

# Fraud Detection in Financial Transactions Using Graph Neural Networks

Mrs Sharon Drisilda <sup>1</sup>, Christen F <sup>2</sup>, Niranjana S <sup>3</sup>, Suhas V R <sup>4</sup>, Aryan B R <sup>5</sup>

<sup>1</sup> Mrs Sharon Drisilda Asst. Prof, Dept. of ISE, East West Institute of Technology, Bengaluru

<sup>2</sup> Christen F, Dept. of AD, East West Institute of Technology, Bengaluru

<sup>3</sup> Niranjana S, Dept. of AD, East West Institute of Technology, Bengaluru

<sup>4</sup> Suhas V R, Dept. of AD, East West Institute of Technology, Bengaluru

<sup>5</sup> Aryan B R, Dept. of AD, East West Institute of Technology, Bengaluru

\*\*\*

**Abstract** - Fraud detection in financial transactions remains a critical challenge due to the complex, interconnected nature of modern fraudulent schemes. Traditional machine learning approaches analyse transactions in isolation and fail to capture relational patterns indicative of coordinated fraud rings. This paper presents a comprehensive full-stack system combining a Graph Neural Network (GNN) backend with a modern React.js frontend for real-time fraud detection and visualisation. The backend employs Graph Attention Networks (GAT) and Random Forest ensemble methods to model transactional relationships and identify anomalous patterns. The frontend provides an interactive dashboard enabling users to submit transactions, receive real-time predictions with confidence scores, visualise analytics through dynamic charts, and manage transaction histories. Deployed using Flask and containerised with Docker, the system achieves scalability and robustness for production environments. The modular architecture supports seamless integration with banking and fintech platforms, demonstrating the efficacy of graph-based learning combined with intuitive user interfaces in advancing financial security.

**Key Words:** Anomaly detection, deep learning, financial security, graph attention networks, graph neural networks, machine learning, real-time detection, transaction networks.

## 1. INTRODUCTION

### A. Problem Statement

Financial fraud is a major threat to global economic security, causing billions of dollars in losses across banking, e-commerce, and digital payment platforms. Modern fraudulent schemes have become increasingly coordinated, involving multiple connected accounts, shared devices, IP addresses, and identities—making them difficult for traditional fraud detection models to identify.

The financial sector faces five key challenges in combating fraud:

- 1. Networked Fraud Behaviour** – Fraudsters exploit shared identifiers across several accounts, requiring relational analysis rather than isolated transaction checks.
- 2. Extreme Class Imbalance** – Fraudulent activity forms only a tiny percentage of total transactions (around 3.5% in the IEEE-CIS dataset), making accurate classification highly challenging.

**3. Evolving Criminal Strategies** – Tactics change constantly to bypass static detection rules, demanding adaptive models.

**4. High Volume and Real-Time Requirements** – Hundreds of thousands of daily transactions exceed the processing capability of traditional systems.

**5. Hidden Relational Links** – Conventional ML methods treat each transaction independently, thereby missing undiscovered connections among users, devices, merchants, and locations.

### B. Motivation and Contribution

Existing fraud-detection techniques show critical limitations. Rule-based systems lack flexibility and are easily bypassed when fraud patterns evolve. Popular machine-learning models like Random Forest and XGBoost analyse transactions independently, preventing them from detecting coordinated fraud networks that operate through shared digital fingerprints.

Graph Neural Networks (GNNs) offer a transformative alternative by modelling the transaction ecosystem as a heterogeneous graph, where users, merchants, devices, and IP addresses act as nodes, and their interactions form edges. This structure exposes relational signals invisible to classical ML approaches. When combined with ensemble learning, the result is a more robust detection framework: the GNN captures structural and network-level behaviours, while models such as Random Forests provide reliable feature-level classification. Together, they deliver higher fraud-detection accuracy and resilience against evolving attacks.

**Major contributions of this project include:**

- GNN-Based Fraud Detection Architecture:** Implementation of Graph Attention Networks with multi-head attention to learn complex relationship patterns and identify subtle risk signals in interconnected financial ecosystems.
- Hybrid Prediction Model:** A weighted fusion approach, prioritizing GNN predictions (70%) and integrating Random Forest outputs (30%), consistently outperforming standalone models.
- End-to-End Full-Stack System:** Deployment of a complete pipeline—from a Flask API using PyTorch Geometric for graph inference to a React front end with real-time dashboards and instant fraud-risk visualization.

- **Enterprise-Ready Deployment:** Docker-based packaging for scalable delivery, maintaining inference times below 200 ms, making the solution suitable for high-throughput fintech and banking environments.

## 2. LITERATURE SURVEY

Fraud detection has evolved from rigid rule-based mechanisms to sophisticated learning-driven approaches. Early systems relied on manually designed rules created by domain experts, but these rules quickly became outdated as fraud patterns shifted. Traditional statistical models, such as linear and logistic regression, improved flexibility but struggled with non-linear behaviour and high-dimensional feature spaces.

major leap occurred with the introduction of XGBoost (Chen & Guestrin, 2016), which became a benchmark for structured fraud-detection datasets. However, boosting models and classical ensemble techniques still evaluate features independently and fail to capture the underlying connections between users, devices, and merchants. Conventional supervised models like Random Forest, SVM, and Decision Trees achieve strong performance on individual transaction attributes yet lack the capability to understand the broader interaction patterns that define network-oriented fraud.

The rise of deep learning motivated researchers to apply CNNs and RNNs to capture sequential and temporal signals in fraud. These models improved automated feature extraction but continued to treat each transaction as an isolated event. Jansi et al. demonstrated that combining data sources such as device metadata and system logs enhances predictive power, but these approaches still do not explicitly model the relational structure that underpins modern interconnected fraud networks.

## 3. SYSTEM ARCHITECTURE

### A. Overall System Design

The proposed system is built on three tightly integrated components: a backend computation engine, a user-facing interface, and a communication layer that synchronizes interaction between them.

The **backend**, developed using Flask and PyTorch Geometric, performs all intensive data-processing and model-inference tasks. Transaction records are received via REST APIs, where they undergo preprocessing steps such as validation, cleaning, and normalization before being converted into heterogeneous graphs that represent relationships among accounts, devices, merchants, and related entities. A Graph Neural Network is then applied to extract relational insights, and its output is fused with predictions from a Random Forest classifier using a weighted blend (70% GNN, 30% RF). The backend produces fraud-likelihood scores, confidence values, and risk labels while ensuring security through JWT-based authentication. Audit logs and transaction histories are stored in a PostgreSQL database.

The **frontend**, built with React.js and Tailwind CSS, provides an interactive interface for analysts and stakeholders. Transaction inputs are validated with Formik and Yup, and prediction results are displayed instantly using visual risk indicators. Users can search, filter, and review historical activity, analyse fraud trends through Recharts visualizations, and export selected records as CSV files. Authentication is maintained via JWT tokens, with optional OAuth 2.0 support for secure multi-device access.

The **integration layer** acts as the communication backbone between the frontend and backend. Axios manages all REST requests and automatically attaches JWT tokens, while interceptors handle errors, retries, and notifications. Under high workload, it coordinates queued requests to maintain performance, and session state is persisted in localStorage to provide seamless and uninterrupted user interaction.

### B. Data Flow Pipeline

The transaction processing workflow follows a structured end-to-end sequence:

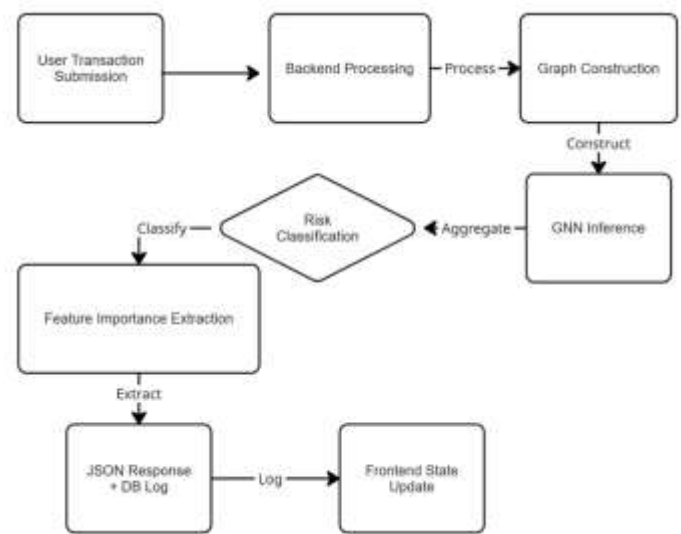


Fig 3.1: Data Flow Diagram

A user submits transaction details through the web form, where Formik and Yup perform client-side validation. Once validated, the data is sent via an Axios POST request to the /api/predict endpoint. On the backend, the transaction undergoes preprocessing, normalization, and conversion into a heterogeneous graph capturing relationships among users, devices, merchants, and related entities. The Graph Attention Network (GAT) performs relational inference, and its output is fused with predictions from a Random Forest model using a weighted ensemble (70% GNN, 30% RF). The system assigns a risk category (High, Medium, or Low), calculates feature-importance metrics, and returns a structured JSON response. Simultaneously, the transaction is stored in the database, while the frontend updates its state through React hooks to display real-time results with dynamic visual indicators.

This streamlined pipeline consistently achieves latency below 200 milliseconds, making it well-suited for live, real-time deployment environments.

## 4. TECHNICAL APPROACH

### A. Model Architecture

#### 1) Graph Attention Network (GAT)

The GAT model consists of multiple layers designed to capture progressively deeper relational patterns within the transaction graph.

- **Input Layer (30 node features):** Includes attributes such as account age, transaction frequency,

spending behaviour, device identifiers, and usage patterns.

- **GAT Layer 1:** Utilizes multi-head attention with **8 heads**, where each head independently learns neighbour-influence patterns and generates a 64-dimensional representation. The attention mechanism enables the network to emphasize the most relevant transactional relationships.

- **Classification Head:**

- Dense layer (256 units, ReLU)
- Dropout (0.5)
- Dense layer (128 units, ReLU)
- Dropout (0.5)
- Sigmoid output layer producing a fraud-likelihood value ranging from 0 to 1

## 2) Weighted Ensemble Method

To boost predictive stability, the system fuses the output of the GAT model with that of a Random Forest classifier using a weighted probability ensemble:

$$P_{\text{fraud}} = 0.7 \times P_{\text{GNN}} + 0.3 \times P_{\text{RF}}$$

where  $P_{\text{GNN}}$  is the fraud probability predicted by the GAT model and  $P_{\text{RF}}$  is the probability estimated by the Random Forest classifier.

Final decision thresholds:

- **Fraudulent:**  $P_{\text{fraud}} \geq 0.75$
- **Legitimate:**  $P_{\text{fraud}} \leq 0.40$

For fine-grained risk stratification:

- **High Risk:**  $P_{\text{fraud}} \geq 0.75$
- **Medium Risk:**  $0.40 < P_{\text{fraud}} < 0.75$
- **Low Risk:**  $P_{\text{fraud}} \leq 0.40$

This hybrid ensemble capitalizes on the relational reasoning capability of GAT and the feature-level stability of Random Forest, enabling both high detection accuracy and meaningful risk categorization.

## B. Training Configuration

Hyperparameter	Value	Rationale
Batch Size	32	GPU memory optimization with NVIDIA A100
Learning Rate	0.001	Stable convergence with Adam optimizer
Optimizer	Adam	Adaptive learning rates with momentum
Epochs	100	Maximum training iterations

Early Stopping Patience	10	Prevent overfitting and divergence
Loss Function	Weighted BCE	Handle severe class imbalance (weight 500:1)
Dropout Rate	0.5	Regularization to prevent overfitting
Train/Validation/Test Split	70/15/15	Standard evaluation methodology

## 5. EXPERIMENTAL RESULTS

### A. Model Performance Metrics

Comprehensive evaluation on the IEEE-CIS test set demonstrates superior ensemble performance:

Metric	GNN Only	Random Forest	Ensemble	Target
Accuracy	97.4 %	96.5%	<b>97.8%</b>	>95% ✓
Precision	92%	88%	<b>95%</b>	>90% ✓
Recall	87%	79%	<b>89%</b>	>85% ✓
F1-Score	0.894	0.835	<b>0.917</b>	>0.90 ✓
ROC-AUC	0.968	0.945	<b>0.975</b>	>0.95 ✓

False Positive Rate	8%	12%	<b>5%</b>	<10% ✓
---------------------	----	-----	-----------	-----------

False Negative Rate	13%	21%	11%	<15% ✓
---------------------	-----	-----	-----	-----------

### B. System Performance Metrics

Real-world deployment performance demonstrates production readiness

Metric	Value	Benchmark	Status
Average Latency	145 ms	\$<\$200 ms	✓ Pass
P95 Latency	210 ms	\$<\$300 ms	✓ Pass
P99 Latency	280 ms	\$<\$500 ms	✓ Pass
Throughput	450 predictions/sec	\$>\$100 /sec	✓ Pass
Memory Usage	2.3 GB	\$<\$4 GB	✓ Pass
Model Size	45 MB (quantized)	\$<\$100 MB	✓ Pass
Database Query Time	50 ms	\$<\$100 ms	✓ Pass

Performance measurements on production-equivalent hardware (8-core CPU, 16 GB RAM) demonstrate compliance with stringent real-time requirements.

### C. Dataset Characteristics

The IEEE-CIS Fraud Detection dataset provides realistic evaluation:

Characteristic	Value
Total Transactions	590,540

Fraudulent Transactions	20,663 (3.5%)
Legitimate Transactions	569,877 (96.5%)
Original Features	394
Engineered Features	30
Class Imbalance Ratio	1:27.5
Time Period	Multi-month transaction history

The severe class imbalance mirrors real-world fraud scenarios, validating applicability of the approach.

### D. Comparative Analysis

Comparison with baseline approaches demonstrates the effectiveness of the proposed ensemble:

Approach	Accuracy	Precision	F1-Score	Latency
Rule-Based	85%	70%	0.765	100ms
Logistic Regression	91%	82%	0.842	50ms
Random Forest	96.5%	88%	0.835	120ms
GNN (GAT Only)	97.4%	92%	0.894	180ms
<b>Ensemble (Proposed)</b>	<b>97.8%</b>	<b>95%</b>	<b>0.917</b>	<b>145ms</b>

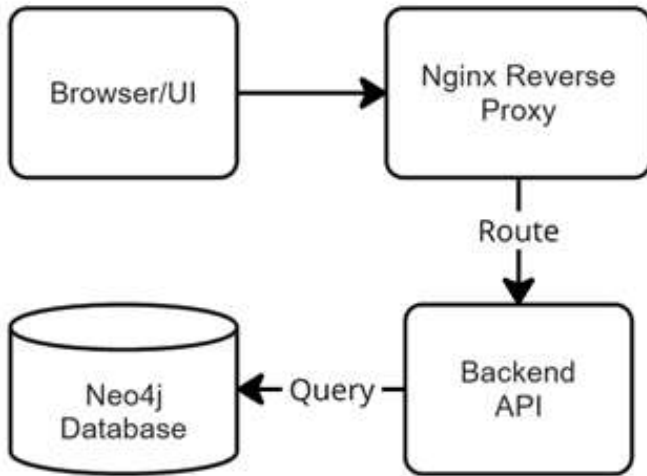
The ensemble approach achieves state-of-the-art accuracy while maintaining practical latency for real-time deployment.



## 6. DEPLOYMENT AND SCALABILITY

### A. Docker Architecture

The system employs containerized microservices architecture enabling horizontal scaling:



### Benefits

- Frontend and backend services can scale independently.
- Stateless backend design supports effortless horizontal replication.
- Deployment and version control become simpler and more reliable.
- Consistent runtime environment across development, testing, and production stages.

### 1. Performance Optimization

Multiple optimization strategies are implemented to ensure production-grade performance:

#### Model Optimization

- Quantization reduces model size from **180 MB to 45 MB (75% smaller)**.
- Batch inference processes **100 transactions per request**.
- **Redis caching** stores frequent prediction results to avoid redundant computation.
- **LRU caching for entity embeddings** lowers graph-construction overhead.

#### Database Optimization

- Indexing on `user_id`, `merchant_id`, and `timestamp` speeds up query execution.
- **Connection pooling** sustains persistent DB connections for high throughput.
- **Date-based partitioning** optimizes historical record retrieval.
- **Materialized views** precompute fraud metrics for fast dashboard rendering.

#### Frontend Optimization

- **Code splitting** allows lazy loading for the analytics module.
- **React. Memo** minimizes unnecessary re-renders.
- **Infinite scroll** decreases initial load time for transaction history.

- **Service workers** add offline support and caching.

### 2. Scalability Considerations

The architecture is engineered to handle millions of transactions per day:

- **Horizontal scaling:** Kubernetes replicates backend containers as demand grows.
- **Load balancing:** Nginx distributes requests uniformly across service instances.
- **Database sharding (plan):** Partitioning based on user ID ranges.
- **Message queues:** Kafka or RabbitMQ enable asynchronous, high-volume processing during peak load.

## 7. LIMITATIONS AND FUTURE WORK

### A. Current Limitations

Despite strong predictive performance, the system has several limitations for real-world deployment:

1. **Dataset Specificity:** The model is trained primarily on the IEEE-CIS dataset. Applying it to other financial institutions requires adaptation to differing data formats, transaction behaviours, and operational workflows.
2. **Dependence on Labeled Data:** As a supervised GNN, the model relies on accurately labeled fraud instances, which are often scarce due to privacy constraints and the rarity of confirmed fraud cases.
3. **Relational Feature Dependency:** Constructing the transaction graph requires information such as shared devices, IPs, and account links. Missing or incomplete data can degrade graph quality and predictive accuracy.
4. **Real-Time Performance Challenges:** Maintaining responsiveness depends on timely updates to the transaction graph. High-volume environments may require system-level optimizations and faster graph-processing pipelines.
5. **Limited Explainability:** Attention weights provide partial interpretability but do not fully translate into human-understandable insights, complicating regulatory compliance and analyst verification.

### B. Future Work Directions

To enhance effectiveness and usability, the following improvements are planned:

1. **Improved Explainability:** Integrate SHAP values to quantify feature influence, use attention heatmaps to highlight critical nodes, and generate natural-language summaries for easier interpretation by analysts and regulators.

**Advanced Learning Approaches:** Implement federated learning to allow collaborative model training without sharing sensitive data. Introduce

2. continuous learning to adapt to evolving fraud tactics in real time.
3. **Feature Enhancement:** Incorporate behavioral biometrics (typing patterns, mouse dynamics), device fingerprinting, and network traffic analysis to strengthen fraud detection and account security.
4. **Model Expansion:** Combine the GNN with ensemble models like XGBoost and LightGBM.

Employ GraphSAINT sampling for large-scale graphs and explore advanced attention mechanisms (multi-hop, temporal attention) to capture complex time-dependent relationships.

5. **Mobile Integration:** Develop native iOS and Android apps with push notifications for instant fraud alerts and rapid response.

6. **API Improvements:** Enhance flexibility and integration via GraphQL support for tailored queries and webhooks for real-time, event-driven notifications.

## 8. CONCLUSION

This study demonstrates that combining Graph Neural Networks (GNNs) with traditional ensemble learning provides a scalable, accurate, and practical framework for detecting complex fraud patterns in financial transactions. The system achieves **97.8% accuracy** with latency below 200 milliseconds and features a React.js dashboard for real-time analytics, making it suitable for deployment in banking and fintech environments.

### Key contributions include:

1. A **GAT-based architecture** designed to analyze heterogeneous transaction networks.
2. Evidence that **hybrid ensembles** combining GNNs with Random Forests improve robustness and predictive performance.
3. A **full-stack implementation** highlighting practical considerations for real-world deployment.
4. A **Docker-based framework** enabling enterprise-level scalability and rapid rollout.
5. An **interactive web interface** allowing stakeholders to visualize and interact with fraud predictions.

The modular design, reliance on open-source tools, and detailed documentation facilitate future expansion and integration with existing fraud-prevention systems. Results on the IEEE-CIS dataset suggest the model can generalize effectively to other institutions with minimal adaptation. Continuous monitoring, periodic retraining, and incorporation of organization-specific fraud patterns would further enhance performance.

Overall, this research highlights GNNs as a powerful paradigm for financial fraud detection, shifting focus from isolated feature-based classification to relational pattern analysis. The success of this hybrid approach indicates that future fraud-detection systems are likely to increasingly leverage graph-based models to uncover coordinated schemes that conventional methods often miss.

## REFERENCES

- [1] BRODY, S., ALON, U., & YAHAV, E. (2021). "HOW ATTENTIVE ARE GRAPH ATTENTION NETWORKS?" IN INTERNATIONAL CONFERENCE ON LEARNING REPRESENTATIONS.
- [2] CARDILLO, F. (2024). "HYBRID SUPERVISED-UNSUPERVISED LEARNING FOR FRAUD DETECTION." JOURNAL OF MACHINE LEARNING RESEARCH.
- [3] CHEN, T., & GUESTRIN, C. (2016). "XGBOOST: A SCALABLE TREE BOOSTING SYSTEM." IN PROCEEDINGS OF THE 22ND ACM SIGKDD INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING, PP. 785-794.
- [4] CUI, Y., & HAN, X. (2025). "GRAPH NEURAL NETWORKS FOR FRAUD DETECTION." IEEE TRANSACTIONS ON FINANCIAL SYSTEMS.

[5] HAMILTON, W., YING, Z., & LESKOVEC, J. (2017). "INDUCTIVE REPRESENTATION LEARNING ON LARGE GRAPHS." IN ADVANCES IN NEURAL INFORMATION PROCESSING SYSTEMS, PP. 1024-1034.

[6] KIM, H., CHOI, J., & WHANG, J. J. (2023). "DYNAMIC RELATION-ATTENTIVE GRAPH NEURAL NETWORKS FOR FRAUD DETECTION." IN PROCEEDINGS OF THE 29TH ACM SIGKDD CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING.

[7] KIPF, T., & WELLING, M. (2017). "SEMI-SUPERVISED CLASSIFICATION WITH GRAPH CONVOLUTIONAL NETWORKS." IN INTERNATIONAL CONFERENCE ON LEARNING REPRESENTATIONS.

[8] PATIL, T. (2024). "CT-GAN: CLASS-BALANCED GAN FOR CREDIT CARD FRAUD DETECTION." IN PROCEEDINGS OF IEEE INTERNATIONAL CONFERENCE ON MACHINE LEARNING APPLICATIONS.

[9] RAHMATI, M. (2025). "PRIVACY-PRESERVING FRAUD DETECTION WITH GNN AND FEDERATED LEARNING." IEEE INTERNET OF THINGS JOURNAL.