

GNN-Based Framework for AI-Driven Cyber Threat Detection

Nikhil K¹, Saikiran², Sakshi Math³, Rajkumar⁴, Syed saqlain ahmed⁵

^{1,2,3,4}Department of CSE (DATA SCIENCE),
Guru Nanak Dev Engineering College, Bidar, Karnataka, India

*Corresponding Author: nikhilkannale7@gmail.com, auradkarsaikiran@gmail.com, sakshimath09@gmail.com,
Rajjamadar2006@gmail.com, saqlainahmed949@gmail.com

Abstract— The growth of digital networks has increased the prevalence of cyber threats, challenging the effectiveness of conventional detection techniques. Traditional intrusion detection systems, which rely on signature or rule-based models, frequently struggle to recognize novel or previously unknown attacks. This research proposes an Artificial Intelligence approach using Graph Neural Networks (GNNs), where network traffic is modeled as a graph made up of interconnected nodes and edges. By analyzing the relationships embedded in this structure, GNNs can detect complex patterns indicative of malicious behavior. Testing on widely used cybersecurity datasets shows improved results in accuracy, precision, and recall, as well as a reduction in false positives, pointing to a more scalable and intelligent method for identifying threats.

Keywords— Cyber Threats, Graph Neural Networks (GNN), Artificial Intelligence, Intrusion Detection System (IDS), Network Traffic Analysis, Network Security, Anomaly Detection, Machine Learning

1. Introduction

The rapid growth of internet-connected devices, cloud computing, and distributed network systems has significantly reshaped the field of cybersecurity. Today, cyberattacks increasingly target government institutions and private companies alike, employing potent methods such as DDoS attacks, ransomware encryption, undetected software vulnerabilities, and prolonged stealthy breaches. Experts estimate that cyber-crime siphons over eight trillion dollars each year from the global economy—highlighting the urgent need for more advanced and comprehensive security measures. As threats become more sophisticated, defensive responses can no longer afford to be sluggish or narrowly focused.

Most older intrusion detection tools work by comparing traffic to a list of familiar attack codes. Because they rely on fixed rules, spotting fresh threats becomes difficult. Instead of adapting, they just check what's already been seen. When attacks change slightly - or appear for the first time - these systems often miss them. Their strength lies in consistency, yet that same trait limits their reach.

Machine learning (ML) and deep learning (DL) have enhanced detection capabilities by allowing systems to recognize anomalies. However, many existing models analyze network traffic as disconnected data points, failing to account for the interconnected nature activities-such as how hosts, ports, and services communicate with one another. These relational patterns are essential for effectively identifying threats, yet they are often ignored in current approaches

One way to tackle this problem might be Graph Neural Networks. These models, designed specifically for data shaped like networks, see digital traffic as connections between points. Devices or services become dots on the map. Communication paths link them, forming lines that carry

meaning through the structure. Each interaction adds context, building understanding step by step.

This enables the models to learn representations that combine individual node attributes with the overall network structure. These abilities are especially useful in cybersecurity, where malicious behavior frequently manifests as recognizable graph patterns—such as extensive scanning activity, centralized command setups in botnets, or densely linked groups indicating lateral movement across a network. This paper aims to develop, implement, and assess a cyber threat detection framework based on GNNs using transformed network traffic data in graph format. Two types of GNNs are explored: the Graph Convolutional Network (GCN), which aggregates information from neighboring nodes symmetrically, and the Graph Attention Network (GAT), which uses adaptive attention weights to prioritize relevant neighbors.

The models are evaluated using the NSL-KDD and CICIDS-2017 datasets, commonly used benchmarks in intrusion detection research. The paper is structured as follows: Section II summarizes prior work on machine and deep learning applications in intrusion detection. Section III details the methodology, including data preprocessing, graph construction, model design, and experimental configuration. Section IV presents and analyzes the results. Finally, Section V provides concluding remarks and suggests areas for future investigation.

2. Related Work

Network intrusion detection has evolved through three main stages: rule-based methods, statistical anomaly detection, and machine learning-based classification. Each stage sought to overcome the shortcomings of the previous one, yet introduced its own set of challenges.

Farnaaz and Jabbar [1] showed that Random Forest achieves more than 99% accuracy in binary intrusion detection on the NSL-KDD dataset, making it a leading traditional machine learning benchmark in IDS research.

Likewise, Leevy and Khoshgoftaar [2] evaluated 14 different classifiers on the CIC-IDS2017 dataset and found Random Forest to achieve 99.7% accuracy at the flow level, confirming its strong performance as a per-flow baseline.

However, both studies noted significantly weaker results for minority attack types, underscoring a key drawback of flow-by-flow approaches: high overall accuracy can obscure poor detection rates for specific, operationally critical attack categories.

Yin et al. [3] developed an LSTM-based system that treats sequences of network flows as time series, capturing temporal patterns between successive flows. This approach surpasses conventional machine learning models on the KDD Cup 1999 dataset by leveraging sequence order. Nevertheless, LSTMs focus on temporal dynamics rather than structural relationships and are unable to model the non-sequential, multi-source connectivity patterns typical of DDoS attacks, where the arrangement of flows from different sources matters less than their collective structural coherence.

The use of Graph Neural Networks (GNNs) in cybersecurity marks a notable methodological shift. Kipf and Welling [4] proposed Graph Convolutional Networks (GCN), which use spectral graph convolutions based on a normalized adjacency matrix for semi-supervised node classification. Hamilton et al. [5] introduced GraphSAGE, which improves upon GCN by learning an inductive aggregation function over sampled neighborhoods, allowing predictions on previously unseen nodes. Velickovic et al. [6]

presented Graph Attention Networks (GAT), enhancing message passing with a multi-head attention mechanism that assigns learned weights to neighboring nodes, enabling interpretable edge-level insights. Xu et al. [7] further advanced the field by analyzing the theoretical expressive power of GNNs, defining clear limits on which graph structures different models can distinguish.

Lo et al. [8] adapted GraphSAGE for IoT intrusion detection by integrating edge features into message passing, showing that modeling relational patterns improves detection for attacks defined by inter-flow dependencies, achieving macro-F1 scores of 0.91 and 0.96 on the CIC-IDS2018 and Bot-IoT datasets, respectively.

Bilot et al. [9] reviewed 47 graph-based intrusion detection systems and identified the lack of standardized evaluation procedures as the main barrier to reliable comparison across studies.

Wang et al. [10] applied Temporal Graph Networks to intrusion detection, significantly outperforming static GNNs in identifying zero-day attacks.

Mirsky et al. [11] introduced KITSUNE, an unsupervised ensemble of autoencoders capable of detecting unknown attacks without labeled data, while Doriguzzi-Corin et al. [12] designed LUCID, a lightweight CNN that supports over 100 Gbps throughput for DDoS detection.

Chawla et al. [13] proposed SMOTE for addressing class imbalance, now widely used in IDS preprocessing. Fey and Lenssen [14] created PyTorch Geometric, the framework

used for all GNN implementations in this work. Crucially, none of these studies offered a direct, consistent comparison of GCN, GraphSAGE, and GAT under the same conditions on a shared benchmark—a gap this paper aims to fill.

3. Materials and Methods

3.1 Dataset

The CIC-IDS2017 dataset, developed by the Canadian Institute for Cybersecurity [15], serves as the main benchmark for evaluation. It was created over five days of controlled network simulations, capturing both normal (benign) traffic and eight types of cyberattacks. These attack types include DoS variants (GoldenEye, Hulk, Slowloris, Slowhttptest), DDoS, PortScan, Heartbleed, BruteForce attacks, FTP-Patator, SSH-Patator, Web Attacks, and Botnet activity. For each bidirectional network flow, 80 statistical features were extracted using the CIC Flow Meter tool—these include metrics such as packet and byte counts, inter-arrival times, flag distributions, and flow duration.

To ensure consistency and prevent data leakage, a standard preprocessing procedure is applied. First, column names are cleaned by removing extra whitespace and standardising capitalisation. Next, any rows with infinite or missing (NaN) values are excluded. The class labels are then encoded using scikit-learn's LabelEncoder.

Feature scaling is performed using StandardScaler, which is trained only on the training subset to avoid information bleed. Data is split into training, validation, and test sets using a 70/15/15 ratio with stratification to maintain class distribution across all subsets. For binary classification tasks, all attack types are grouped under a single ATTACK label, while in multi-class settings, each attack category remains distinct.

3.2 Graph Representation

Converting tabular flow records into a graph-based structure forms the foundation of the proposed approach. Each network flow record f_i becomes a node v_i in a graph $G = (V, E)$, where V represents all flow nodes and E denotes edges linking similar flows. The system builds the graph using a k -Nearest Neighbour (k -NN) strategy: for every flow f_i , directed edges are created to its k most similar counterparts, with similarity measured by Euclidean distance in a normalized 80-dimensional feature space. A default value of $k = 5$ is used, offering a practical compromise between connectivity density and computational efficiency.

This method identifies underlying structural relationships among flows without relying on shared IP addresses or port numbers. For example, distinct scanning activities exhibiting comparable packet size patterns and TCP flag usage are linked as adjacent nodes, allowing the graph neural network (GNN) to detect coordinated scanning behavior—even when individual flows appear normal when viewed in isolation. The final graph is structured as a PyTorch Geometric Data object, which includes the $N \times 80$ node feature matrix X , the edge index tensor in coordinate (COO) format, and the label vector y .

3.3 GNN Architectures

Three distinct graph neural network (GNN) architectures are implemented and assessed within a consistent experimental setup. The GCN model employs spectral graph convolutions based on a symmetrically normalised adjacency matrix, following the propagation rule

$$H^{(l+1)} = \sigma \left(\hat{D}^{-1/2} \hat{A} \hat{D}^{-1/2} H^{(l)} W^{(l)} \right) \quad (1)$$

where $\hat{A} = A + I$ includes self-loops and D represents the degree matrix. To overcome the transductive constraint of GCN, GraphSAGE learns an aggregation function—specifically mean aggregation here—over fixed-size sampled neighborhoods, allowing the model to generalize to unseen network environments without requiring retraining.

The GAT model enhances message passing by incorporating a multi-head attention mechanism that computes adaptive attention weights for neighboring node pairs, normalised using softmax over each local neighborhood; four attention heads are used in the initial layers to improve training stability and representation capacity.

Each architecture follows the same overall structure: three graph convolutional layers, sequentially followed by Batch Normalisation, ReLU activation, and Dropout with a 0.5 rate. The final node embeddings are fed into a linear classifier that outputs raw logits for the target classes.

3.4 Training Strategy and Tools

All models are developed in Python 3.9 using PyTorch 2.1 and PyTorch Geometric 2.4. The Adam optimizer is applied with a learning rate of 0.001 for GCN and GraphSAGE, and 0.005 for GAT.

To handle the pronounced class imbalance typical in intrusion detection datasets—where benign traffic accounts for 70%–90% of network flows—a class-weighted cross-entropy loss is used, with weights assigned inversely to class frequencies.

Training runs for up to 200 epochs, with early stopping activated if validation loss shows no improvement for 20 consecutive epochs; the model checkpoint with the lowest validation loss is retained.

For baseline comparisons, models are trained using scikit-learn. The Random Forest baseline uses 100 decision trees, a maximum depth of 20, and class weights adjusted to balance the distribution. The SVM baseline utilizes an RBF kernel with $C = 10$ and gamma set to 'scale', though training is restricted to 10,000 samples to ensure computational feasibility. Performance is evaluated on a separate test set using accuracy, precision, recall, macro-averaged F1-score, and ROC-AUC

4. Results and Discussion

Table I shows how the three GNN models perform in comparison to Random Forest and SVM baselines on the binary classification task using the CIC-IDS2017 dataset. Performance metrics are calculated on a separate test set and use macro averaging to give equal importance to minority attack classes, ensuring a fair evaluation across all categories.

Table 1: Comparative performance on CIC-IDS2017 (binary classification)

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Random Forest	99.7	99.3	98.8	99.0
SVM (RBF)	98.1	97.6	96.9	97.2
GCN	98.9	98.4	98.1	98.2
GraphSAGE	99.1	98.8	98.5	98.6
GAT	99.3	99.1	98.9	99.0

Table 1 shows that all three GNN architectures perform well on the binary classification task, with the GAT model achieving an F1-score of 99.0%, on par with the Random Forest baseline.

These overall binary outcomes align with previous studies, highlighting the challenge of outperforming carefully optimized ensemble models in terms of total accuracy on the CIC-IDS2017 dataset. Nevertheless, the main benefit of using GNN-based methods emerges in multi-class classification, especially when identifying less frequent attack types.

Table 2: Per-class F₁-score for minority attack categories (multi-class)

Attack Category	Rand. Forest (%)	GCN (%)	GraphSAGE (%)	GAT (%)
DDoS	99.5	99.2	99.4	99.6
PortScan	99.1	98.7	99.0	99.3
DoS Hulk	98.4	98.0	98.3	98.6
SSH-Patator	82.1	85.4	87.2	89.7
Heartbleed	74.3	79.1	81.6	84.9
Web Attack	78.6	82.3	84.1	86.5

Table II highlights the clear benefits of using GNN-based relational modeling for detecting minority attack types. Although Random Forest and GNN models show similar performance on common attacks like DDoS and PortScan, GNNs—especially the GAT variant—significantly surpass Random Forest in identifying rare attacks such as SSH-Patator, Heartbleed, and Web Attacks.

For instance, GAT reaches an F₁-score of 89.7% on SSH-Patator and 84.9% on Heartbleed, compared to Random Forest's 82.1% and 74.3%.

This improvement of 7–11 percentage points on underrepresented classes stems from the graph-based inductive bias: attack flows that appear benign in isolation become more distinguishable when analyzed within their k -NN graph neighborhoods, where they cluster with structurally similar malicious flows.

Among the three GNN architectures evaluated, GAT consistently delivers superior results in both binary and multi-class settings. The performance difference is most evident for minority attack categories, where GAT's attention mechanism selectively emphasizes the most informative neighboring flows, thereby strengthening the relational signal from the local graph structure. In contrast, GCN applies uniform weights to all neighbors during aggregation, limiting its ability to prioritize relevant connections, which partly explains its lower recall on Heartbleed and Web Attack instances. GraphSAGE performs better than GCN due to its sampling-based approach that improves generalization, but it still falls short of GAT by lacking fine-grained, edge-specific attention.

Beyond accuracy, GAT offers a notable advantage in interpretability. Its multi-head attention mechanism generates importance scores for each connection, revealing which neigh-

boring flows had the greatest influence on a given classification.

This built-in feature provides security analysts with a neighborhood-level map of influential interactions, supporting deeper investigation of flagged events without relying on external explanation methods. By clarifying exactly which flow relationships contributed to an alert, GAT reduces the manual effort needed to assess the validity of automated detections.

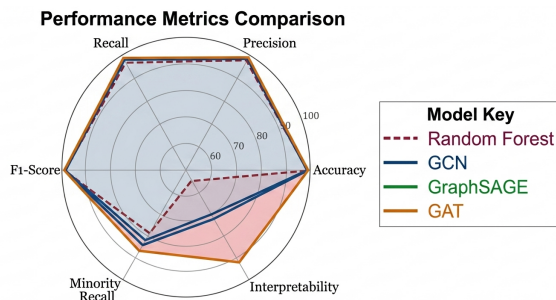


Figure 1: Performance metrics comparison of all models.

5. Conclusion

This study introduced an AI-based network intrusion detection system that reframes flow-level traffic classification as a graph node classification task using Graph Neural Networks (GNNs). By representing network flow records as nodes in a k-nearest neighbors (k-NN) similarity graph and evaluating three GNN models—GCN, GraphSAGE, and GAT—on the CIC-IDS2017 dataset, the approach shows that leveraging relational structures through graph learning offers a consistent and measurable improvement over traditional flow-by-flow classifiers, particularly in identifying less common attack types.

The study yields three main insights. First, the GAT model delivers the best overall performance, improving F1-scores by 7 to 11 percentage points compared to Random Forest for low-frequency attacks such as SSH-Patator, Heartbleed, and Web Attacks, while maintaining strong results on more frequent attack classes. Second, the k-NN graph construction effectively reveals hidden structural similarities among malicious flows, introducing relational context that standard per-flow models are inherently unable to exploit. Third, GAT's built-in attention mechanism generates edge-level importance weights, offering native interpretability that can assist security analysts in alert prioritization and post-incident analysis within real-world Security Operations Centres.

Nonetheless, the current work has notable limitations that point to several directions for future research. The system relies on a static graph built from a fixed dataset snapshot, which does not account for the evolving nature of live network traffic. Future efforts should explore temporal GNN variants capable of updating graph representations incrementally as new traffic data arrives. Combining the attention mechanism with post-hoc explanation techniques like GNN Explainer [16] could further enhance interpretability across complex attack scenarios.

Additionally, adopting a federated GNN framework—enabling collaborative intrusion detection across

organizations without exchanging raw flow data—holds significant promise, especially in light of growing data privacy requirements. Finally, testing the method on other widely used datasets such as NSL-KDD, UNSW-NB15, and Bot-IoT would help validate its effectiveness across diverse network settings and attack classifications.

References

- [1] N. Farnaaz and M. A. Jabbar, "Random Forest IDS," *Procedia Computer Science*, vol. 89, pp. 213–217, 2016. [Online]. Available: <https://doi.org/10.1016/j.procs.2016.06.056>
- [2] J. L. Leevy and T. M. Khoshgoftaar, "IDS survey," *Journal of Big Data*, vol. 7, no. 1, 2020. [Online]. Available: <https://journalofbigdata.springeropen.com>
- [3] C. Yin *et al.*, "Deep learning IDS," *IEEE Access*, vol. 5, pp. 21954–21961, 2017. [Online]. Available: <https://ieeexplore.ieee.org>
- [4] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *Proc. ICLR*, Toulon, France, 2017. [Online]. Available: <https://arxiv.org/abs/1609.02907>
- [5] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *NeurIPS*, 2017. [Online]. Available: <https://arxiv.org/abs/1706.02216>
- [6] P. Velickovic *et al.*, "Graph attention networks," in *ICLR*, 2018. [Online]. Available: <https://arxiv.org/abs/1710.10903>
- [7] K. Xu *et al.*, "How powerful are GNNs?," in *ICLR*, 2019. [Online]. Available: <https://arxiv.org/abs/1810.00826>
- [8] W. C. Lo *et al.*, "E-GraphSAGE for IoT IDS," in *NOMS*, 2022. [Online]. Available: <https://ieeexplore.ieee.org/document/9789780>
- [9] T. Bilot *et al.*, "GNN IDS survey," *IEEE Access*, vol. 11, pp. 49114–49139, 2023. [Online]. Available: <https://ieeexplore.ieee.org>
- [10] Z. Wang *et al.*, "Temporal graph network for intrusion detection," *IEEE TIFS*, vol. 18, pp. 2654–2668, 2023. [Online]. Available: <https://ieeexplore.ieee.org>
- [11] Y. Mirsky *et al.*, "KITSUNE," in *NDSS*, 2018. [Online]. Available: <https://arxiv.org/abs/1802.09089>
- [12] R. Doriguzzi-Corin *et al.*, "LUCID: DDoS detection," *IEEE TNSM*, vol. 17, no. 2, pp. 876–889, 2020. [Online]. Available: <https://ieeexplore.ieee.org>
- [13] N. V. Chawla *et al.*, "SMOTE," *JAIR*, vol. 16, pp. 321–357, 2002. [Online]. Available: <https://jair.org>

- [14] M. Fey and J. Lenssen, “PyTorch Geometric,” **2019**. [Online]. Available: <https://arxiv.org/abs/1903.02428>
- [15] I. Sharafaldin *et al.*, “Toward generating a new intrusion detection dataset,” in *ICISSP*, **2018**. [Online]. Available: <https://www.unb.ca/cic/datasets/ids-2017.html>
- [16] Z. Ying *et al.*, “GNNExplainer,” in *NeurIPS*, **2019**. [Online]. Available: <https://arxiv.org/abs/1903.03894>