which are concatenated into graph-level representations,

$$h_{\mathrm{G}} = \mathrm{CONCAT}\left(h_{\mathrm{G}}^{(l)} \mid l = 1, ..., L\right). \qquad (6)$$

This concatenation introduces information from various hierarchical levels (Xu *et al.*, 2018b) into the graph representation. Our empirical study in Sec. 4 shows that the choice of the readout function (which determines how the node representations are aggregated into graph representations) is particularly important to detect anomalies reliably.

**Anomaly Scoring with OCGTL**
OCGTL is an end-to-end methods for graph-level AD. During training the GNNs are trained on Eqn. (1). During test, $\mathcal{L}_{\mathrm{OCGTL}}$ (Eqn. (1)) is used directly as the score function for detecting anomalous graphs. A low loss on a test sample means that the graph is likely normal, whereas a high loss is indicative of an anomaly. One advantage of OCGTL is that its loss makes it more sensitive to different types of anomalies by considering both angles between embeddings and Euclidean distances. In contrast, OCC-based methods typically rely on the Euclidean distance only.

Another advantage of OCGTL over OCC-based approaches is that its training is more robust and the AD model can be more flexible. We prove this next.

## 3.2 A Theory of OCGTL

A known difficulty for training OCC-based deep anomaly detectors (such as deep SVDD and OCGIN) is *hypersphere collapse* (Ruff *et al.*, 2018). Hypersphere collapse is a trivial optimum of the training objective

$$\mathcal{L}_{(\mathrm{Ruff}\,et\,al.,\,2018)}(\mathrm{G}) = ||f(\mathrm{G}) - \theta||_2^2, \qquad (7)$$

which occurs when the feature extractor $f$ maps all inputs exactly into the center $\theta$. The hypersphere then has a radius of zero, and AD becomes impossible. Ruff *et al.* (2018) recommend fixing $\theta$ and avoiding bias terms for $f$ and show good results in practice. However, there is no guarantee that a trivial solution can be avoided under any architecture for $f$. Here we prove that OCGTL overcomes this.

We first show that our one-class term (Eqn. (3)) is also prone to hypersphere collapse when all the feature extractors are constant. However, we then show that this trivial solution for minimizing Eqn. (3) is not optimal under the OCGTL loss. Our method provably avoids hypersphere collapse even when the center $\theta$ is a trainable parameter. This result makes OCGTL the first deep one-class approach where the center can be trained.

**Proposition 1.** *The constant feature extractors, $f_k(G) = \theta$ for all $k$ and all inputs G, minimize $\mathcal{L}_{OCC}$ (Eqn. (3)).*

*Proof.* $0 \leq \mathcal{L}_{\mathrm{OCC}}$ is the squared $\ell_2$ norm of the distance between the embedding of G and the center $\theta$. Plugging in $f_k(\mathrm{G}) = \theta$ attains the minimum 0. $\qquad \square$

In contrast, regularization with transformation learning can avoid hypersphere collapse. Under the constant encoder, all the latent views are the same and hence at least as close to each other as to the reference embedding, leading to $\mathcal{L}_{\mathrm{GTL}} \geq K \log K$. However, the transformation learning objective

aims at making the views predictive of the reference embeddings, in which case $\mathcal{L}_{\mathrm{GTL}} < K \log K$. The following proposition shows that if there is a parameter setting which achieves this, the constant feature extractors do not minimize the OCGTL loss which proves that hypersphere collapse can be avoided.

**Proposition 2.** *If there exists a parameter setting such that $\mathcal{L}_{GTL} < K \log K$ on the training data, then the constant feature extractors $f_k(G) = \theta$ do not minimize the combined loss $\mathcal{L}_{OCGTL}$ (Eqn. (1)).*

*Proof.* For constant feature extractors $f_k(\mathrm{G}) = \theta$, $\mathcal{L}_{\mathrm{OCGTL}} = \mathcal{L}_{\mathrm{GTL}} \geq K \log K$, where $K$ is the number of transformations and $K \log K$ is the negative entropy of randomly guessing the reference embedding. Assume there is a constellation of the model parameters s.t. $\mathcal{L}_{\mathrm{GTL}} < K \log K$. Since $\theta$ is trainable, we can set it to be the origin. The loss of the optimal solution is at least as good as the loss with $\theta = 0$. Set $\epsilon = K \log K - \mathcal{L}_{\mathrm{GTL}}$. The encoders can be manipulated such that their outputs are rescaled and as a result all the embeddings have norm $||f_k(\mathrm{G})||_2 < \epsilon/K$. As the norm of the embeddings changes, $\mathcal{L}_{\mathrm{GTL}}$ remains unchanged since the cosine similarity is not sensitive to the norm of the embeddings. By plugging this into Eqn. (1) we get $\mathcal{L}_{\mathrm{OCGTL}} = \sum_{k=1}^{K} ||f_k(\mathrm{G})||_2 + \mathcal{L}_{\mathrm{GTL}} < K \log K$, which is better than the performance of the best set of constant encoders. $\qquad \square$

Props. 1 and 2 demonstrate that our method is the first deep one-class method not prone to hypersphere collapse. The assumption of Prop. 2, that $\mathcal{L}_{\mathrm{GTL}} < K \log K$ can be tested in practice by training graph transformation learning (GTL) and evaluating the predictive entropy on the training data. In all scenarios we worked with $\mathcal{L}_{\mathrm{GTL}} << K \log K$ after training.

## 3.3 Newly Developed Baselines

The main contribution of our work is OCGTL. To study the effectiveness of OCGTL we implement the following graph-level AD methods as ablations. These methods have not been studied on graphs before, so their implementation is also one of our contributions that paves the way for future progress.

**One-class pooling (OCPool).** As a shallow method, OCPool uses pooling to construct a graph representation:

$$h_{\mathrm{G}} = \mathrm{POOLING}\left(x_v \mid v \in \mathrm{G}\right). \qquad (8)$$

This feature extractor does not have parameters and hence requires no training. Anomalies can be detected by training an one-class SVM (OCSVM) (Manevitz and Yousef, 2001) on these features. This novel approach for graph-level AD is a simple baseline and achieves solid results in our empirical study (even though it does not use the edge sets $\mathcal{E}$ of the graphs). Another reason for studying OCPool, is that it helps us understand which pooling function might work best as readout function (Eqn. (5)) for GNN-based AD methods.

**Graph transformation prediction (GTP).** GTP is an end-to-end self-supervised detection method based on transformation prediction. It trains a classifier $f$ to predict which transformation has been applied to a samples and uses the cross-entropy loss to score anomalies. We implement GTP with

six graph transformations (node dropping, edge adding, edge dropping, attribute masking, subgraph, and identity transformation) originally designed in You *et al.* (2020).

**Graph transformation learning (GTL).** GTL is an end-to-end self-supervised detection method using neural transformations (Qiu *et al.*, 2021). $K$ GNNs, $f_k$ for $k = 1, \cdots, K$ in addition to the reference feature extractor $f$ are trained on $\mathcal{L}_{\text{GTL}}$ (Eqn. (2)). The loss is used directly to score anomalies. While this method works well in practice, it is not sensitive to the norm of the graph embeddings in Eqn. (2). The normalization step in computing the cosine similarity makes mean and add pooling equivalent when aggregating the graph representations. This may put GTL at a disadvantage compared to the other methods, which profit from add pooling.

# 4 Experiments

This section details our empirical study. We benchmark nine algorithms on nine real-world graph classification datasets from different domains using various evaluation measures. First, we describe the datasets and how the AD benchmark is set up. Second, we present all methods we compare, including baselines and their implementation details. Third, the evaluation results are presented and analyzed. In summary, OCGTL achieves the best performance on real-world datasets from various domains and raises the anomaly detection accuracy significantly ($+11.8\%$ in terms of AUC on average compared to OCGIN of Zhao and Akoglu (2021)). Finally, we present our findings about preferable design choices that are also beneficial for other deep methods in graph-level AD.

## 4.1 Datasets and Experimental Setting

We benchmark nine methods on nine graph classification datasets that are representative of three domains. In addition to financial and social networks security, health organizations need an effective graph-level AD method to examine proteins (represented as graphs) to monitor the spread and evolution of diseases. Targeting these application domains, we study three bioinformatics datasets: DD, PROTEINS, and ENZYMES, three molecular datasets: NCI1, AIDS, and Mutagenicity, and three datasets of social networks: IMDB-BINARY, REDDIT-BINARY, and REDDIT-MULTI-5K. The datasets are made available by Morris *et al.* (2020), and the statistics of the datasets are given in Appendix A.

We follow the standard setting of previous work to construct an AD task from a classification dataset (Ruff *et al.*, 2018; Golan and El-Yaniv, 2018; Zhao and Akoglu, 2021). A classification dataset with $N$ classes produces $N$ experimental variants. In each experimental variant, one of the classes is treated as "normal"; the other classes are considered as anomalies. The training set and validation set only contain normal samples, while the test set contains a mix of normal samples and anomalies that have to be detected during test time. For each experimental variant, $10\%$ of the normal class is set aside for the test set, and $10\%$ of each of the other classes is added to the test set as anomalies. (The resulting fraction of anomalies in the test set is proportional to the class balance in the original dataset. The remaining $90\%$ of the normal class is used for training and validation. We use

10-fold cross-validation to estimate the model performance. In each fold, $10\%$ of the training set is held out for validation. We train each model three times separately and average the test results of three runs to get the final test results in each fold. Training multiple times ensures a fair comparison as it favors methods that are robust to the random initialization.

**Evaluation.** Results will be reported in terms of the area under the ROC curve (AUC) (%), averaged over 10 folds with standard deviation. We also report the results in terms of F1-score in Appendix C. In addition, all methods will be evaluated in terms of their susceptibility to performance flip.

Zhao and Akoglu (2021) coined the term "performance flip" for AD benchmarks derived from binary classification datasets. We generalize their definition to multiple classes:

**Definition 1.** *(Performance flip.) A model suffers from performance flip on an anomaly detection benchmark derived from a classification dataset if it performs worse than random on at least one experimental variant.*

## 4.2 Baselines and Implementation Details

Many deep AD approaches that have achieved impressive results in other domains have not yet been adapted to graph-level AD. There has been no comprehensive study of various GNN-based graph-level AD approaches. An additional contribution of our work is that we adapt recent advances in deep AD to graphs. In our empirical study, we compare OCGTL both to GNN-based methods and to non-GNN-based methods, which we outline below.

**GNN-based Baselines.** Our study includes OCGTL, OCGIN (Zhao and Akoglu, 2021) and the self-supervised approaches graph transformation prediction (GTP) and graph transformation learning (GTL) described in Sec. 3.3. We use GIN as the feature extractor for all GNN-based baselines to compare with OCGIN fairly. In particular, we use 4 GIN layers, each of which includes a two-layer MLP and graph normalization (Cai *et al.*, 2020). The dimension of the node representations is 32. The readout function of almost all methods consists of a two-layer MLP and then an add pooling layer. In GTP, the final prediction is obtained by summing the layer-wise predictions, and the readout function is composed of an add pooing layer followed by a linear layer. In GTP, we employ six hand-crafted transformations. For a fair comparison, GTL and OCGTL use six learnable graph transformations in all experiments. Additional hyperparameter settings are recorded for reproducibility in Appendix B.

**Improved Implementation of OCGIN.** With these implementation details for the GNNs we can significantly improve the performance of OCGIN over the implementation in Zhao and Akoglu (2021). For this reason, our empirical study includes both OCGIN (the original version with mean pooling and batch normalization) and OCGIN[†] (our improved version with add pooling and graph normalization).

**Non-GNN-based Baselines.** Besides OCPool, we include four two-stage detection methods proposed by Zhao and Akoglu (2021). Two of them use unsupervised graph embedding methods, Graph2Vec (G2V) (Narayanan *et al.*, 2017) or

|  |  | DD | PROT | ENZY | NCI1 | AIDS | Mutag | Rank |
|---|---|---|---|---|---|---|---|---|
| Baselines (prev. work) | WLK | 50.2±0.3* | 49.7±0.5* | 52.1±2.0* | 49.6±0.4* | 51.3±0.8* | 52.3±0.6 | 8.1±1.7 |
|  | PK | 51.2±2.3* | 50.8±1.5* | 51.3±1.2* | 51.4±1.7* | 59.5±2.3* | 52.5±1.6 | 8.4±1.1 |
|  | G2V | 49.5±2.2* | 53.2±3.0* | 52.0±3.9* | 50.5±0.7* | 48.4±0.8* | 50.2±0.7* | 8.9±1.0 |
|  | FGSD | 66.0±2.3 | 58.5±1.8* | 52.7±3.4* | 55.4±0.7 | 91.6±3.7 | 51.3±0.8* | 5.3±2.9 |
|  | OCGIN | 50.7±1.2* | 54.2±1.2* | 62.4±2.7 | 53.6±1.3* | 60.8±2.2* | 59.3±1.4 | 5.4±1.7 |
| Ablations (ours) | OCGIN† | 61.4±1.6 | 57.2±2.3 | 63.5±3.9 | 62.2±1.3 | **97.5±2.0** | 61.5±1.8* | 2.7±0.9 |
|  | OCPool | 61.1±3.3* | **61.9±2.1** | 53.1±2.9* | 57.0±1.3 | **97.6±1.4** | 53.2±0.5* | 4.3±2.2 |
|  | GTP | 54.2±1.9 | **61.9±2.9** | 55.0±2.0* | 55.3±1.2* | 77.2±2.9 | 54.7±1.8* | 4.8±1.5 |
|  | GTL | 51.7±0.9* | 56.2±2.5* | 60.4±1.6 | 59.8±1.0* | 67.8±3.3* | 61.8±1.0 | 3.8±1.7 |
| ours | OCGTL | **69.9±2.6** | 60.7±2.4 | **65.5±3.8** | **63.7±1.2** | 97.5±2.0 | **65.7±2.1** | 1.4±0.7 |

† OCGIN is the original implementation from Zhao and Akoglu (2021), while OCGIN† denotes our improved implementation with the same GIN architecture choices (add pooling etc.) as OCGTL, GTP, and GTL.

Table 1: Average AUCs (%) with standard deviations of 9 methods on 6 of 9 datasets. (For the remaining 3 datasets, see Tab. 3.) The performance rank averaged on all nine datasets is provided in the last column. Results marked * perform worse than random on at least one experimental variant (performance flip). OCGTL outperforms the other methods and has no performance flip.

|  | DD |  | PROT |  | NCI1 |  | AIDS |  |
|---|---|---|---|---|---|---|---|---|
| Outlier class | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
| OCGIN | 26.3±2.7 | 75.2±3.4 | 42.5±4.4 | 65.9±4.5 | 64.4±2.5 | 42.9±3.0 | 26.0±3.9 | 95.5±2.2 |
| OCGTL (ours) | 66.8±4.6 | 73.0±2.9 | 63.2±5.4 | 58.1±6.1 | 71.2±3.0 | 56.2±2.5 | 99.3±0.9 | 95.7±3.6 |

Table 2: Average AUCs (%) with standard deviations of OCGIN (Zhao and Akoglu, 2021) and OCGTL (ours) on both experimental variants of four datasets, where the performance flip is observed. The results that are worse than random are marked in red. OCGIN suffers from performance flip, while OCGTL not.

|  |  | IMDB-B | RDT-B | RDT-M |
|---|---|---|---|---|
| Baselines (prev. work) | WLK | 62.0±2.0 | 50.2±0.2* | 50.3±0.1* |
|  | PK | 53.5±2.0 | 50.0 | 50.0 |
|  | G2V | 54.3±1.6 | 51.5±0.6* | 50.0±0.2* |
|  | FGSD | 57.1±1.8 | - | - |
|  | OCGIN | 60.4±2.8 | 67.1±3.5 | 62.4±1.3 |
| Ablations (ours) | OCGIN† | 63.7±2.4 | 74.5±3.4 | 70.4±1.5 |
|  | OCPool | 56.5±0.8 | 65.3±2.2 | 62.4±0.9 |
|  | GTP | 57.6±1.1 | 64.4±2.2 | 62.4±1.3 |
|  | GTL | **65.2±1.9** | 71.6±2.3 | 67.8±0.9 |
| ours | OCGTL | **65.1±1.8** | **77.4±1.9** | **71.5±1.1** |

† OCGIN is the original implementation from Zhao and Akoglu (2021), while OCGIN† denotes our improved version.

Table 3: Average AUCs (%) with standard deviations of nine methods on three datasets to complement Tab. 1.

FGSD (Verma and Zhang, 2017), to extract graph-level representations. The other two of them make use of graph kernels (Weisfeiler-Leman subtree kernel (WLK) (Shervashidze *et al.*, 2011) or propagation kernel (PK) (Neumann *et al.*, 2016)), which measure the similarity between graphs. For all two-stage detection baselines, we use OCSVM (with $\nu = 0.1$) as the downstream outlier detector.

The number of iterations specifies how far neighborhood information can be propagated. By setting the number of iterations to 4, we get a fair comparison to the GNN-based methods, which all have 4 GNN layers. All other hyperparameters correspond to the choices in Zhao and Akoglu (2021).

### 4.3 Experimental Results

**Summary.** We compare OCGTL with all existing baselines on nine real-world datasets. The detection results in terms of average AUC (%) with standard deviation are reported in Tabs. 1 and 3. The results in terms of F1-score are reported in Appendix C. We can see that OCGTL achieves competitive results on all datasets and has the best average rank of 1.4. On average over nine datasets, OCGTL outperforms OCGIN of Zhao and Akoglu (2021) by 11.8%. We can conclude that OCGTL raises the detection accuracy in graph-level AD on various application domains significantly, namely by 9.6% on the bioinformatics domain, by 17.7% on the molecular domain, and by 8% on the social-networks domain.

Moreover, methods with performance flip are marked with a * in Tab. 1. In Tab. 2 we report the results of OCGIN (Zhao and Akoglu, 2021) and our method OCGTL on both experimental variants of datasets where the performance flip is observed. We can see that all existing baselines suffer from the performance flip issue, while OCGTL is the only model without performance flip on any of the datasets.

**Ablation Study of Methods.** Here we discuss the results in Tab. 1 from the perspective of an ablation study to understand if and how the advantages of combing deep OCC and neural transformation learning complement each other. From results in Tab. 1, we can see that OCGTL improves over OCGIN† (with our improved implementation) on 8 of 9 datasets by adding $\mathcal{L}_{GTL}$ as the regularization term and outperforms GTL on 8 of 9 datasets by utilizing the Euclidean distance for detection. We can conclude, that the two terms in the loss function of OCGTL complement each other and offer two metrics for detecting anomalies. As a result, OCGTL consistently outperforms OCGIN and GTL. This is aligned with our theoretical results in Sec. 3.2.

GTP applies hand-crafted graph transformations. Its performance varies across datasets since it is sensitive to the choice of transformations. Even though it works well on the
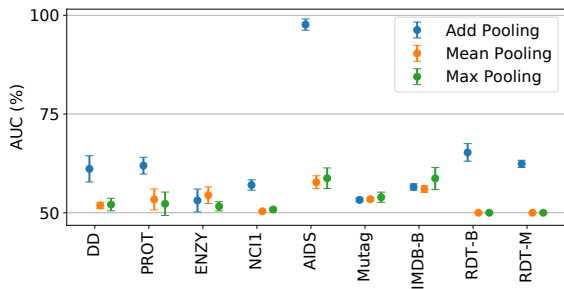
Figure 3: OCPool with add pooling (blue) outperforms alternative choices (mean (orange), max (green)). The results on nine datasets are reported in terms of AUC (%).

PROTEINS dataset with the chosen graph transformations, its performance on other datasets is not competitive to OCGTL. Finding the right transformations for each dataset requires domain knowledge, which is not the focus of this work. In comparison, OCGTL learns data-specific transformations automatically and performs consistently well on various datasets.

**Study of Design Choices.** To raise the bar in deep AD on graphs, we have to understand the impact of the design choices associated with the GNN architecture. Here we discuss the type of pooling layer for the readout function (Eqn. (5)) and normalization of the GNN layers.

First, we study the impact of different pooling layers. OCPool is the ideal testbed to compare add pooling, mean pooling, and max pooling due to its simplicity. Results of running the three options on nine datasets are reported in Fig. 3. Add pooling outperforms the other options. This may result from add pooling injecting information about the number of nodes into the graph representations. OCPool with add pooling is a simple yet effective method for AD. It provides a simple heuristic for aggregating node attributes into graph representations. As shown in Tab. 1, it does well on many datasets (particularly on the PROTEINS and AIDS datasets), even though it does not account for graph structure (edges).

Second, to study the combined impact of the pooling layer and normalization layers on the *deep* GNN-based methods, we compare add pooling with graph normalization (AP + GN) and mean pooling with batch normalization (MP + BN). In Fig. 4, we visualize in a scatter plot the performance of all GNN-based methods on all nine datasets, contrasting the AP + GN result with the MP + BN result in terms of average AUCs (%) with standard deviations. Almost all points fall above the diagonal, meaning that AP + GN is preferable to MP + BN, which has been the design choice of Zhao and Akoglu (2021) for OCGIN. With the new design choices, we are able to significantly raise the bar in graph-level AD.

## 5 Conclusion

We develop a novel end-to-end graph-level AD method, OCGTL, that combines the best of deep OCC and neural transformation learning. OCGTL mitigates the shortcomings of deep OCC by using graph transformation learning as regularization and complements graph transformation learning by introducing a sensitivity to the norm of the graph representa-
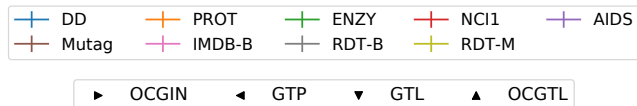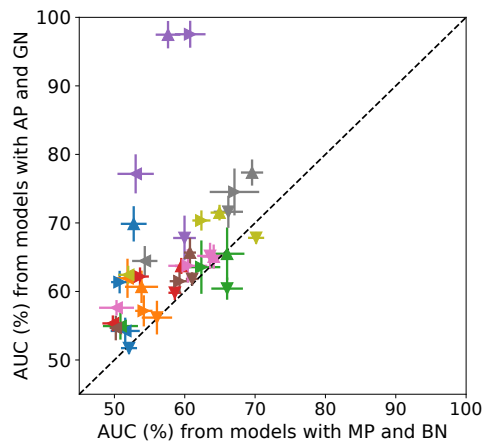


Figure 4: A comparison of two design choices in *deep* GNN-based methods, namely, add pooling with graph normalization (AP + GN, y-axis) and mean pooling with batch normalization (MP + BN, x-axis). Each point compares the detection results of the two variants of one model on one dataset. Most points falling above the diagonal indicate that AP + GN is the preferred design choice for GNN-based AD methods.

tions. Our comprehensive empirical study supports our claim and theoretical results. It shows that OCGTL performs best in various challenging domains and is the only method not struggling with performance flip.

## Acknowledgements

## References

Leman Akoglu, Hanghang Tong, and Danai Koutra. Graph based anomaly detection and description: a survey. *Data mining and knowledge discovery*, 29(3):626–688, 2015.

Tianle Cai, Shengjie Luo, Keyulu Xu, Di He, Tie-yan Liu, and Liwei Wang. Graphnorm: A principled approach to accelerating graph neural network training. *arXiv preprint arXiv:2009.03294*, 2020.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learn-

ing of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020.

Lucas Deecke, Robert Vandermeulen, Lukas Ruff, Stephan Mandt, and Marius Kloft. Image anomaly detection with generative adversarial networks. In *ECML and KDD*, pages 3–17. Springer, 2018.

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*, pages 594–602. SIAM, 2019.

Kaize Ding, Jundong Li, Nitin Agarwal, and Huan Liu. Inductive anomaly detection on attributed networks. In *IJ-CAI*, pages 1288–1294, 2020.

Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. A fair comparison of graph neural networks for graph classification. In *ICLR*, 2020.

Izhak Golan and Ran El-Yaniv. Deep anomaly detection using geometric transformations. In *Neural Information Processing Systems*, pages 9781–9791, 2018.

Martin Jullum, Anders Løland, Ragnar Bang Huseby, Geir Ånonsen, and Johannes Lorentzen. Detecting money laundering transactions with machine learning. *Journal of Money Laundering Control*, 2020.

Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.

Larry M Manevitz and Malik Yousef. One-class svms for document classification. *Journal of machine Learning research*, 2(Dec):139–154, 2001.

Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. In *ICML GRL+ Workshop*, 2020.

Annamalai Narayanan, Mahinthan Chandramohan, Rajasekar Venkatesan, Lihui Chen, Yang Liu, and Shantanu Jaiswal. graph2vec: Learning distributed representations of graphs. *arXiv preprint arXiv:1707.05005*, 2017.

Marion Neumann, Roman Garnett, Christian Bauckhage, and Kristian Kersting. Propagation kernels: efficient graph kernels from propagated information. *Machine Learning*, 102(2):209–245, 2016.

Zineb Noumir, Paul Honeine, and Cedue Richard. On simple one-class classification methods. In *2012 IEEE International Symposium on Information Theory Proceedings*, pages 2022–2026. IEEE, 2012.

Chen Qiu, Timo Pfrommer, Marius Kloft, Stephan Mandt, and Maja Rudolph. Neural transformation learning for deep anomaly detection beyond images. In *ICML*, pages 8703–8714. PMLR, 2021.

Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In *ICML*, pages 4393–4402. PMLR, 2018.

Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*, 2021.

Nino Shervashidze, Pascal Schweitzer, Erik Jan Van Leeuwen, Kurt Mehlhorn, and Karsten M Borgwardt. Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research*, 12(9), 2011.

Kihyuk Sohn, Chun-Liang Li, Jinsung Yoon, Minho Jin, and Tomas Pfister. Learning and evaluating representations for deep occ. In *ICLR*, 2020.

Saurabh Verma and Zhi-Li Zhang. Hunt for the unique, stable, sparse and fast feature learning on graphs. In *Neural Information Processing Systems*, pages 87–97, 2017.

Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *ICLR*, 2018.

Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In *ICML*, pages 5453–5462. PMLR, 2018.

Minji Yoon, Bryan Hooi, Kijung Shin, and Christos Faloutsos. Fast and accurate anomaly detection in dynamic graphs with a two-pronged approach. In *Knowledge Discovery & Data Mining*, pages 647–657, 2019.

Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph contrastive learning with augmentations. *Neural Information Processing Systems*, 33:5812–5823, 2020.

Lily Zhang, Mark Goldstein, and Rajesh Ranganath. Understanding failures in out-of-distribution detection with deep generative models. In *ICML*, pages 12427–12436. PMLR, 2021.

Lingxiao Zhao and Leman Akoglu. On using classification datasets to evaluate graph outlier detection: Peculiar observations and new insights. *Big Data*, 2021.

Li Zheng, Zhenpeng Li, Jian Li, Zhao Li, and Jun Gao. Addgraph: Anomaly detection in dynamic graph using attention-based temporal gcn. In *IJCAI*, pages 4419–4425, 2019.

Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Daeki Cho, and Haifeng Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *ICLR*, 2018.

# Appendix

## A    Datasets Statistics

We select nine graph classification datasets for the evaluation from three domains (bioinformatics, molecules, and social networks). For each dataset, we report the number of graphs, the dimension of node attributes, the average number of nodes, and the average number of edges in each class as the statistics. We list the statistics of these nine datasets in Tab. 4.

Table 4: The statistics of used datasets. We report the number of graphs, the dimension of node attributes, the average number of nodes, and the average number of edges for each class in each dataset.

| Dataset | Category | Class | #Graphs | #NodeAttrs | Avg.#Nodes | Avg.#Edges |
|---|---|---|---|---|---|---|
| DD | Bioinformatics | 0 | 691 | 89 | 355.2 | 1806.6 |
|  |  | 1 | 487 | 89 | 183.7 | 898.9 |
| PRTOEINS | Bioinformatics | 0 | 663 | 3 | 50.0 | 188.1 |
|  |  | 1 | 450 | 3 | 22.9 | 83.0 |
| ENZYMES | Bioinformatics | 0 | 100 | 3 | 36.2 | 132.7 |
|  |  | 1 | 100 | 3 | 29.9 | 113.8 |
|  |  | 2 | 100 | 3 | 28.9 | 111.2 |
|  |  | 3 | 100 | 3 | 38.2 | 148.8 |
|  |  | 4 | 100 | 3 | 31.4 | 119.6 |
|  |  | 5 | 100 | 3 | 31.2 | 119.6 |
| NCI1 | Molecules | 0 | 2053 | 37 | 25.7 | 55.3 |
|  |  | 1 | 2057 | 37 | 34.1 | 73.9 |
| AIDS | Molecules | 0 | 400 | 38 | 37.6 | 80.5 |
|  |  | 1 | 1600 | 38 | 10.2 | 20.4 |
| Mutagenicity | Molecules | 0 | 2401 | 14 | 29.4 | 60.6 |
|  |  | 1 | 1936 | 14 | 31.5 | 62.7 |
| IMDB-B | Social networks | 0 | 500 | 136 | 20.1 | 193.6 |
|  |  | 1 | 500 | 136 | 19.4 | 192.6 |
| REDDIT-B | Social networks | 0 | 1000 | 1 | 641.3 | 1471.9 |
|  |  | 1 | 1000 | 1 | 218.0 | 519.1 |
| REDDIT-M | Social networks | 0 | 1000 | 1 | 799.5 | 2035.5 |
|  |  | 1 | 1000 | 1 | 852.1 | 1940.4 |
|  |  | 2 | 1000 | 1 | 374.1 | 856.5 |
|  |  | 3 | 1000 | 1 | 249.6 | 534.0 |
|  |  | 4 | 1000 | 1 | 267.0 | 581.7 |

$^\star$ In IMDB-B, the one-hot degree is used as node attributes. In REDDIT-B and REDDIT-M, the constant one is used as node attributes.

## B    Additional Implementation Details

**Training hyperparameters:** We use the Adam optimizer with an initial learning rate of 0.001 and decay the learning rate by 0.5 every 100 epochs. We set the maximum epochs as 500 and the batch size as 128. We use the early stopping based on validation loss (without access to the true labeled anomalies) for the training. The early stopping is implemented with a patience parameter of 100 epochs to ease the sensitivity to fluctuations in the validation loss. An early stopping without access to the true anomalies is critical for an unbiased model evaluation in the AD tasks.

   **Hardware and Software:** All models are trained in our GPU cluster, which consists of NVIDIA GeForce GTX TITAN X GPUs, and NVIDIA TITAN X Pascal GPUs. All GNN-based models have been implemented by means of the Pytorch Geometrics library. The implementations of G2V and FGSD are from Karate Club library, while the implementations of WLK and PK are from GraKel library.

## C    Additional Experimental Results

Here we provide additional experimental results, which include two additional comparisons of design choices in Fig. 5 and the evaluation results of nine methods on nine datasets in terms of F1-score in Tab. 5.

Table 5: Average F1-scores with standard deviations of 9 methods on 9 datasets. OCGTL performs best generally.

| Datasets | DD | PROT | ENZY | NCI1 | AIDS | Mutag | IMDB-B | RDT-B | RDT-M |
|---|---|---|---|---|---|---|---|---|---|
| WLK | 0.56±0.016 | 0.58±0.036 | 0.84±0.004 | 0.49±0.008 | 0.48±0.008 | 0.52±0.01 | 0.57±0.032 | 0.5±0.007 | 0.8±0.001 |
| PK | 0.54±0.019 | 0.57±0.025 | 0.83±0.002 | 0.51±0.017 | 0.49±0.016 | 0.52±0.014 | 0.52±0.018 | 0 | 0 |
| G2V | 0.48±0.026 | 0.5±0.036 | 0.84±0.01 | 0.5±0.007 | 0.76±0.032 | 0.5±0.006 | 0.53±0.024 | 0.5±0.006 | 0.8±0.001 |
| FGSD | 0.63±0.019 | 0.59±0.025 | 0.83±0.01 | 0.54±0.009 | 0.95±0.021 | 0.52±0.009 | 0.56±0.018 | - | - |
| OCGIN | 0.55±0.022 | 0.55±0.021 | 0.86±0.007 | 0.52±0.013 | 0.49±0.01 | 0.57±0.007 | 0.57±0.025 | 0.63±0.033 | 0.82±0.003 |
| OCGIN† | 0.59±0.018 | 0.56±0.029 | 0.86±0.011 | 0.59±0.013 | **0.97±0.012** | 0.60±0.017 | 0.60±0.02 | 0.68±0.031 | 0.84±0.004 |
| OCPool | 0.59±0.028 | **0.61±0.024** | 0.84±0.008 | 0.55±0.009 | **0.97±0.017** | 0.52±0.009 | 0.53±0.012 | 0.66±0.023 | 0.82±0.002 |
| GTP | 0.52±0.014 | 0.59±0.028 | 0.84±0.005 | 0.54±0.011 | 0.63±0.033 | 0.54±0.018 | 0.55±0.013 | 0.61±0.016 | 0.82±0.003 |
| GTL | 0.57±0.020 | 0.60±0.036 | 0.86±0.007 | 0.57±0.009 | 0.53±0.024 | 0.60±0.007 | 0.60±0.024 | 0.66±0.02 | 0.84±0.003 |
| OCGTL | **0.66±0.025** | 0.60±0.025 | **0.87±0.012** | **0.60±0.011** | 0.97±0.011 | **0.63±0.02** | **0.62±0.019** | **0.7±0.02** | **0.85±0.004** |

[†] OCGIN is the original implementation from Zhao and Akoglu (2021), while OCGIN† denotes our improved implementation.



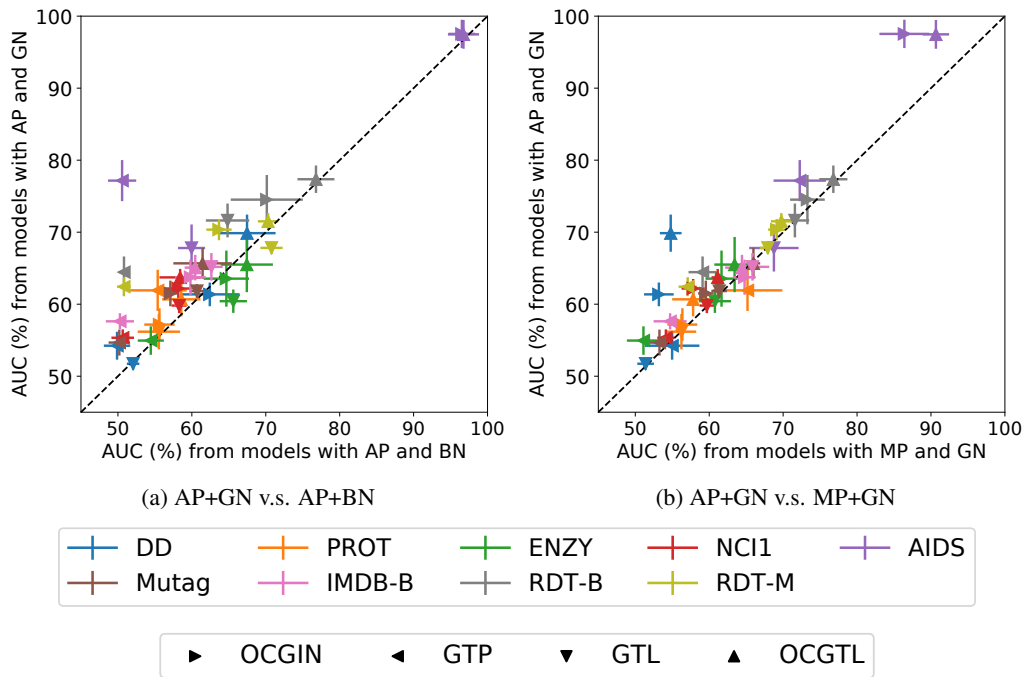(a) AP+GN v.s. AP+BN

(b) AP+GN v.s. MP+GN

Figure 5: A comparison of design choices in *deep* GNN-based methods. (a) Results of AP+BN on x-axis against AP+GN on y-axis. (b) Results of MP+GN on x-axis against AP+GN on y-axis. In conclusion, add pooling with graph normalization is preferable.