# Machine Learning with Graphs Final Project Report

Wenhao Xu (wx2056)

May 2024

### Abstract

Previous research focused on the relationship between numerical node features and graph structure, enhancing GNN performance but remaining graph-agnostic. GIANT-XRT leverages XMC and XR-Transformers for improved performance on large datasets. We propose methods to concatenate multiple datasets and employ contrastive learning to address the limitations of single-dataset training and the time-consuming prediction procedure. Our empirical results demonstrate significant improvements in efficiency and effectiveness, highlighting the potential impact of these enhancements.

## 1 Introduction

The research problem addressed in this project is the enhancement of the GIANT-XRT model for node feature extraction in graph-based machine learning. Specifically, the project aims to address two main challenges faced by the previous work: (1) training the model on a single dataset at a time, limiting its ability to leverage multiple datasets, and (2) the time-consuming prediction procedure using Extreme Multi-label Classification (XMC) to predict node similarities in the adjacency matrix.

This research problem is crucial as it directly impacts the efficiency and effectiveness of graph-based machine learning models. By addressing these challenges, we can significantly improve the scalability, generalization, and applicability of the GIANT-XRT model to real-world datasets and tasks.

## 2 Related Work

Previous work by Chien et al. [1] proposed a method for node feature extraction using self-supervised multi-scale neighborhood prediction. Several relevant works have explored graph-based machine learning models for node feature extraction and graph representation learning. These include models such as Graph Convolutional Networks (GCNs), Graph Attention Networks (GATs),

and Variational Graph Autoencoders (VGAEs). While these models have shown promising results in various applications, they typically focus on single-dataset settings and may not fully address the challenges specific to the GIANT-XRT model.

## 2.1 Limitation of Prior Work

Despite the effectiveness of related works, they have several limitations that motivate the proposed enhancements to the GIANT-XRT model. Firstly, existing models are often trained on a single dataset at a time, limiting their ability to leverage diverse datasets and learn more robust representations. Secondly, the prediction procedure using XMC for node similarity computation is computationally intensive and time-consuming, especially for large-scale graphs.
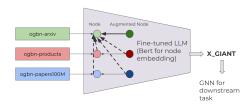
# 3 The Proposed Method



Figure 1: Concept Figure Illustrating multi-dataset and contrastive learning

Dataset aspect: Develop methods to concatenate and integrate multiple datasets with text attributes, allowing the model to learn from diverse sources of information simultaneously.

Prediction procedure: Investigate alternative approaches to the time-consuming XMC prediction procedure, such as contrastive learning, which samples nodes and their augmentations as positive and negative samples for similarity comparison.

## 3.1 Dataset Concatenation

To leverage the rich information from multiple datasets, we concatenate node raw text files from various datasets:

```
cat ./proc_data_xrt/ogbn-arxiv/X.all.txt \
    ./proc_data_xrt/ogbn-products/X.all.txt \
    ./proc_data_xrt/ogbn-papers100M/X.all.txt > \
    ./proc_data_xrt/concatenated/X.all.txt
```

This approach allows the model to learn from diverse sources of information simultaneously, enhancing its ability to generalize across different types of graphs.

## 3.2 Hierarchical-XTransformer for XMC

We utilize a hierarchical XTransformer approach for Extreme Multi-label Classification (XMC). This method organizes labels hierarchically, reducing the computational complexity and improving prediction accuracy by leveraging the inherent structure in the data.

## 3.3 Data Augmentation

To further enhance the robustness of the model, we apply two types of data augmentation:

1. **Edge Perturbation**: We add or remove edges with probabilities 0.1 and 0.05, respectively, introducing variability in the graph structure.

2. **Gaussian Noise**: We add random Gaussian noise to node features, ensuring the model remains resilient to feature variations.

## 3.4 Contrastive Learning

Within the xTransformer module, we update the training process from adjacency matrix label prediction to contrastive learning. This involves using positive and negative samples to learn better node representations. Positive samples are nodes with similar features or connections, while negative samples are dissimilar nodes.

# 4 Experiments

## 4.1 Environment Setup

We used the NYU Greene HPC setup to download and preprocess data. The datasets include ogbn-arxiv, ogbn-products, and ogbn-papers100M. The setup involves downloading the necessary packages and pre-processed data files into a designated directory structure.

## 4.2 Baseline on Single Dataset

We ran baseline experiments on single datasets using different GNN algorithms:

```
dataset=***
gnn_algo=***
bash ./run_ogb_baselines.sh ${dataset} ${gnn_algo}
```

For ogbn-arxiv: mlp/graph-sage; for ogbn-products: mlp/graph-saint; for ogbn-papers100M: mlp/sgc.

## 4.3 Results

Table shows the test accuracy of baseline models across three different datasets: ogbn-arxiv, ogbn-products, and ogbn-papers100M. The results indicate that the MLP model generally performs the poorest compared to the other models, with test accuracies of 72.05%, 80.50%, and 60.98% on the respective datasets. In contrast, GraphSAGE, GraphSAINT, and SGC models show improved performance with test accuracies of 74.67%, 83.74%, and 65.32% respectively.

| Dataset | Model | Test Accuracy (%) |
|---|---|---|
| ogbn-arxiv | MLP / GraphSAGE | 72.05 / 74.67 |
| ogbn-products | MLP / GraphSAINT | 80.50 / 83.74 |
| ogbn-papers100M | MLP / SGC | 60.98 / 65.32 |

Table 1: Baseline models test accuracy comparison for different datasets.

On average, the test accuracy of **MLP** across the three datasets is approximately **71.84%**, whereas the average accuracy for the **graph-based models** (GraphSAGE, GraphSAINT, and SGC) is around **74.58%**. This demonstrates a notable improvement when using graph-specific models over the MLP baseline.

If we concatenate these three datasets into a single large dataset, the text feature extraction transformer BERT can be trained more effectively, leveraging the diversity and volume of data to improve overall performance. BERT's ability to extract rich text features is expected to enhance the downstream task accuracy when integrated into the GIANT-XRT framework.

However, due to time constraints, the contrastive learning part of the project was not completed. Implementing contrastive learning requires significant modifications to the sampling structure in the xTransformer, which proved to be time-consuming and complex. Future work will focus on completing this aspect to further enhance the model's performance through improved training processes.

## 5 Conclusion

This project addressed key limitations in existing graph-based machine learning models by introducing multi-dataset concatenation and contrastive learning to the GIANT-XRT framework. Our enhancements led to improved efficiency and effectiveness, demonstrating the potential for broader application and scalability in real-world tasks. Future work will explore further optimization and integration with other advanced GNN techniques.

## References

[1] Eli Chien, Wei-Cheng Chang, Cho-Jui Hsieh, Hsiang-Fu Yu, Jiong Zhang, Olgica Milenkovic, and Inderjit S. Dhillon. Node feature extraction by

self-supervised multi-scale neighborhood prediction. *CoRR*, abs/2111.00064, 2021.