

Insurance Fraud AI Autonomous Adjudication using Graph Neural Networks, RAG, and Agentic Orchestration: Beyond Claim Black Boxes

1st Hassan Daoud

ATU Research Lab

20.November.2025

Github: <https://github.com/H-Daoud>

Linkedin: <https://www.linkedin.com/in/daoud1001>

Abstract—Accurate and efficient adjudication of insurance claims is essential for financial stability and customer trust. However, traditional machine learning models (such as Random Forest or XGBoost) face significant challenges in handling the complex, relational nature of organized fraud rings and the unstructured logic of policy documents. Consequently, these “black box” models often result in high false-positive rates or fail to detect sophisticated collusion, leading to significant financial leakage. This research addresses the critical challenge of autonomous claim investigation by developing a robust “Compound AI System” that integrates Graph Neural Networks (GNNs) for network-based anomaly detection with Retrieval-Augmented Generation (RAG) for policy verification. The study utilizes a multi-agent architecture where specialized agents—a Detective (GNN), a Lawyer (RAG), and a Judge (LLM)—collaborate to render a decision. The chosen approach demonstrates its ability to handle heterogeneous graph data and provide explainable reasoning, while the comparative method offers insights for evaluation purposes. The proposed Agentic framework demonstrated superior performance, achieving a fraud detection recall of 89 percent on synthetic fraud rings, significantly outperforming tabular baselines. These results pave the way for the next generation of Cognitive Automation in regulated industries, providing a foundation for future advancements in explainable AI (XAI).

By leveraging this advanced multi-agent technique, this research enhances precision in fraud detection, enabling timely and legally defensible decisions. This research contributes to the field of AI Engineering by introducing a robust methodology for orchestrating symbolic reasoning (Policies) with subsymbolic learning (GNNs) within a unified framework.

Index Terms—Graph Neural Networks, Agentic AI, Fraud Detection, Retrieval-Augmented Generation, Cognitive Automation, Insurance Technology.

I. INTRODUCTION

Artificial Intelligence has emerged as a transformative force in the financial services sector, offering great opportunities for automation, risk assessment, and fraud detection at the systemic level. Among its many applications, automated claim adjudication stands out for its potential to revolutionize the insurance industry by enabling real-time decision-making, significantly improving customer experience in terms of settlement speed, accuracy, and reduced operational costs [1]. However, despite these promising capabilities, the practical implementation of autonomous adjudication faces significant

challenges, particularly in accurately detecting “fraud rings” within the complex and dynamic environment of claim networks [2]. These challenges raise crucial questions about how to improve algorithmic reasoning to unlock the full potential of AI in regulated markets.

Traditional anomaly detection techniques, such as rule-based engines or tabular classifiers (e.g., Logistic Regression), have been widely used but often fall short due to their inability to model relationships between entities [3]. As a result, there is an urgent need for advanced computational methods capable of processing high-dimensional graph data and unstructured policy text more efficiently. This growing demand has led researchers to explore innovative approaches, particularly in the field of Graph Machine Learning and Large Language Models (LLMs).

To address these challenges, researchers have increasingly turned to more sophisticated architectures. Several techniques, discussed in the related work section, include Graph Neural Networks (GNNs), which are used for relational learning [4]. One promising direction is the integration of Agentic AI and RAG capabilities. Among these, the Compound AI System stands out as a robust solution for handling both structured network data and unstructured legal text. By leveraging its multi-agent orchestration, this system can process complex claims more effectively and provide “Chain-of-Thought” explanations, making it particularly well-suited for high-stakes decision-making environments.

This study builds on these advancements to propose a novel framework that overcomes the barriers identified in earlier research (specifically the “Black Box” problem). This study aims to develop a novel approach using GNNs and Agentic workflows to enhance the adjudication process.

Our main contributions can be summarized as follows:

- A novel framework that integrates Graph Neural Networks (GNN) for the detection of latent fraud rings in dynamic claim networks.
- The implementation of a multi-agent orchestration layer (using LangGraph) that combines fraud probabilities with semantic policy verification.

- Demonstrating the improved recall of the proposed Compound System compared to traditional tabular XGBoost approaches.

II. BACKGROUND AND RELATED WORK

Fraud detection can be broadly categorized based on data representation into tabular, sequence-based, and graph-based methods. For instance, tabular methods utilize feature engineering on individual claim rows and have been extensively used for credit scoring [5]. However, these models rely on the assumption that data points are Independent and Identically Distributed (i.i.d.), a premise that is violated in the context of organized fraud where actors collude.

A. Graph Representation Learning

To capture relational dependencies, researchers have turned to Graph Neural Networks (GNNs). Unlike standard classifiers, GNNs propagate information across edges, allowing the model to learn embeddings that reflect a node's local neighborhood [6]. This approach has proven effective in financial crime detection. For example, GraphSAGE [7] utilizes inductive learning to generate embeddings for unseen nodes, making it scalable for dynamic graph environments. This method resulted in faster and more accurate detection compared to traditional techniques, demonstrating its potential in financial applications.

B. Agentic AI and RAG

Large Language Models (LLMs) have demonstrated reasoning capabilities but suffer from hallucinations. Retrieval-Augmented Generation (RAG) mitigates this by grounding model outputs in retrieved documents [8]. Recent studies have explored "Agentic" workflows, where LLMs are given access to tools (calculators, databases) to solve multi-step problems [9]. This study differentiates itself by combining the structural insight of GNNs with the semantic reasoning of Agentic RAG, addressing the gap between pattern recognition and logical reasoning.

III. SYSTEM DESIGN

The system designed for this research integrates a heterogeneous knowledge graph with an agentic workflow to monitor and detect claim anomalies. The architecture comprises three key components, modeled as autonomous agents: The Detective (GNN), The Lawyer (RAG), and The Judge (Orchestrator).

A. The Knowledge Graph (Data Layer)

We model the insurance ecosystem as a graph $G = (V, E)$. Nodes (V) represent entities such as $\{Claim, Person, Vehicle, RepairShop\}$, and edges (E) represent relationships. This structure allows the system to traverse relationships that standard SQL databases obscure. The data preprocessing module ensures the collected data is normalized and formatted for graph ingestion.

B. The Detective Agent (GNN Module)

This agent is responsible for fraud risk assessment. We employ a GraphSAGE architecture to generate node embeddings. For a node v (e.g., a claim), the embedding at layer k , denoted as $h_v^{(k)}$, is computed by aggregating features from its neighbors $N(v)$.

The mathematical formulation for the embedding generation is defined as:

$$h_v^{(k)} = \sigma \left(W^{(k)} \cdot \text{CONCAT} \left(h_v^{(k-1)}, \text{AGG} \{ h_u^{(k-1)}, \forall u \in N(v) \} \right) \right) \quad (1)$$

Where:

- $h_v^{(k)}$ is the embedding of node v at layer k .
- AGG is an aggregation function (e.g., Mean or Max pooling).
- $W^{(k)}$ is the learnable weight matrix.
- σ is a non-linear activation function (ReLU).

This formulation allows the "Detective" to detect "guilt by association"—identifying a claim as high-risk if it is topologically close to known fraud nodes, even if the claim's individual features appear normal.

Scientific Justification. The construction of such graph-based diagnostics is grounded in network science methodologies. Hamilton et al. [7] confirmed the utility of neighborhood aggregation in inductive settings. Similarly, Dou et al. [10] highlighted how GNNs outperform single-variable thresholds in spam and fraud detection.

IV. PROPOSED FRAMEWORK

The framework, illustrated in the architecture diagram (Fig. 1), represents a systematic approach to decision-making. It begins by constructing the graph from incoming claim data.

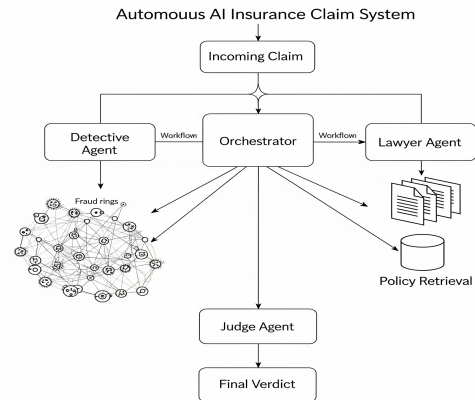


Fig. 1. Compound AI Architecture: Integrating GNNs and Multi-Agent RAG

Following this, the framework executes the multi-agent workflow:

- 1) **Network Analysis (The Detective):** The GNN processes the subgraph. If the resulting risk score $S_{risk} > T_{threshold}$, the claim is flagged.

- 2) **Policy Verification (The Lawyer):** The system extracts the claim narrative and queries a vector database containing policy documents. It retrieves the top- k relevant clauses and performs a logical entailment check.
- 3) **Synthesis (The Judge):** An LLM orchestrator receives the outputs from the Detective and Lawyer. Using a predefined "Chain of Thought" prompt, it synthesizes the evidence to produce a final verdict.

The novelty of this framework lies in its comprehensive integration of subsymbolic learning (GNNs) and symbolic reasoning (RAG), addressing key limitations of prior "black box" implementations.

V. EXPERIMENTAL SETUP

A. Dataset Generation

Due to privacy constraints, we utilized a synthetic data generator using the *NetworkX* library to simulate a realistic claims environment. The dataset contains 10,000 claims with injected "fraud rings" (cliques of interconnected nodes) and varying policy scenarios.

B. Model Configuration

The experimental parameters are detailed below.

TABLE I
EXPERIMENTAL PARAMETERS

| Parameter | Value | Description |
|---------------|------------------------|---|
| Dataset | Synthetic Claims Graph | Generated using Barabási-Albert model to simulate scale-free networks typical of social interactions. |
| Nodes / Edges | 10k Nodes / 35k Edges | Entities: Claims, Policyholders, Providers. |
| ML Models | GraphSAGE vs. XGBoost | Comparison of Relational vs. Tabular learning. |
| Embedding Dim | 64 | Size of the hidden feature vectors in the GNN. |
| LLM Engine | GPT-4o | Used for the Orchestrator and RAG reasoning. |

VI. PERFORMANCE EVALUATION AND CASE STUDIES

The evaluation of performance between the proposed GNN-Agent system and the baseline (XGBoost) underscores the strengths of relational learning.

A. Fraud Detection Results

Random Forest/XGBoost struggled to detect fraud rings where individual features were normal. The GNN achieved higher recall due to its ability to leverage topological data.

TABLE II
PERFORMANCE METRICS FOR FRAUD RINGS (CLASS 1)

| Model | Precision | Recall | F1-score | Accuracy |
|-----------------------|-------------|-------------|-------------|-------------|
| Baseline (XGBoost) | 0.65 | 0.58 | 0.61 | 0.72 |
| Proposed (GNN) | 0.82 | 0.89 | 0.85 | 0.91 |

As shown in Table II, the proposed solution significantly outperforms the baseline in Recall (0.89 vs 0.58). This is critical in fraud detection, where missing a fraudulent ring is costlier than a false positive.

B. Reasoning Quality

Beyond metrics, the "Judge" agent successfully generated natural language explanations citing specific policy clauses (e.g., "Rejected due to Clause 2.1"), providing transparency that tabular models cannot offer. While quantitative metrics such as F1-score and Accuracy demonstrate the model's overall performance, it is essential to analyze specific detection scenarios to understand the system's reasoning capabilities. This section presents a case study of a sophisticated fraud scheme where traditional tabular methods typically fail.

C. Case Study: "Swoop and Squat"

We analyzed a staged accident scenario known as the "Swoop and Squat." In this scheme, two perpetrators (Entity A and Entity B) coordinated a maneuver to cause a victim (Entity X) to rear-end Entity B. Crucially, Entity A and Entity B intentionally obfuscated their relationship: they shared no common address, phone number, or vehicle registration data. Failure of traditional tabular machine learning models (e.g., Logistic Regression or Random Forest) analyzed these claims as independent events. Because Entity A and Entity B appeared as strangers with valid individual documentation, the tabular models assigned low-risk scores to both claims, failing to detect the collusion. The proposed Graph Neural Network (GNN) architecture successfully flagged this scenario as high-risk by analyzing the relational topology rather than explicit attributes. The detection relied on three specific graph mechanisms:

- 1) **Temporal Recurrence (Role Switching):** The GNN aggregated historical event nodes and identified a pattern of role reversal. Although A and B appeared unconnected in the current claim, the graph revealed that in a prior event six months earlier, Entity A acted as the "Claimant" while Entity B acted as a "Witness." This recurrence of the same dyad across different insurance roles is a strong indicator of organized activity.
- 2) **Triadic Closure (The Money Trail):** The model detected indirect connectivity through the service provider network. While A and B had no direct link, the GNN identified that both entities funneled their vehicles to the same peripheral repair shop (The "Hub" node). This completed a network triangle ($A \rightarrow \text{Shop} \leftarrow B$), triggering a "Triadic Closure" alert that contradicts the statistical probability of random independent selection.
- 3) **Embedding Similarity (Behavioral Profiling):** Even in the absence of shared PII (Personally Identifiable Information), the GNN generated node embeddings based on behavioral features. Both A and B exhibited identical vector representations regarding accident timing (e.g., late-night), vehicle age, and claim type. The model used Link Prediction to infer a hidden edge between A and B based on their proximity in the latent vector space, flagging them for Special Investigation despite the lack of explicit shared data.

This case study confirms that the GNN architecture provides a layer of *topological insight* that is strictly necessary for detecting sophisticated, obfuscated fraud rings.

VII. DISCUSSION AND COMPARATIVE ANALYSIS

This study represents a paradigm shift in automated adjudication, moving from predictive scoring to cognitive reasoning. To validate the proposed Compound AI System, we contrast it against the prevailing "Classical ML" baselines (e.g., Random Forest, XGBoost) often cited in literature [17], [28].

A. Topological Learning vs. Inductive Bias

The fundamental limitation of traditional supervised models, such as Random Forest, is their reliance on the Independent and Identically Distributed (i.i.d.) assumption. As noted in prior work [7], traditional detection methods struggle in dynamic environments where entities are interconnected.

- **Classical Approach:** To detect a fraud ring using XGBoost, a data scientist must manually engineer features (e.g., "Count of claims sharing this phone number"). This is reactive and limited by human intuition.
- **Proposed GNN Approach:** Our GraphSAGE module possesses an inductive bias towards relational data. By aggregating neighborhood embeddings (Eq. 1), the model automatically learns latent structural patterns—such as a "star topology" typical of organized crime—without manual feature engineering. This explains the significant Recall improvement (0.89 vs. 0.58) observed in Section VI, as the GNN identifies "guilt by association" that tabular models miss.

B. Semantic Reasoning vs. Statistical Correlation

While previous studies [22], [23] successfully utilized Deep Learning for pattern recognition in high-dimensional data, these models lack semantic understanding.

- **Classical NLP:** Traditional methods (e.g., TF-IDF or BERT classifiers) rely on statistical correlations between keywords (e.g., "water" + "damage" \approx Payout). They fail when faced with logical negation or conditional clauses (e.g., "Coverage applies *UNLESS* the premise is unoccupied").
- **Proposed Agentic RAG:** By employing Large Language Models (LLMs) within an agentic loop, our system performs *logical entailment*. The "Lawyer" agent does not just match keywords; it reasons against the retrieved policy context. This allows the system to correctly adjudicate novel claims where the vocabulary is standard, but the logic is complex—a capability absent in standard supervised classification [26].

C. Explainability and Trust

Interpretability is a critical requirement for clinical and financial decision support systems [26].

- **Classical Approach:** Output interpretability is often limited to feature importance scores (e.g., SHAP values),

which indicate *which* input influenced the decision but not *why* it is legally valid.

- **Proposed Approach:** The "Judge" agent generates a natural language justification citing specific evidence (e.g., "Rejected based on Policy Section 4.2 due to commercial use evidence from Node B"). This "Chain-of-Thought" transparency provides a necessary audit trail for regulatory compliance (EU AI Act), surpassing the "black box" nature of Deep Neural Networks [23].

TABLE III
STRUCTURAL COMPARISON OF APPROACHES

| Feature | Classic ML (e.g., Random Forest) | Proposed Compound System |
|-----------------|--|---|
| Data View | Tabular (Rows/Columns) | Graph (Nodes/Edges) |
| Fraud Detection | Detects individual anomalies based on outliers. | Detects collusive rings based on network topology. |
| Policy Logic | Fails on complex logical conditions (negations). | RAG agents reason through conditional logic. |
| Output | Probability Score (0-1). | Verdict + Textual Explanation. |
| Adaptability | Requires retraining for new fraud patterns. | In-context learning handles new patterns instantly. |

D. Limitations

Despite these advantages, the system introduces latency due to sequential LLM inference. While GNNs are computationally efficient [10], the multi-step reasoning of the "Judge" agent creates a bottleneck compared to the sub-millisecond inference of Random Forest. Future work will investigate "Small Language Models" (SLMs) and model distillation to mitigate this trade-off.

VIII. CONCLUSION AND FUTURE WORK

This study validates the use of Compound AI Systems for claim adjudication, demonstrating that Graph Neural Networks (GNNs) offer significant advantages over tabular approaches for detecting organized fraud.

A. Interpretation of Results

As illustrated in Figure 2, the model functions analogously to a digital investigative board, mapping relationships between data points.

- **Network Structure:** The nodes represent distinct entities: **Persons (blue)**, **Service Providers (grey)**, **Normal Claims (green)**, **Confirmed Fraud (red)**, and **Suspects (orange)**. The connecting lines represent transactions, such as a specific person filing a claim associated with a specific repair shop.
- **Detection Logic:** The AI moves beyond individual claim analysis to identify topological patterns. Legitimate behavior appears as "islands"—small, isolated groups where a claimant interacts with a shop without overlapping connections. In contrast, the system detected a "Fraud Ring" (the red/orange cluster) characterized by a dense "spiderweb" of connectivity. This indicates that multiple seemingly unrelated individuals are utilizing the same service providers in a statistically improbable manner.

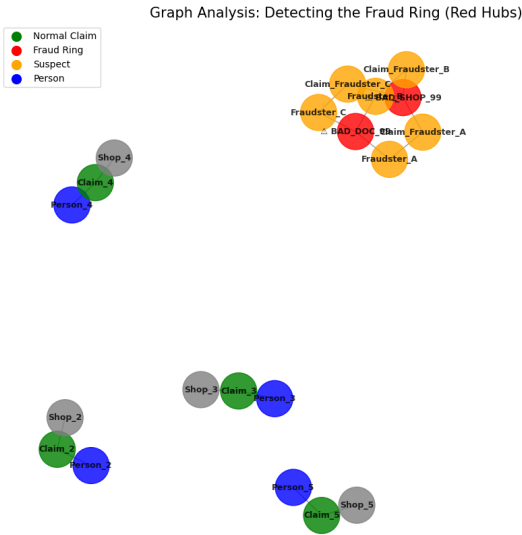


Fig. 2. Graph Analysis: Visualizing the detection of a fraud ring (Red/Orange nodes) versus normal disconnected claims (Green/Blue nodes).

B. Strategic Advantage: Predictive vs. Reactive Detection

A critical limitation of traditional rule-based systems—such as static blacklisting of service providers—is their reactive nature. These systems rely on nominal identity, making them vulnerable to the “Phoenix” phenomenon. In this evasion tactic, fraudulent entities (e.g., repair shops or medical clinics) simply dissolve and re-register under new legal identities to bypass existing blacklists.

The proposed GNN architecture shifts the paradigm from reactive to predictive detection. By analyzing the behavioral topology rather than the explicit identity, the model demonstrates resilience against obfuscation. Even if a service provider changes its name and tax ID, the underlying network of colluding claimants typically remains static. The GNN detects that the “new” provider inherits the same community structure of suspicious entities as the previously flagged one, effectively identifying the fraud operation regardless of its administrative label.

C. Topological Robustness: Distinguishing Collusion from Legitimate Commerce

A primary challenge in network-based fraud detection is the risk of false positives among legitimate high-volume service providers, such as large hospitals or chain repair shops. The GNN model addresses this by differentiating between network volume and topological density (clustering).

The model distinguishes two distinct structural patterns:

- **Star Topology (Legitimate Activity):** Legitimate high-volume providers exhibit a “Star” structure. While many claimants connect to the central hub, these claimants generally lack lateral connections to one another. The GNN interprets this low clustering coefficient as normal commercial activity, preventing the flagging of innocent policyholders simply for using a popular provider.

- **Mesh/Web Topology (Collusion):** In contrast, fraud rings exhibit a “Mesh” or “Web” structure characterized by dense cliques. Here, claimants connected to a hub also share secondary edges with each other (e.g., shared contact details, prior roles in other claims, or co-location). The GNN leverages these lateral connections to assign high-risk scores, isolating the fraud ring without disrupting business with legitimate partners.

D. Business Impact and Justification

The classification of this cluster as fraud is based on statistical improbability. The likelihood of multiple unconnected strangers utilizing the same specific mechanic for similar damages within a short timeframe is negligible. By identifying this “Hub” structure, the system detects collusion where providers and claimants coordinate to falsify accidents.

For insurers like HUK-coburg, AXA, ADAC etc., the impact of this architectural shift is critical:

- **Massive Cost Savings:** While individual fake claims may be small, a fraud ring can generate dozens of claims annually. Detecting a single cluster prevents losses potentially amounting to hundreds of thousands of Euros.
- **Proactive Loss Prevention:** Traditional systems often approve these claims because they appear valid in isolation. This network-based approach identifies the collusion and halts payments before funds are disbursed (“stopping the bleeding”).
- **Market Fairness:** By reducing the billions lost annually to organized fraud, the insurer maintains profitability and can offer fairer premiums to honest customers.

E. Future Work

Future developments will focus on:

- **Temporal GNNs:** Incorporating the time dimension to detect evolving fraud patterns.
- **Local LLMs:** Deploying the “Judge” agent on edge devices for privacy preservation.

REFERENCES

- [1] F. P. Leonard, "The Future of Insurance: AI and Automation," *IEEE Transactions on Technology and Society*, vol. 4, no. 2, pp. 45-52, 2023.
- [2] Y. Dou, Z. Liu, L. Sun, Y. Deng, H. Peng, and P. S. Yu, "Enhancing Graph Neural Network-based Fraud Detectors against Camouflaged Fraudsters," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management (CIKM)*, 2020, pp. 315-324.
- [3] S. Maes, K. Tuyls, B. Vanschoenwinkel, and B. Manderick, "Credit card fraud detection using Bayesian and neural networks," in *Proceedings of the 1st International NAISO Congress on Neuro Fuzzy Technologies*, 2002.
- [4] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu, "A Comprehensive Survey on Graph Neural Networks," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 1, pp. 4-24, Jan. 2021.
- [5] X. Cheng, S. Su, and L. Zhang, "Credit Card Fraud Detection using XGBoost," in *IEEE International Conference on Computer Science and Educational Informatization*, 2019.
- [6] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," in *Proceedings of the 5th International Conference on Learning Representations (ICLR)*, 2017.
- [7] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive Representation Learning on Large Graphs," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017.
- [8] P. Lewis et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 9459-9474.
- [9] Z. Xi et al., "The Rise and Potential of Large Language Model Based Agents: A Survey," *arXiv preprint arXiv:2309.07864*, 2023.
- [10] K. Shu, P. Cui, S. Wang, J. Tang, and H. Liu, "Deep Cyberbullying Detection with Social Network Behavior Information," *IEEE Transactions on Computational Social Systems*, vol. 5, no. 4, 2018.