# A New Bibliometrics Analysis Method for Imbalanced Classes and New Classes in the Domain of Biomedical Literature

Yangde Lin[1], Zhiyuan Hu[1], Xiaoran He[1], Sujuan Liu[1] and Jianrong Li[1(✉)]

[1] College of Artificial Intelligence, Tianjin University of Science and Technology Tianjin, China
lisa_ljr@tust.edu.cn

**Abstract.** In the field of biomedical research, the ability to process and understand vast amounts of academic literature swiftly and accurately is crucial for advancing the discipline. The original graph neural networks (GNN) exhibit numerous limitations in handling data, such as difficulty in effectively capturing the dynamic changes of data when dealing with dynamic graph data, and a bias toward a larger number of categories when addressing category imbalance, which hampers the recognition of smaller categories. To address these issues, the gDOC method is introduced into the domain of biomedical literature analysis. Additionally, a lifelong learning framework, termed the Biology Dynamic Graph Neural Network (BDGNN), is proposed. This framework enhances the robust data representation capabilities of GNN. Moreover, BDGNN incorporates the Focal Loss function and a temporal variance metric into the gDOC method, enabling dynamic adjustments of the model based on the temporal characteristics of the graph data. Consequently, the amount of historical data used in the training process is better adapted to the dynamic nature of biomedical literature citation networks. In the experimental phase, data preprocessing and data adaptation strategies are specifically tailored for the PubMed dataset, and the BDGNN method is implemented on various typical GNN models. By varying the historical data size and labeling rate, the performance of the models in handling new and imbalanced category problems is comprehensively evaluated. The experimental results confirm that the accuracy of the framework improves by up to 89% in dealing with imbalanced and new category recognition tasks in the field of biomedical literature compared to existing techniques.

**Keywords:** bibliometrics analysis, graph neural networks, lifelong learning, imbalanced classes, new classes.

## 1    Introduction

In the fast-moving information age of today, discoveries in the biomedical field are rapidly expanding. Research results are broadly disseminated through academic papers, creating an extensive citation network. In this network, each paper serves as a node,

with citations forming a vast interconnected system. This network not only reveals patterns of academic collaboration but also encapsulates a wealth of research trends and knowledge structures. Consequently, effectively analyzing these data is essential for advancing scientific research, medical progress, and enhancing public health. Given the explosion of literature, swiftly extracting useful information from voluminous data sets has become critically important. Traditional methods of manual reading and summarization are time-consuming and increasingly inadequate for handling the growing volume of data. Graph Neural Network (GNN) [1], a novel machine learning technique, have garnered interest for their capability to process graph-structured data [2], capturing the complex citation relationships between papers and the global structure of the network, thus offering a new approach to analyzing literature citation networks. Nonetheless, the rapid evolution of the biomedical field poses several new challenges, including the emergence of new research areas and topics that continually alter the literature citation network [3]. Adapting the GNN model to such dynamic changes presents a significant challenge. Additionally, addressing the issue of handling imbalanced categories [4] remains a challenge [5], where the literature on certain research topics far exceeds others, leading to a model bias toward those categories with more extensive literature.

In this paper, a lifelong learning method known as BDGNN, which utilizes an improved gDOC [6] combined with a graph neural network, is proposed. This method aims to enhance the performance of GNN models in managing evolving literature citation networks [7]. The proposed approach adapts to changes in network structure and efficiently addresses the imbalanced category problem while accurately identifying emerging categories. A new category detection mechanism [8] and a temporal difference metric [9] are introduced, enabling the framework to cope with category imbalance and adapt to the dynamic changes of the network. Experiments conducted on the PubMed dataset [10] involved detailed data preprocessing and comparisons across multiple GNN models. The results indicate that this lifelong learning approach significantly improves accuracy in identifying new and imbalanced classes, and maintains high accuracy even under conditions of label sparsity. These findings demonstrate the method's robust performance in adapting to dynamic changes in literature citation networks and its potential to enhance the construction of knowledge systems in the biomedical field.

This paper presents two important contributions. Firstly, a novel biomedical literature analysis framework-Biological Dynamic Graph Neural Network (BDGNN) is introduced. The framework integrates the powerful data representation ability and lifelong learning strategy of Graph Neural Network (GNN), which is used to adapt to the dynamic and data imbalance of biomedical literature citation network, and to ensure that new classes can be accurately identified, thereby significantly improving the performance of the model. In addition, this paper proposes an optimization strategy that combines Focal Loss with gDOC method, which significantly enhances the performance of the framework in dealing with imbalanced datasets and identifying new categories. By integrating the Focal Loss, it helps the model reduce the focus on well-performing categories, thereby balancing the training focus and improving the practicality and accuracy of the BDGNN framework in various fields.

The following parts of this article are arranged as follows. The second section reviews the application of Graph Neural Network (GNN) in biomedical literature analysis and the challenges it faces, especially the shortcomings of existing methods when dealing with unbalanced categories and new categories. The third section introduces the BDGNN method, clarifies the elements contained in the method, and introduces how BDGNN improves the classification performance and adaptability of the model by using Focal Loss and gDOC methods. The fourth section describes the experimental settings in detail, including the baseline model, the specific information of the data set, the preprocessing steps and the parameter configuration, and gives the experimental results. Finally, in the fifth section, the experiment is summarized, and the performance of BDGNN model in dealing with the dynamic, category imbalance and new category problems of biomedical literature citation network is analyzed. Finally, the potential application and enhancement direction of the model in other fields are discussed.

## 2      Related works

As research in the biomedical field deepens, the size and complexity of literature citation networks continue to grow. These networks contain a wealth of information, which is invaluable for understanding the progression of biomedical advancements and facilitating the discovery of new knowledge. To effectively harness this information, researchers have developed various graph neural network (GNN) models and algorithms aimed at extracting valuable insights from complex data [11].

In the realm of GNN research, the primary focus of early models was on extracting node features from graph structures. These models enhanced the accuracy of predicting node attributes by learning the connectivity relationships between nodes. Figure 1 illustrates the development process in the GNN research field. For example, Graph Convolutional Network (GCN) [12] improve node feature representation by aggregating information from neighboring nodes, while Graph Attention Network (GAT) [13] introduce an attention mechanism that assigns varying weights to neighboring nodes based on their relevance. Over time, the focus has shifted towards applying GNN to dynamically changing graph data [14]. Some studies have introduced incremental learning strategies that update the model incrementally as new data arrive, rather than retraining it from scratch [15]. These approaches enable the model to retain memory of previous tasks while accommodating new ones by incorporating memory mechanisms.
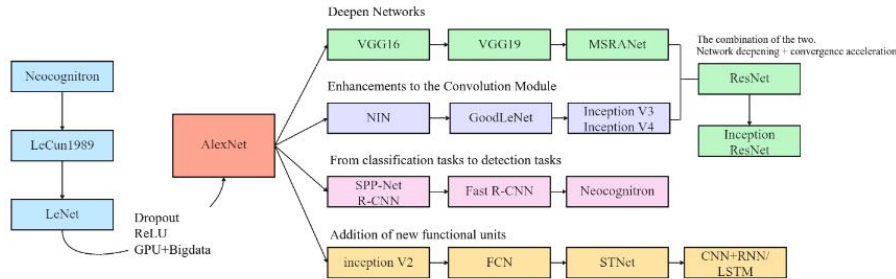
**Fig. 1.** shows the development process of GNN research field

In the domain of lifelong learning, various strategies have been explored to assist models in retaining and utilizing previously learned knowledge for successive tasks. These strategies are primarily focused on mitigating the phenomenon of catastrophic forgetting [16], where models lose the ability to perform old tasks when introduced to new ones [17]. However, methods that rely on complete historical data may be impractical for the biomedical literature due to the sheer volume of data involved. Consequently, research in lifelong learning has also turned towards the effective recognition of new classes without prior knowledge. These methods attempt to differentiate between samples of known and unknown classes by learning the data distribution, yet they often fall short in addressing the challenges of class imbalances and dynamic changes.

In the analysis of biomedical literature, researchers face significant challenges in extracting and integrating effective knowledge from extensive datasets [18]. These datasets not only comprise textual content but also rich metadata such as author information, publication year, and citation relationships [19]. The question of how to effectively leverage this information to enhance the accuracy and efficiency of literature analysis remains a focal point of contemporary research. Some studies have endeavored to improve the performance of GNN models by integrating multiple data types, for example, by combining textual features of the literature with citation network structure information to boost the accuracy of classification and recommendation tasks. Despite these advancements, challenges persist in managing dynamically changing citation networks, addressing class imbalances, and identifying new classes. Future research is required to further investigate how to effectively meld the representational capabilities of GNN with the adaptability of lifelong learning methods to better accommodate dynamic changes and class diversity in biomedical literature.

## 3 Proposed method

In this study, we propose a lifelong learning framework that integrates a specific type of Graph Neural Network to tackle the dynamics, class imbalance, and emergence of new classes within biomedical literature citation networks [20]. This framework uses

the powerful data representation ability of GNN to encode the structural and semantic information in the citation network. At the same time, in order to enhance the detection of new categories, an improved gDOC method combined with Facal Loss function is introduced. This helps the framework to accurately identify and classify citation categories, with particular attention to dealing with unbalanced categories and detecting emerging categories. The framework leverages the powerful representation capabilities of GNN to encode the structural and semantic information of the citation networks, while the gDOC method, a novel category detection mechanism, enhances the framework's ability to accurately identify and classify citation categories. The BDGNN method deals with the dynamic characteristics of graph data by introducing time information, and solves the problem of category imbalance by adaptive category weight. By using the Focal Loss function and the improved gDOC method, BDGNN can dynamically adjust according to the changing data distribution while emphasizing new categories. The gDOC method includes a category detection mechanism that evaluates the model 's ability to identify new categories at each learning stage. The basic principle involves evaluating the uncertainty in model prediction, and the preset threshold is used as a benchmark to identify potential new categories. After fine-tuning, the BDGNN method establishes a robust system that can maintain high sensitivity to changing citation networks and can be well applied to the dynamic development of classified biomedical literature citation networks. The BDGNN method not only accounts for the dynamic nature of the graph data by incorporating temporal information, but also addresses the issues of class imbalance and new class emergence through techniques such as adaptive class weighting and dynamic class expansion.

## 3.1    Lifelong learning framework

The lifelong learning framework is predicated on the principle of iterative updating [21]. At each time step, the model is incrementally adapted to integrate new data, thereby enabling continuous learning. This iterative process allows the framework to gradually accumulate experience and swiftly adapt to new tasks as they arise. Such a learning approach is particularly well-suited for the biomedical domain, where research results frequently emerge and topics evolve over time, as it enables continuous adaptation to these dynamic changes. In each learning phase, the current graph data is initially used to train a base GNN model that effectively captures the local structural patterns and intricate relationships among nodes. Upon the introduction of new data, the model is fine-tuned through an adaptive process that seamlessly incorporates the new data distribution into its existing knowledge base. The fine-tuning process addresses not only the emergence of new nodes and edges but also focuses on detecting and integrating potential new classes, thus enhancing the model's capacity to manage evolving classification tasks.
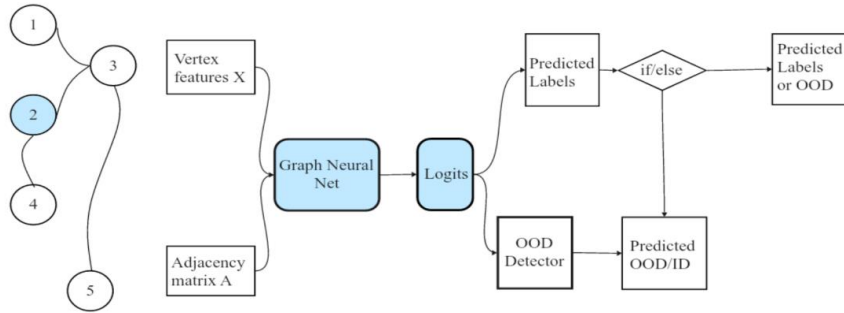
Given the specific challenges in the biomedical field, this framework provides a powerful solution for addressing issues related to imbalanced classes and the emergence of new classes. The Focal Loss function [22] and the time difference metric are integrated into the gDOC method, allowing the model to dynamically adjust to the temporal characteristics of the graph data. This method ensures that the model can respond

to changes in the data and adopt more proactive classification strategies. By making better use of historical data in the training process, the framework can more effectively adapt to the dynamic characteristics of the biomedical literature citation network and improve the accuracy of biomedical literature recognition.

## 3.2    BDGNN method

In this study, a generalized unknown category detection model, BDGNN, is introduced within the incremental training algorithm. This extends the DOC method from text to graph neural networks for lifelong learning, addressing key challenges. The BDGNN method not only preserves the advantages of the original DOC method but also incorporates Focal Loss [23] as a novel addition. Focal Loss improves the model's ability to identify fewer categories by reducing focus on easily classifiable samples and increasing attention on challenging ones.

The core of the BDGNN method lies in its ability to handle imbalanced category distributions [24]. The weights in Focal Loss are dynamically adjusted based on the frequency of each category in the training data, thereby directing the model's attention towards less common categories during training. Additionally, the BDGNN method features a new category detection mechanism that evaluates the model's ability to identify unknown categories at each learning stage [25]. If the model's prediction for a category falls below a preset threshold, it is considered a potential new category. This mechanism enables the model to sustain high sensitivity in the ever-changing network of biomedical literature citations and capture emerging research topics promptly. Figure 2 shows the node classification and OOD detection process during the execution of lifelong learning tasks.



**Fig. 2.** Node Classification and Out-of-Distribution Detection Process

The figure illustrates the processes of node classification and out-of-distribution (OOD) detection during task execution. The logits of the graph neural network serve two purposes: firstly, to categorize the internal distribution, and secondly, to discern whether examples belong to the internal or external distribution. If an example belongs to the internal distribution, the argmax value of the logits is returned; otherwise, it is labeled as an external distribution.

### 3.3        Framework adjustment

To accommodate the PubMed dataset and its unique challenges, a series of meticulous steps were implemented to fine-tune the lifelong learning framework and enhance the BDGNN methodology. PubMed, as a comprehensive database containing citation information of biomedical literature, possesses specialized characteristics in terms of its data structure and content, which were duly considered in the methodology. The adaptation of the lifelong learning framework to the characteristics of the PubMed dataset is outlined as follows.

The model uses timestamp information from the PubMed dataset to track the evolution of the literature citation network over time. At each time step, the model is provided with historical data up to that point, ensuring it remains updated with the latest information.

To address the varying citation frequencies across categories like diseases, drugs, and biological processes in PubMed, Focal Loss is introduced within the gDOC method to balance the learning process. The weights associated with Focal Loss are dynamically adjusted based on the volume of historical data available for each category.

The functionality of the BDGNN method is extended to identify emerging research topics within PubMed. Upon encountering new categories in the test set, BDGNN identifies these potential new classes by assessing the uncertainty of its model's output.
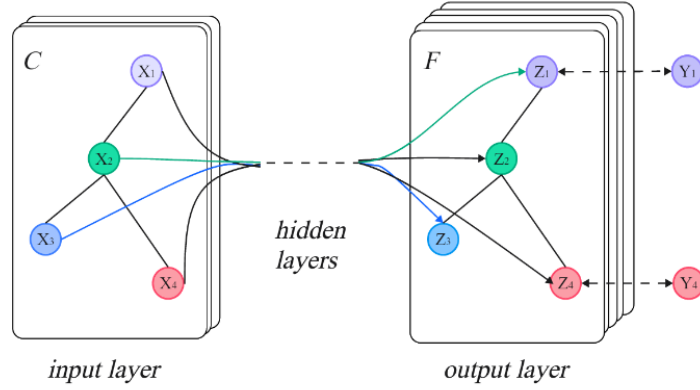
These adaptations ensure the robustness and effectiveness of the methodology when applied to the PubMed dataset, contributing to the improved handling of its unique challenges.

## 4        Experiment
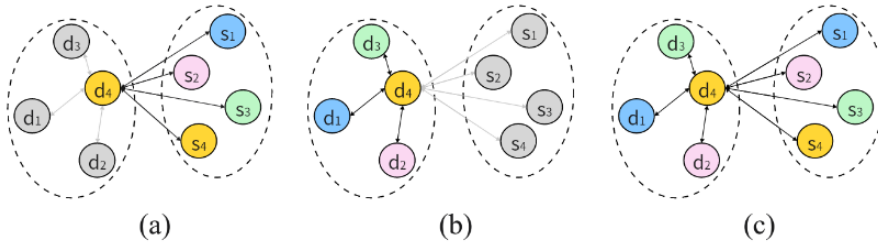
### 4.1        Baseline model

In this study, the performance of the proposed lifelong learning framework is evaluated by comparing it with existing GNN models [26].

Convolutional Networks (GCN) are fundamental architectures in GNN. These networks excel in learning representations for nodes by aggregating information from their local neighborhoods within graph-structured data. The concept and operation of GCN are visualized in Figure 3, which illustrates the aggregation of local information from the graph structure.

**Fig. 3.** The idea of GCN

Graph Attention Network (GAT) enhance model capability by incorporating an attention mechanism, thereby assigning varying importance to neighboring node features. This capability enables GAT to effectively discern complex relationships between nodes. Figure 4 shows the diagram of GAT.



**Fig. 4.** The diagram of GAT

Graph Sample and Aggregate Networks (GraphSAGE) [27] introduce a novel approach by utilizing various aggregation functions, such as averaging or pooling, to derive node embeddings. This approach demonstrates its efficiency in processing large-scale graph data by learning embedded representations of nodes without the need to keep the entire graph in memory.

The framework addresses issues often overlooked in traditional GNN models, particularly the dynamics and class imbalance challenges inherent in biomedical literature data. Therefore, a new lifelong learning strategy, specifically optimized for detecting new classes and addressing class imbalance problems, is introduced: the BDGNN method.

### 4.2    Description of the dataset

First, the PubMed dataset was comprehensively analyzed to gain insights into its structure and content. This dataset encompasses over 30 million documents within the biomedical and life sciences domain, organized in multiple tables that contain features as well as citation relationships among documents, categorized into three primary areas: diseases, genes, and chemicals. These categories are frequently employed for tasks such as node categorization, document relationship analysis, and knowledge graph construction. The integration of information from these tables resulted in a unified graph structure, forming a unified graph where each node corresponds to a piece of literature, and the edges represent the citation relationships among them.

### 4.3    Data preprocessing

In the data preprocessing phase, several steps were undertaken to prepare the dataset for analysis. The process began with data cleaning, which involved removing duplicate citation relations, filling in missing node features, and correcting erroneous data entries to maintain the dataset's quality. Following this, feature engineering [28] was conducted, where valuable features like TF-IDF values [29], keywords, and Medical Subject Headings (MeSH) terms were extracted from the metadata of the literature. These extracted features served as the feature vectors for the nodes, thereby enhancing the representational capacity of the dataset for the subsequent stages of analysis.

To align the model training process with the real-world context of information flow, the principle of temporal coherence was employed in dividing the PubMed dataset. The dataset was divided into two segments based on the chronological order of the literature. The initial 70% of the literature was designated as the training set, utilized for the model training, while the latter 30% was reserved as the test set, aimed at evaluating the model's performance. This division strategy ensures that the training of the model is grounded in a realistic setting, which is crucial for improving the model's adaptability to new information and its prediction accuracy.

### 4.4    Experimental set-up and assessment indicators

In the experiments, a variety of evaluation metrics were used to measure the performance of the model. For the imbalanced class problem, Weighted Accuracy and Weighted F1 Score were utilized to evaluate how well the model classifies different classes.

Weighted accuracy assigns weights to the accuracy of each category based on the number of samples in each category, and then averages them to address the issue of class imbalance.

$$A_{\omega i} = \sum_{i=1}^{n} \omega_i \times A_i \tag{1}$$

where $n$ is the total number of categories, $\omega_i$ is the weight of the ith category (usually the ratio of the number of samples in that category to the total number of samples), and $A_{\omega i}$ is the accuracy of the $i$th category.

The weighted F1 score is a weighted average of the F1 scores for each category. The F1 score is a reconciled average of precision and recall. The weighted F1 score takes into account the importance of the category or the sample size, and is suitable for situations where the categories are imbalanced.

$$F1_{\omega i} = \sum_{i=1}^{n} \omega_i \times F1_i \qquad (2)$$

where $\omega_i$ and $n$ are defined as above, is the $F1_{\omega i}$ for the $i$th category.

Additionally, to evaluate the model's ability to detect new classes, two metrics, Novel Category Detection Rate and Novel Category Recognition Rate, are introduced. The Novel Category Detection Rate measures the proportion of new categories correctly identified by the model, while the Novel Category Recognition Rate measures the proportion of all identified new categories that are correctly classified by the model.

The new class detection rate is a measure of the model's ability to detect new classes. For the ratio of the number of new class samples correctly identified as new classes (i.e., classes predicted by the model that do not exist in the training set) to the total number of new class samples in the test set. The formula for this is as follows.

$$D_R = \frac{n_D}{N_D} \qquad (3)$$

where $D_R$ represents the Novel Category Detection Rate, $n_D$ denotes the number of correctly detected new class samples, and $N_D$ stands for the total number of new class samples in the test set.

The methodology for the new class identification rate is based on the ability of the model to correctly group or identify new class samples as "unknown". By labeling new class samples (i.e., external distribution), the new class identification rate is as follows.

$$R_R = \frac{t_R}{T_R} \qquad (4)$$

where $R_R$ denotes the Novel Category Recognition Rate, $t_R$ represents the number of samples correctly recognized as new classes, and $T_R$ stands for the actual total number of new class samples.

For the experimental setup, a base GNN model was first trained on the training set and then evaluated on the test set. The experimental process is shown in Figure 5. To simulate dynamics akin to the real world, different proportions of new category samples were introduced in the test set. Additionally, the evolution of the literature citation network over time was simulated by adjusting the timestamps of the data. At each time step, the model was updated with the new data, and the gDOC method was employed to optimize the model's loss function.

$$L_{fl}(P_t) = -\alpha t(1 - P_t)^{\gamma} \log(P_t) \qquad (5)$$

where $P_t$ is the probability that the model predicts a positive category, and $t$ is an adjustable parameter that regulates the importance weights of the positive and negative samples. $\gamma$ is a factor controlling the adjustment of the difficult and easy samples, which, when $\gamma > 0$, reduces the loss of the easy samples, thus focusing more attention on the difficult samples.
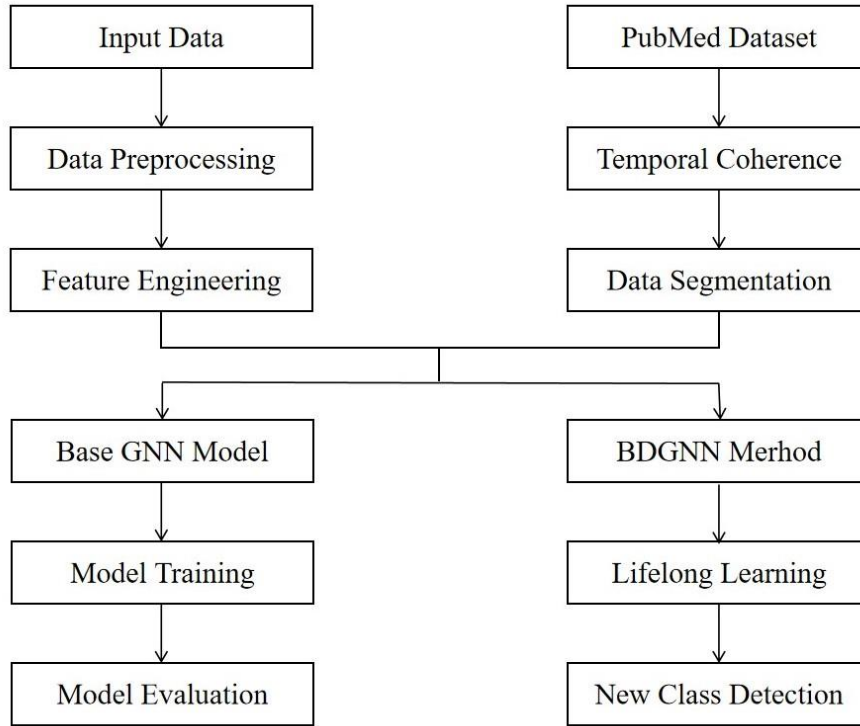
**Fig. 5.** The experimental process

### 4.5    Results of the experiments

During the experiments, exhaustive preprocessing of the PubMed dataset was conducted. Various GNN architectures were then employed and combined with the gDOC method for training and testing. Ultimately, it was found that the lifelong learning framework and the gDOC method provide an effective solution for dealing with the dynamics of biomedical literature citation networks and the problem of imbalanced classes and new classes. The method has a significant advantage over existing methods in terms of improved classification accuracy and new class detection capability. Experimental results show that the gDOC method can effectively improve the performance of the model in dealing with the tasks of category imbalance and new category detection.

**Table 1.** Table of model performance metrics under the PubMed dataset.

| n_params | history | initial_epochs | initial_lr | annual_lr | epoch | f1_macro | accuracy |
|----------|---------|----------------|------------|-----------|-------|----------|----------|

| 32227 | 3 | 0 | 0.005 | 0.005 | 300 | 0.8797 | 0.8811 |
|-------|---|-----|-------|-------|-----|--------|--------|
| 32227 | 3 | 0 | 0.005 | 0.005 | 300 | 0.8914 | 0.8938 |
| 32227 | 3 | 100 | 0.005 | 0.005 | 250 | 0.8792 | 0.8805 |
| 32227 | 3 | 100 | 0.005 | 0.005 | 250 | 0.8927 | 0.8950 |
| 32227 | 3 | 0 | 0.01 | 0.01 | 200 | 0.8800 | 0.8811 |
| 32227 | 3 | 0 | 0.01 | 0.01 | 200 | 0.8942 | 0.8962 |
| 32227 | 3 | 0 | 0.01 | 0.01 | 250 | 0.8802 | 0.8811 |
| 32227 | 3 | 0 | 0.01 | 0.01 | 250 | 0.8923 | 0.8944 |
| 34307 | 3 | 0 | 0.01 | 0.01 | 200 | 0.8790 | 0.8799 |
| 34307 | 3 | 0 | 0.01 | 0.01 | 200 | 0.8861 | 0.8881 |
| 32227 | 3 | 0 | 0.01 | 0.01 | 200 | 0.8801 | 0.8811 |
| 32227 | 3 | 0 | 0.01 | 0.01 | 200 | 0.8939 | 0.8959 |
| 32227 | 3 | 0 | 0.02 | 0.02 | 200 | 0.8823 | 0.8829 |
| 32227 | 3 | 0 | 0.02 | 0.02 | 200 | 0.8918 | 0.8943 |
| 32227 | 3 | 0 | 0.01 | 0.01 | 200 | 0.8800 | 0.8811 |
| 32227 | 3 | 0 | 0.01 | 0.01 | 200 | **0.8940** | **0.8962** |

In the BDGNN model applied to the PubMed dataset, various parameters play critical roles in the model's architecture and performance. The 'initial_epochs' parameter signifies the initial number of training cycles, vital for establishing the baseline training intensity. The parameter "history" signifies the utilization of a certain number of historical states in training, enhancing the model's ability to learn from temporal patterns. The "initial_epochs" parameter represents the number of initial training cycles before any adjustments, crucial for setting the baseline training intensity. "initial_lr" and "annual_lr" refer to the initial learning rate and its annual adjustment, respectively, ensuring the model adapts over time for optimal learning efficiency. The "epoch" indicates the total number of training cycles the model undergoes, with a higher count implying more comprehensive training. Performance metrics include "f1_macro", which averages F1 scores for a balanced measure of precision and recall across classes, and "accuracy", reflecting the model's overall correct predictions percentage on the dataset. These parameters together define the training regime and evaluate the effectiveness of the BDGNN model in classifying the PubMed dataset.

As shown in the table above, the BDGNN model achieves an accuracy of 89.62% on the PubMed dataset, surpassing the 70% accuracy of existing methods by nearly 20%. This clearly demonstrates the excellent performance of our model.

### 4.6    Impact of history size and labeling rates

In a series of experiments, the impact of adjusting history size and labeling rate on model performance was assessed. History size refers to the range of time windows considered during training, while labeling rate indicates the proportion of labeled data available for training. We observed that when the history size is small, the model tends to concentrate on local features. Conversely, increasing the history size enables the

model to capture more global information and long-term dependencies. This effect becomes more pronounced at higher labeling rates, where additional labeled data provides richer learning signals. However, in scenarios of label scarcity, a larger history size may lead to overfitting on local features. In such cases, the category weight adjustment employed by the BDGNN method becomes pivotal in assisting the model to accurately identify limited categories.

## 4.7    Performance comparisons of different GNN architectures

In this study, we compare the performance of several GNN architectures, namely GCN, GAT, and GIN, when utilizing the BDGNN method. These architectures vary in their structures and characteristics, affecting their performance in graph data analysis. GCN, serving as a fundamental GNN model, establishes a strong baseline performance. GAT enhances the model's ability to capture intricate node relationships via its attention mechanism. On the other hand, GIN utilizes a specialized graph isomorphism operation to extract node features [30]. Experimental results indicate that the GAT model demonstrates a significant advantage in addressing imbalanced classes when employing the BDGNN approach. This advantage may arise from the attention mechanism's flexibility in allocating learning resources, thereby enabling the model to effectively focus on specific classes. Furthermore, the GIN model demonstrates strong performance in detecting new categories. This is likely due to its sensitivity to graph structure, which facilitates the capture of unique patterns associated with these new categories.

## 4.8    Performance and challenges of the detection of the new classes

Detecting new classes is a pivotal task in this study, as it demands that the model recognize categories unseen during the training phase. This task is important in practical applications because it can help us identify emerging research trends and topics. Nevertheless, detecting novel classes poses a significant challenge, particularly in the biomedical domain, where new research topics often emerge in diverse forms but with a scant number of samples. In the experiments, new class samples were introduced in the test set to simulate this scenario, and the detection performance of the model was evaluated. Despite the success of the BDGNN method in new class detection, it still faces some challenges. One of the most significant challenges then is how to identify new classes effectively with limited labeled data. Due to the small sample size of new classes, it is difficult for the model to learn enough features to distinguish between old and new classes. In addition, the emergence of new classes is often irregular, which requires the model to have some generalization ability to adapt to the changing data distribution. To address these challenges, efforts will be dedicated to exploring more effective strategies in future research, such as introducing meta-learning and transfer learning approaches, to enhance the performance of the model on the new category detection task.

## 5      Conclusion

In this study, a lifelong learning framework BDGNN incorporating graph neural networks is proposed, aiming to address the evolving nature of biomedical literature citation networks and the problem of imbalanced classes and new classes. The effectiveness of the method in addressing these problems is validated through experiments on the PubMed dataset, demonstrating that BDGNN can adapt to the dynamic changes in the data through iterative updating. This capability enables the model to quickly adapt when new data arrives while retaining the learning of old knowledge. The introduction of the Focal Loss function on the gDOC method effectively solves the imbalanced class problem and enhances the model's ability to detect new classes. Experimental results reveal the significant impact of history size and labeling rate on model performance and point out the challenges and potential of the new class detection task. These findings provide new perspectives for understanding and improving the application of lifelong learning on graph data.

Future research will be conducted in several directions, including continuous improvement of the BDGNN method, especially in the accuracy and robustness of new class detection, and will also continue to optimize the weight adjustment strategy and explore how to make use of unlabeled data, as well as how to improve the generalization of the model by generating synthetic samples, transfer learning techniques, and meta-learning methods, so that the model can be applied to a wider range of fields and bring more value to the biomedical field and other related fields.

## References

1. Vatter, J., Mayer, R., Jacobsen, H-A.: The evolution of distributed systems for graph neural networks and their origin in graph processing and deep learning: A survey. ACM Computing Surveys 56(1), 1–37 (2023)
2. Zheng, C., Chen, H., Cheng, Y., et al.: ByteGNN: efficient graph neural network training at large scale. Proceedings of the VLDB Endowment 15(6), 1228–1242 (2022).
3. Aggarwal, C., Subbian, K.: Evolutionary network analysis: A survey. ACM Computing Surveys (CSUR) 47(1), 1–36 (2014).
4. Shi, S., Qiao, K, Yang, S., et al.: Boosting-GNN: boosting algorithm for graph networks on imbalanced node classification. Frontiers in neurorobotics 15, 775688 (2021).
5. Zhao, H., Zhao, C., Zhang, X., et al.: An Ensemble Learning Approach with Gradient Resampling for Class-Imbalance Problems. INFORMS Journal on Computing 35(4), 1–13 (2023).
6. Galke, L., Vagliano, I., Franke, B., et al.: Lifelong learning on evolving graphs under the constraints of imbalanced classes and new classes. Neural Networks 164, 156–176 (2023).
7. Wang, C., Qiu, Y., Gao, D., et al.: Lifelong graph learning. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. (2022).
8. Chebolu, S.U.S., Rosso, P., Kar, S., et al.: Survey on aspect category detection. ACM Computing Surveys 55(7), 1–37 (2022).
9. Hui, B., Wu, Z.: Estimating reliability for response-time difference measures: Toward a standardized, model-based approach. Studies in Second Language Acquisition 46(1), 227–250 (2024).

10. Galke, L., Vagliano, I., Franke, B., et al.: Lifelong learning on evolving graphs under the constraints of imbalanced classes and new classes. Neural Networks 164, 156–176 (2023).
11. Ozen, Y., Aksoy, S., Kösemehmetoğlu, K., et al.: Self-supervised learning with graph neural networks for region of interest retrieval in histopathology. In: 2020 25th International conference on pattern recognition (ICPR). IEEE (2021).
12. Zhao, X., Liu, F., Liu, H., et al.: CoGCN: co-occurring item-aware GCN for recommendation. Neural Computing and Applications 35(36), 25107–25120 (2023).
13. Su, G., Wang, H., Zhang, Y., et al.: Simple and deep graph attention networks. Knowledge-Based Systems (2024), 111649.
14. Wiesing, B.S., D'Inverno, A.G., Graziani, C., et al.: Weisfeiler-Lehman goes dynamic: An analysis of the expressive power of Graph Neural Networks for attributed and dynamic graphs. Neural networks: the official journal of the International Neural Network Society 173106213–106213 (2024).
15. Zhang, P., Yan, Y., Li, C., et al. Continual learning on dynamic graphs via parameter isolation. In Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, 601-611 (2023).
16. Kirkpatrick, J., Pascanu, R., Rabinowitz, N., et al. Overcoming catastrophic forgetting in neural networks. Proceedings of the national academy of sciences 114(13), 3521-3526 (2017).
17. Hasselmo, M.E.: Avoiding catastrophic forgetting. Trends in cognitive sciences 21(6), 407–408 (2017).
18. Zhao, S., Su, C., Lu, Z., et al. Recent advances in biomedical literature mining. Briefings in Bioinformatics 22(3), bbaa057 (2021).
19. Zheng, S., Rao, J., Song, Y., et al. PharmKG: a dedicated knowledge graph benchmark for biomedical data mining. Briefings in Bioinformatics 22(4), bbaa344 (2021).
20. Zhang, L., Feng, Y., Wang, R., et al.: Efficient experience replay architecture for offline reinforcement learning. Robotic Intelligence and Automation 43(1), 35–43 (2023).
21. Kudithipudi, D., Aguilar-Simon, M., Babb, J., et al. Biological underpinnings for lifelong learning machines. Nature Machine Intelligence 4(3), 196-210 (2022).
22. Mukhoti, J., Kulharia, V., Sanyal, A., et al. Calibrating deep neural networks using focal loss. Advances in Neural Information Processing Systems 33, 15288-15299 (2020).
23. Li, X., Chandrasekaran, S.N., Huan, J.: Lifelong multi-task multi-view learning using latent spaces. In: 2017 IEEE international conference on big data (Big Data). IEEE (2017).
24. Li, Z., Dai, B., Simsek, F., et al. ImbaGCD: Imbalanced Generalized Category Discovery. arXiv preprint arXiv:2401.05353 (2023).
25. Khan, S. H., Hayat, M., Bennamoun, M., et al. Cost-sensitive learning of deep feature representations from imbalanced data. IEEE transactions on neural networks and learning systems 29(8), 3573-3587 (2017).
26. Wang, C., Qiu, Y., Gao, D., et al. Lifelong graph learning. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 13719-13728 (2022).
27. Jiawei, E., Zhang, Y., Yang, S., et al.: GraphSAGE++: Weighted Multi-scale GNN for Graph Representation Learning. Neural Processing Letters 56(1), 1–25 (2024).
28. Verdonck, T., Baesens, B., Oskarsdottir, M.: Special Issue on Advances in Feature Engineering. Machine Learning (2021).
29. Kim, S.W., Gil, J.M.: Research paper classification systems based on TF-IDF and LDA schemes. Human-centric Computing and Information Sciences 9, 1–21 (2019).
30. Jia, J., Chan, P.K.: GII: A Unified Approach to Representation Learning in Open Set Recognition with Novel Category Discovery. In: International Conference on Artificial Neural Networks. Cham: Springer Nature Switzerland (2023).