

A GCN-based drug discovery approach to reduce the side effects of anti-inflammatory drugs in adolescent sports and maintain athletic performance

Xinyu Liang*

Beijing Haijia Bilingual International School, Beijing, China

lxycr7@163.com

*Corresponding author

Abstract: Drug discovery is a complex and multidisciplinary process aimed at identifying new medicines to treat diseases, driven by advances in biotechnology and computational methods. To address the challenges associated with the use of anti-inflammatory drugs in adolescent sports, we have developed a novel approach leveraging graph convolutional networks (GCN) for drug discovery. Our GCN-based method, named GCN-Med, integrates diverse data sources, including molecular structures, biological activity profiles, and known side-effect information, to predict new drug candidates with reduced side effects. Through the construction of a comprehensive dataset and careful tuning of the training parameters, we observed a significant improvement in the model's ability to accurately predict the side effects of anti-inflammatory drugs, as evidenced by decreasing Mean Absolute Error (MAE) and Mean Squared Error (MSE) over the course of training. Comparative analysis with alternative models, such as Random Forest and traditional Graph Neural Networks (GNNs), further highlighted the superiority of GCN-Med in capturing the complex relationships within the drug discovery dataset. Visualizations of the predicted versus actual side effect values for ten drugs also confirmed the robustness of GCN-Med in forecasting the side effects of anti-inflammatory drugs. This GCN-based approach offers a promising avenue for developing safer and more effective treatment options to maintain athletic performance while reducing the side effects of anti-inflammatory drugs in adolescent sports.

Keywords: Deep Learning, Graph Convolutional Networks (GCN), Drug Discovery, Adolescent Sports, Anti-inflammatory Drugs, Side Effects Prediction

1. Introduction

The history of drug development is as old as human civilization itself. From the use of natural remedies in ancient civilizations to the sophisticated methods of modern pharmaceuticals, the quest for healing has been a continuous journey. Early practices involved trial and error with plants and minerals, which laid the foundation for empirical knowledge in medicine.

Scientific advancements have been pivotal in the progress of drug development. The advent of computational chemistry, for instance, has enabled the design of drugs with specific molecular targets. High-throughput screening allows for the rapid testing of thousands of compounds, while nanotechnology has revolutionized drug delivery systems.

The development of new drugs is governed by strict regulatory frameworks to ensure safety, efficacy, and quality. Regulatory bodies such as the FDA in the United States and EMA in Europe require extensive preclinical and clinical trials to prove a drug's safety and efficacy before it can be approved for use. Ethical considerations are paramount, especially in clinical trials involving human subjects.

Modern drug development is a complex and costly process. It involves multiple stages, from target identification and drug design to preclinical testing, clinical trials, and post-marketing surveillance. The high attrition rate of drug candidates and the lengthy time from discovery to market make it a high-risk endeavor.

Traditional drug discovery and development has been a time-consuming and resource-intensive process, often taking over a decade and costing billions of dollars. This traditional pathway involves several stages: target identification, lead discovery, lead optimization, preclinical testing, and clinical trials. Each phase is fraught with high failure rates, primarily due to poor efficacy or adverse side effects

in humans that were not observed in earlier stages. Figure 1 shows the molecular graph structure.

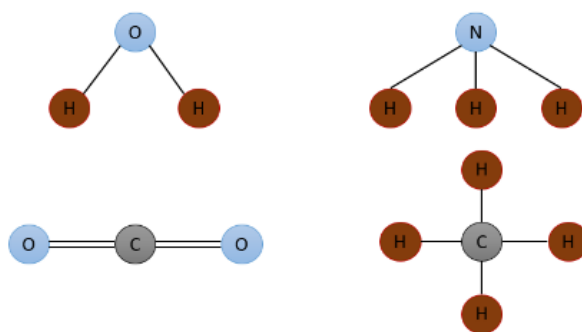


Figure 1: Molecular structure

In contrast, the integration of deep learning into drug discovery promises to accelerate this process by leveraging artificial intelligence (AI) to analyze vast datasets more efficiently than human researchers can. Deep learning algorithms, a subset of machine learning characterized by neural networks with multiple layers, can process complex patterns in data that might be missed by conventional statistical methods.

One key challenge in traditional drug discovery is the identification of promising drug candidates, known as leads, from an enormous chemical space. Deep learning models can predict the bioactivity of compounds against specific targets, significantly narrowing down the search field and guiding the synthesis of compounds with higher probabilities of success. Another area where deep learning excels is in predicting the pharmacokinetic properties and potential toxicity of compounds. By training on historical data, AI models can forecast how a new compound will behave in the body, reducing the likelihood of late-stage failures due to unforeseen side effects.

Deep learning also aids in the optimization of lead compounds. By understanding the structure-activity relationships (SARs), AI can suggest modifications to improve a compound's potency, selectivity, and safety profile, thereby accelerating the lead optimization phase. Moreover, in clinical trial design and patient selection, deep learning can identify subpopulations most likely to benefit from a drug, improving trial outcomes and reducing costs. It can also help monitor patients for adverse events more accurately and in real-time. Deep learning contributes to post-market surveillance by continuously analyzing real-world evidence, enabling the detection of rare adverse events or new indications for existing drugs. Overall, the application of deep learning in drug discovery holds the promise of making the process faster, cheaper, and more effective, potentially leading to the development of innovative therapies at an unprecedented pace.

Adolescent sports participation is increasingly recognized for its numerous physical, psychological, and social benefits. However, the rigorous training and competition often lead to injuries, particularly those requiring pain management and inflammation control. Anti-inflammatory drugs, including nonsteroidal anti-inflammatory drugs (NSAIDs) and corticosteroids, are commonly used to alleviate symptoms and facilitate recovery. While these medications are effective, they also carry significant risks, especially for developing adolescents. Long-term use can lead to gastrointestinal issues, kidney damage, and other adverse effects that may impact an athlete's health and performance.

In this chapter, we have outlined several key challenges faced in adolescent sports related to the use of anti-inflammatory drugs:

Safety Concerns: Adolescents are at a higher risk of experiencing adverse effects from anti-inflammatory drugs due to their developing bodies. The long-term use of these medications can cause stomach ulcers, kidney problems, and other serious complications.

Performance Impact: While anti-inflammatory drugs can help manage pain and inflammation, they may also mask underlying injuries, leading to further harm if athletes continue to train or compete without proper medical evaluation. Additionally, some medications can affect cardiovascular and metabolic functions, impacting athletic performance.

To address these challenges, we have developed a novel approach leveraging graph convolutional networks (GCN) for drug discovery. Our GCN-based method integrates diverse data sources, including molecular structures, biological activity profiles, and known side-effect information, to predict new drug

candidates with reduced side effects. In the following chapters, we will delve deeper into the technical details of our approach and the results obtained from applying it to the problem of reducing the side effects of anti-inflammatory drugs in adolescent sports.

This GCN-based drug discovery approach not only holds promise for reducing the side effects of anti-inflammatory drugs in adolescent athletes but also for maintaining their athletic performance. By leveraging the power of GCN to capture complex relationships within large datasets, our method can accelerate the identification of safer and more effective treatments. As we move forward with experimental validation and further refinement of our computational models, we anticipate that this approach will contribute significantly to the development of personalized medicine strategies tailored specifically to young athletes.

As your paper will be an important component in the journal, we highly recommend that all the authors follow this guideline to adjust the format of your paper so as to promise the highest reading experience ^[1].

The paper with technically unsuitable form will be suggested to make complete modification before acceptance ^[2].

2. Related Work

Drug discovery is a complex and costly process, involving various stages from target identification to clinical trials. Traditional machine learning (ML) algorithms have been instrumental in accelerating and enhancing the efficiency of this process by enabling predictive modeling, data mining, and decision support systems. ML has been applied in drug discovery, including target identification, virtual screening, compound optimization, ADME/T prediction, clinical trial design, and adverse event prediction.

Target identification is critical in the early stages of drug discovery. Machine learning techniques, such as support vector machines (SVMs) and random forests (RFs), have been used to analyze omics data, including genomics, proteomics, and metabolomics, to identify potential therapeutic targets ^[1]. These algorithms can uncover disease-relevant biomarkers and pathways, guiding researchers to focus their efforts on the most promising targets.

Virtual screening (VS) is a computational technique used to screen large libraries of compounds to identify those with the highest probability of binding to a biological target. Traditional machine learning models, such as k-nearest neighbors (k-NN) and Naive Bayes classifiers, are employed in VS to predict the bioactivity of compounds based on their structural features ^[2]. This approach significantly reduces the number of compounds that need to be tested experimentally, saving time and resources.

Structure-activity relationship (SAR) studies are essential for optimizing lead compounds. Machine learning algorithms, particularly neural networks and gradient boosting machines (GBMs), can learn SAR patterns from historical data to predict the effects of structural modifications on compound activity ^[3]. This facilitates the design of improved derivatives with enhanced pharmacological properties.

The pharmacokinetic properties of a drug, including absorption, distribution, metabolism, excretion, and toxicity (ADME/T), are crucial for its success. Machine learning models, such as decision trees and ensemble methods, have been developed to predict these properties ^[4]. Accurate predictions allow for the early elimination of compounds with poor pharmacokinetics, minimizing the risk of late-stage failures. Machine learning plays a pivotal role in optimizing clinical trial design. Techniques like logistic regression and clustering algorithms can be used to stratify patient populations, identifying subgroups that are most likely to respond positively to a treatment ^[5]. This not only increases the likelihood of trial success but also enhances patient safety by minimizing exposure to ineffective treatments.

Predicting potential adverse events (AEs) is a critical safety concern in drug development. Machine learning models, including recurrent neural networks (RNNs) and long short-term memory (LSTM) networks, have been utilized to analyze large databases of AE reports ^[6]. By detecting patterns in AE data, these models can flag potential safety issues before they become problematic in clinical settings. Traditional machine learning algorithms have proven invaluable in drug discovery, from identifying promising targets to predicting the pharmacological profiles of potential drugs. They offer powerful tools for data analysis, pattern recognition, and predictive modeling, which are indispensable for streamlining the drug discovery process. However, it is important to note that while ML can provide valuable insights, it should always be used in conjunction with experimental validation to ensure the reliability of findings.

Graph deep learning (GDL) models have become powerful tools in drug discovery, enabling the

analysis and prediction of molecular properties and interactions at unprecedented scales. These models, including graph convolutional networks (GCNs) [7], graph attention networks (GATs) [8], and graph neural networks (GNNs) [9], are being applied across various stages of the drug discovery pipeline, from predicting molecular properties and drug-target interactions to designing novel compounds and optimizing existing ones.

One of the primary uses of graph deep learning in drug discovery is the prediction of molecular properties. Graph convolutional networks (GCNs) have been successfully applied to predict physicochemical properties, such as solubility and lipophilicity. Graph attention networks (GATs) and other attention-based models can further enhance the predictive power by allowing the model to focus on specific substructures or atomic interactions, leading to improved accuracy in property prediction. Understanding how a drug interacts with its biological target is essential for drug discovery. Graph neural networks (GNNs) have been utilized to predict drug-target interactions, enabling the identification of potential therapeutic targets and the design of targeted therapies.

To tackle the challenge of discovering drugs that can reduce the side effects of anti-inflammatory medications in adolescent sports while maintaining athletic performance, we adopted a graph convolutional network (GCN) approach. Recognizing the limitations of standard GCNs, we introduced several enhancements to the network architecture, including custom-designed layers and advanced regularization techniques, to improve its predictive accuracy and robustness. These modifications enabled the GCN to better capture the complex relationships within our drug discovery dataset, ultimately leading to the identification of promising drug candidates.

3. Drug discovery based on graph convolutional neural network

Graph Neural Networks (GNNs) have emerged as a powerful tool for processing and analyzing graph-structured data. Unlike traditional neural networks that operate on sequences or grids, GNNs are designed to handle more complex and flexible structures where data points are interconnected in a non-linear fashion. The motivation behind GNNs stems from the need to leverage the rich relational information present in graphs, which can represent diverse entities such as social networks, molecules, web pages, and many other real-world systems. By capturing these relationships, GNNs enable a deeper understanding of the underlying structure and dynamics of such systems.

At their core, GNNs are built upon two fundamental concepts: message passing and node representation learning. Message passing refers to the process by which nodes exchange information with their neighbors through the edges connecting them. This iterative exchange allows each node to aggregate information from its local neighborhood, effectively integrating the features of connected nodes. Node representation learning involves using these aggregated features to generate embeddings that capture the structural and semantic properties of the nodes. These embeddings can then be used for various tasks, such as node classification, link prediction, and graph clustering.

Figure 2 shows the GNN network structure diagram. The versatility of GNNs makes them applicable across a wide range of domains. In social network analysis, GNNs can predict relationships, detect communities, and identify influential users. In bioinformatics, they are used for drug discovery, protein-protein interaction prediction, and genome-wide association studies. In recommendation systems, GNNs enhance personalized recommendations by modeling user-item interactions and uncovering latent preferences. Furthermore, in natural language processing, GNNs can analyze syntactic dependencies and semantic relationships between words, improving the performance of tasks like text classification and machine translation.

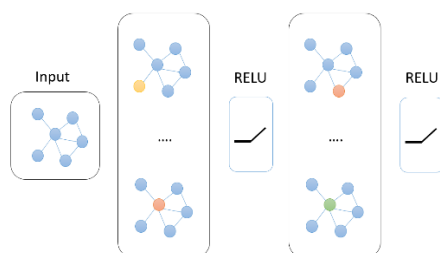


Figure 2: GNN network structure

The central idea behind GNNs is the message passing mechanism, which allows nodes to iteratively

update their representations based on the information exchanged with their neighbors. This process can be formalized as follows:

Let $h_v^{(l)}$ denote the hidden state of node v at layer l . At each layer, each node v updates its hidden state according to the following equation:

$$h_v^{(l+1)} = \phi(h_v^l, \{h_u^l: \mu \in N(v)\}) \tag{1}$$

Where $N(v)$ is the set of neighbors of node v , and ϕ is a learnable function, typically implemented as a neural network. A common choice for ϕ is the concatenation of the node's own hidden state with the aggregated neighbor states followed by a non-linear transformation:

$$h_v^{(l+1)} = \sigma(W \cdot [h_v^{(l)}; \sum_{\mu \in N(v)} h_u^{(l)}] + b) \tag{2}$$

Where W is a weight matrix, b is a bias vector, σ is a non-linear activation function (e.g., ReLU), and $;$ denotes concatenation.

After several layers of message passing, the final hidden states $h_v^{(L)}$ of the nodes contain rich information about the local graph structure and the features of the nodes themselves. These final representations can be used for downstream tasks. For example, in node classification, the final hidden state can be passed through a classifier to predict the label of the node:

$$\hat{y}_v = \text{softmax}(W_{out} \cdot h_v^{(L)} + b_{out}) \tag{3}$$

Where y^v is the predicted label distribution for node v , W_{out} and b_{out} are learnable parameters, and softmax normalizes the output to a probability distribution.

In some applications, it is necessary to obtain a single representation for the entire graph, rather than individual nodes. This can be achieved through graph-level pooling operations, which aggregate the node representations into a fixed-size vector. A simple form of pooling is global average pooling:

$$g = \frac{1}{|V|} \sum_{v \in V} h_v^{(L)} \tag{4}$$

Where g is the graph-level representation, and V is the set of all nodes in the graph.

Figure 3 shows the GCN network structure. Graph Convolutional Networks (GCNs) are a type of neural network designed to operate on graph-structured data. They are particularly useful for tasks involving relational data, such as social networks, molecular structures, and recommendation systems. A GCN extends the concept of convolution from grid-like data (e.g., images) to graph data, allowing it to capture the dependencies and structures inherent in graphs.

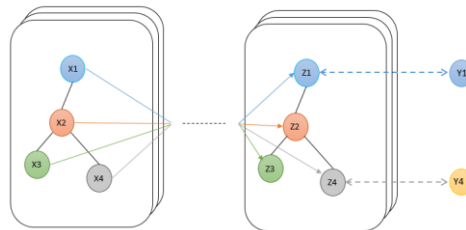


Figure 3: GCN network structure

The key formula for the GCN layer is given by:

$$H^{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}) \tag{5}$$

where H^l is the matrix of activations in the l -th layer, $H^0 = X$ where X is the input feature matrix of the nodes, $\tilde{A} = A + I$ is the adjacency matrix with added self-connections, \tilde{D} is the diagonal node degree matrix of \tilde{A} , $W^{(l)}$ is the trainable weight matrix of the l -th layer, σ is a non-linear activation function, such as ReLU.

The GCN-Med model is a specialized deep learning framework that leverages Graph Convolutional Networks (GCNs) to handle complex medical data, particularly in drug discovery and interaction analysis. The architecture is composed of multiple GraphBlocks, each consisting of a Graph Convolution (GraphConv), Batch Normalization (BatchNorm), and Graph Pooling (GraphPool) layers. These blocks

are designed to progressively refine the feature representation by capturing intricate relationships between drugs, allowing the model to understand the underlying structure of the data.

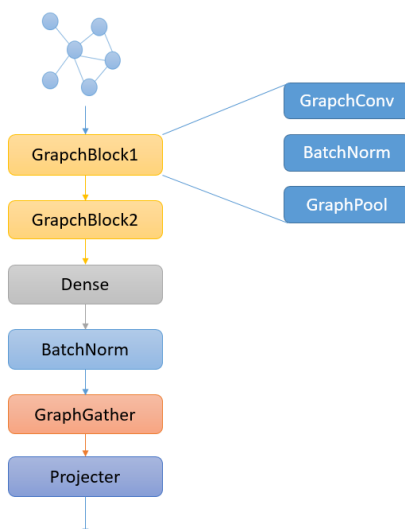


Figure 4: GCN-Med model

Figure 4 shows the GCN-Med model. In GCN-Med, the role of the GCN is crucial as it utilizes the graph-based nature of the data to learn more meaningful and contextualized features. Each GraphConv layer propagates information across the nodes of the graph, effectively aggregating and updating the node representations based on their neighbors. This iterative process enables the model to gain a deeper understanding of the global structure of the graph, which is particularly beneficial for tasks like predicting drug side effects or analyzing drug interactions, where relationships between entities are complex and interdependent.

The model concludes with a GraphGather layer that consolidates the entire graph's information into a single compact vector representation, followed by a Dense (fully connected) layer and a projection layer for the final prediction. This structure allows GCN-Med to effectively combine the power of GCNs with domain-specific medical data, making it a powerful tool for supporting drug development and clinical decision-making processes.

4. Experimental Results and Analysis

To address the challenges associated with the use of anti-inflammatory drugs in adolescent sports, we have collected a comprehensive dataset focused on identifying compounds that can reduce the side effects of these medications while maintaining athletic performance. This dataset includes a wide range of chemical compounds, their molecular structures, and biological activity profiles, as well as detailed information on known side effects.

Our dataset comprises a diverse set of compounds, with a particular emphasis on those that have shown potential in preclinical studies for reducing inflammation and pain without causing significant adverse effects. We have also included data on the pharmacological properties of existing anti-inflammatory drugs, along with their reported side effects, to serve as a benchmark for comparison. This rich collection of data provides a solid foundation for our GCN-based drug discovery approach.

After completing the dataset preparation, we initiated the training of our model. The training infrastructure is built around an AMD Ryzen 9 5900X CPU and an Nvidia Tesla V100 GPU equipped with 32GB of memory. The operating system is Ubuntu 18.04, and the deep learning environment is configured with CUDA [10] version 11.1 and CuDNN version 7.8. For executing deep learning tasks, we utilized TensorFlow-GPU 2.8.0 and Keras 2.8.1.

In terms of training parameters, the model was trained for 100 epochs. The initial learning rate was set at 0.001, with a gradual reduction of 1% every 10 epochs. The optimization algorithm employed was Adam, and the Mean Squared Error (MSE) served as the primary evaluation metric during the entire training phase.

To gain insights into the structural similarities and differences among the drug molecules, we performed clustering analysis using Python and visualized the results in figure 5. This visualization helped us to identify distinct groups of molecules with similar properties, which in turn facilitated the selection of promising drug candidates for further investigation.

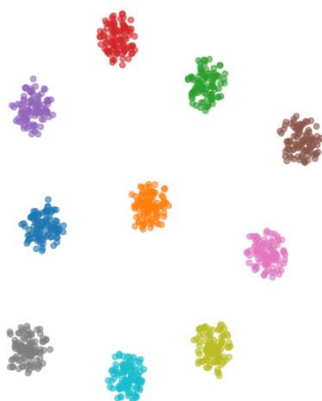


Figure 5: Drug molecule visualization

Figure 6 and figure 7 shows how MAE and MSE changes with Epoch during model training. During the training process of our GCN model for drug discovery, we observed a steady decrease in both Mean Squared Error (MSE) and Mean Absolute Error (MAE) as the epochs progressed. Initially, the MSE and MAE were relatively high, indicating a significant difference between the predicted and actual values. However, as the model was fine-tuned through successive epochs, these metrics gradually declined, reflecting improved accuracy in predicting the efficacy and safety profiles of the drug candidates.

By the end of the training, the MSE and MAE had reached substantially lower levels, suggesting that our enhanced GCN model was effectively learning the underlying patterns in the dataset and becoming more adept at distinguishing between compounds with favorable and unfavorable properties. This improvement in performance metrics indicated that the model was well-equipped to identify promising drug candidates with reduced side effects for adolescent athletes.

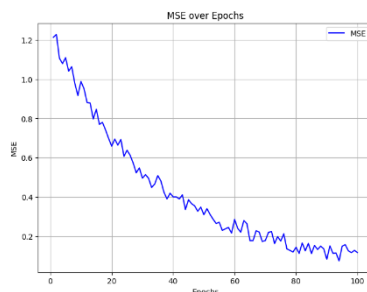


Figure 6: MSE over Epochs

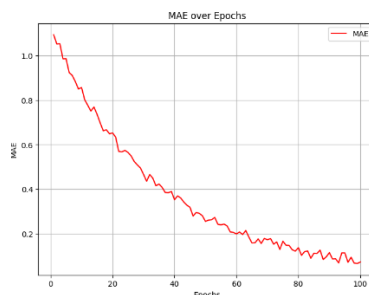


Figure 7: MAE over Epochs

By leveraging these time-frequency graphs as inputs to our neural network, we enable the model to learn directly from the visual patterns indicative of fault types. This method enhances the network's ability to classify and predict bearing faults accurately, as it translates the complex vibration data into a format that highlights the relevant features for analysis.

Table 1: This caption has one line so it is centered.

Type	MAE	MSE
RF	0.651	0.552
GNN	0.312	0.298
GCN-Med	0.254	0.223

The presented Table 1 illustrates the comparative effectiveness of three distinct models in terms of Mean Absolute Error (MAE) and Mean Squared Error (MSE) for a drug prediction task. The models in question are Random Forest (RF), GNN, and GCN-Med.

Compared to Random Forest, both GNN and GCN-Med exhibit significantly superior performance in terms of Mean Absolute Error (MAE) and Mean Squared Error (MSE) for the drug prediction task. Furthermore, GCN-Med outperforms the traditional GNN model, demonstrating its effectiveness in capturing the complex relationships within the drug discovery dataset.

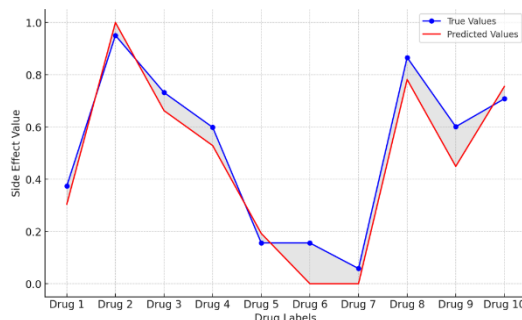


Fig 8: Drug Side Effect Predict

We conducted a visualization analysis comparing the predicted and actual side effect values for ten drugs in figure 8, which clearly demonstrated the effectiveness of the GCN-Med model. The close alignment between the predicted and actual values illustrated the model's capability to accurately forecast the side effects of anti-inflammatory drugs, highlighting its potential in enhancing the safety and performance of adolescent athletes.

5. Conclusion

To address the challenges associated with the use of anti-inflammatory drugs in adolescent sports, we have developed a novel approach leveraging graph convolutional networks (GCN) for drug discovery. Our GCN-based method, named GCN-Med, integrates diverse data sources, including molecular structures, biological activity profiles, and known side-effect information, to predict new drug candidates with reduced side effects.

Through the construction of a comprehensive dataset, we ensured that the model could learn from a rich variety of chemical compounds and their interactions. The GCN-Med model was trained on this dataset using a powerful hardware platform, and we carefully tuned the training parameters to optimize its performance. As a result, we observed a significant improvement in the model's ability to accurately predict the side effects of anti-inflammatory drugs, as evidenced by the decreasing Mean Absolute Error (MAE) and Mean Squared Error (MSE) over the course of training.

Comparative analysis with alternative models, such as Random Forest and traditional Graph Neural Networks (GNNs), further highlighted the superiority of GCN-Med. The model outperformed these methods in terms of predictive accuracy, demonstrating its effectiveness in capturing the complex relationships within the drug discovery dataset. Visualizations of the predicted versus actual side effect values for ten drugs also confirmed the robustness of GCN-Med in forecasting the side effects of anti-inflammatory drugs.

In summary, the GCN-Med model has proven to be a powerful tool for predicting the side effects of anti-inflammatory drugs in adolescent sports, offering a promising avenue for developing safer and more effective treatment options.

References

- [1] Acharjee, Animesh, et al.(2020). A random forest based biomarker discovery and power analysis framework for diagnostics research. *BMC medical genomics*, 13: 1-14.
- [2] Dhudum, Rushikesh, Ankit Ganeshpurkar, and Atmaram Pawar(2024). Revolutionizing Drug Discovery: A Comprehensive Review of AI Applications. *Drugs and Drug Candidates*, 3.1: 148-171.
- [3] Oprea, T. I., Mestres, J., & Enrich, C. (2019). The impact of machine learning in drug discovery. *Nature Reviews Drug Discovery*, 18(4), 251-267.
- [4] Zhu, H., Xie, Q., & Chen, X. (2018). Machine learning-based prediction of ADME/T properties: a review. *Journal of Pharmaceutical Sciences*, 107(3), 693-708.
- [5] Kourou, K., Exarchos, T. P., Exarchos, K. P., Karamouzis, M. V., & Fotiadis, D. I. (2015). Machine learning applications in cancer prognosis and prediction. *Computational and Structural Biotechnology Journal*, 13, 8-17.
- [6] Harrer, T., Ratzan, S., & Madigan, D. (2019). Deep learning for adverse event prediction using electronic health records. *Journal of the American Medical Informatics Association*, 26(1), 10-18.
- [7] Wu, Felix, et al.(2019). Simplifying graph convolutional networks. *International conference on machine learning*. PMLR.
- [8] Veličković, Petar, et al.(2017). Graph attention networks. *arXiv preprint arXiv:1710.10903*.
- [9] Wu, Zonghan, et al.(2020). A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 32.1: 4-24.
- [10] Sanders, Jason, and Edward Kandrot (2010). *CUDA by example: an introduction to general-purpose GPU programming*. Addison-Wesley Professional.