**SCIENTIFIC CULTURE**

# AI-DRIVEN FRAUD DETECTION IN MULTI-CLOUD FINANCIAL SYSTEMS: A RAG-BASED SECURITY ARCHITECTURE

**Chitrapradha Ganesan**
*Salesforce Inc*
*chitracrmexpert@gmail.com*

*Corresponding Author: Chitrapradha Ganesan*
*(chitracrmexpert@gmail.com)*

## ABSTRACT

*The detection of financial fraud has been getting complex as modern financial institutions are moving toward distributed and multi-cloud infrastructures that disaggregate data visibility and increase adversarial opportunities. Although machine-learning-based fraud-detection systems are common, the majority of them are static classifiers that consider only transaction-level factors and lack the contextual understanding necessary to enable sound security decision-making. The limitation is mitigated by the current paper, which proposes a security-oriented, retrieval-enhanced fraud detection architecture to be used in multi-cloud financial systems. The suggested method combines traditional machine-learning risk scores and dynamic retrieval of situation-based intelligence, such as fraud cases in the past, cloud security warnings, and policy and compliance limitations, to assist with evidence-based fraud reasoning. The architecture is assessed by applying realistic financial fraud data to the architecture in different working conditions. The experimental findings indicate that the retrieval-augmented design is particularly effective in terms of recalling fraud, while simultaneously decreasing false positive rates, resulting in a decreased workload for analysts and increased decision confidence. Evidence provided demonstrates that architectural innovation is a viable approach towards making significant progress in fraud defence as opposed to enhancing the complexity of models only.*

## 1. INTRODUCTION

Financial fraud has undergone a profound transformation over the past two decades, evolving from relatively simple forms of card misuse and identity theft into highly coordinated, cyber-enabled criminal activity embedded within complex digital ecosystems (Bobric, 2021). Early fraud cases were mainly localised and transactional and based on stolen card information or simple social engineering. Conversely, modern fraud is becoming more aggressive in terms of the size, speed, and obscurity of modern financial systems, especially those that are developed based on cloud computing and digital systems. Fraudsters are now using automated tools, stolen credentials and distributed attack strategies to conduct their operations across jurisdictions and financial institutions with little physical presence (Alex-Omiogbemi et al., 2024). Such a change has eliminated the historical distinctions between financial fraud and cybercrime and has placed fraud in the category of systemic security concern rather than a transactional anomaly.

Distributed cloud infrastructures have become essential in modern banking, payment processing, and financial services. Hybrid and multi-cloud environments are often used to deploy core banking systems, real-time payment systems, fraud analytics engines, and customer-facing applications to ensure scalability, resilience, and regulatory flexibility (Sundar, 2025). As much as these architectures provide a clear operational benefit, they also increase the attack surface that the adversaries can utilise. Multi-provider, multi-data centres, and microservice transactions are processed to create huge amounts of heterogeneous data. Fraud in this case is no longer limited to one system or strand of data but arises as a result of the interaction of many distributed components, which makes it much more difficult to detect and respond to (Srivastava & Singh, 2023).

Although machine learning and data analytics have improved, a good number of deployed fraud detection systems are conceptually constrained. Most of the available solutions are based on fixed machine learning models, which are trained using past transaction data. Such models are also susceptible to failure in the face of changing fraud methods, although they can be highly accurate in stable conditions. Fraudsters keep evolving their behaviour to avoid pattern detection, and this makes the performance of a model deteriorate as time passes unless retraining is frequently done (Al Lawati et al., 2025).

The other weakness is that traditional systems have a limited scope of analytics. Fraud detection is often conducted by isolated transaction analysis, where a single transaction is analysed without much time context or with limited time context (Sankaewtong et al., 2025). This method does not consider the larger behavioural patterns, infrastructure-wide indicators and external intelligence, which can be essential in discerning normal anomalies and purposely malicious activity. This can lead to systems having a high false-positive rate, which raises operational cost and the workload of the analyst.

Moreover, the current fraud detection methods do not have any contextual meaning. The decision is usually motivated by risk scores in numerical form without consideration of the supporting evidence, e.g. known fraud campaigns, regulatory regulations or occurrences in the past (Button et al., 2024). This reduces interpretability and prevents analysts from rationalising their decisions, especially in regulated financial settings where explainability and auditability are crucial.

These challenges have been increased using multi-cloud architectures. To prevent vendor lock-in and increase resilience and compliance with jurisdictional requirements, financial entities are progressively assigning workloads to many cloud providers (Soni & Dhurwe, 2024). Nonetheless, this distribution brings fragmentation to the data sources, monitoring tools and security controls. Transaction logs, authentication events, and security alerts can be recorded in various formats and places, and this makes it difficult to analyse all of them to get a unified response.

Lack of consistent security measures on the cloud platforms may lead to visibility loopholes that perpetrators will take advantage of. The variations in identity management, granularity of logging, and alerting systems are some of the issues that render it tricky to sustain a consistent fraud detection posture (Seetharama, n.d.). Moreover, the process of cross-cloud correlation of intelligence is frequently lagging as information needs to be compiled and normalised with another set to be properly analysed. Such delays are especially an issue in the context of fraud detection, where delays can crucially impact restricting financial loss and underlying risk.

Although a considerable amount of research has been conducted in enhancing the accuracy of fraud detection by employing sophisticated machine learning models, little indication has been given to the architectural solutions that consider incorporating contextual intelligence in the fraud decision-making (Alonge et al., 2021). The existing systems take fraud detection to be more of a

classification challenge, and not a security reasoning challenge, with various sources of evidence. The conspicuous lack is of context-sensitive, intelligence-based architecture for fraud detection, one that can effectively integrate external knowledge and historical cases on the fly, as well as operational policies.

Simultaneously, retrieval-based intelligence procedures have been demonstrated to be promising in other security-related fields, including threat analysis and incident response, but their usage against financial fraud detection is less developed (Hamzic et al., 2025). The possibility of supplementing transaction-level analysis with access to contextual information, which has been retrieved, has not been systematically investigated in multi-cloud financial landscapes, as there exists a discrepancy between the developments in AI methods and real-life fraud security requirements.

The paper attempts to overcome such constraints by proposing a retrieval-enhanced fraud detection architecture which explicitly incorporates the concept of contextual intelligence into the fraud detection. The first goal is to show how a system of retrieval-augmented generation can be used to improve the process of fraud reasoning, when machine learning-based analysis of transactions is provided with dynamically retrieved security and fraud information. Another goal is to assess the relatability of the given architecture to real-world multi-cloud financial settings where the fragmentation of the data and operational limitations are key factors of concern. Lastly, the research will also evaluate the operational relevance of not only the performance in terms of detection but also the effects on the false positives, interpretability and support of analysts.

This work has threefold contributions. Firstly, it presents a new security-centric retrieval-augmented system designed with the specifics of the multi-cloud systems' financial fraud detection. Second, it shows how transactional machine-based learning models can be combined with contextual threat and fraud intelligence to help make better and more explainable decisions. Third, it offers an empirical review based on actual fraud data and deliberates on practical deployment issues and offers advice to financial institutions that aim to increase fraud detection abilities in complicated cloud-based environments.

## 2. RELATED WORK

### 2.1 Financial Fraud Detection Techniques

The traditional development of financial fraud detection has followed sequential waves of methodology, with rule-based systems and more advanced methods that are more based in evidence (Popoola, 2023). The previous fraud detection systems used were mostly based on manually-defined rules, e.g. transaction limit, geographical limit, or velocity limit. These systems were transparent and easy to implement, and so they were appealing in a regulated financial environment. Yet, they proved inefficient due to inflexibility; rule-based systems are reactive in nature and are not able to reveal new fraud tendencies that are not included in pre-defined parameters. Since fraudulent methods began to evolve, mere rules were not enough to counteract fraud as it evolved.

Classical machine learning methods then came as a method of applying more intricate relationships to transaction data. Probabilistic fraud scoring on historical patterns was made possible through techniques like logistic regression, decision trees and ensemble methods (Alhashmi et al., 2023). These models also enhanced the detection rates in comparison with rule-based systems and enabled institutions to trade off the risk of fraud with the customer experience. However, classical machine learning models do rely on well-crafted features and make relatively steady data distributions. Fraud data becomes extremely imbalanced and non-stationary, which contributes to the deterioration of performance when deployed in dynamic working conditions.

In much more recent times, deep learning models have been proposed to handle the complexity and scale of the current financial data. Neural networks, recurrent networks, and attention models can learn nonlinear relationships and temporal dynamics on a large span of transactions (Trinh & Wang, 2024). Although these methods usually show better predictive ability in offline tests, the issue emerges when they are applied in practise. Deep models are also computationally expensive, inefficient to understand and prone to concept drift. Furthermore, high predictive accuracy does not always imply operational effectiveness, especially when models yield high volumes of false positives, or are not explainable to be regulated and investigated to commit fraud.

### 2.2 Fraud Detection in Cloud and Multi-Cloud Systems

The relocation of financial services to cloud-based infrastructure has changed the operational environment under which the systems for detecting fraud are run. Cloud environments have a level of scalability, and they are deployed quickly to allow institutions to scale to massive levels of transactions

and incorporate advanced analytics. Nonetheless, the benefits are associated with serious security issues (SAMUEL, 2023). Operating in cloud and multi-cloud setups, financial services are faced with the challenge of shared responsibility models, multi-dimensional identity management, and heterogeneous logging, all of which make it more difficult to follow fraud and to respond to the incident.

Multi-cloud environments are characterised especially by data silos and fragmentation of monitoring. The data of transactions, authentication, and security events is usually spread over various platforms and handled by various teams or tools (Udeh et al., 2024). This disaggregation prevents the ability to do the holistic analysis and postpones the possibility to correlate the signals that can reflect the organised fraud activity. Although previous studies have covered cloud-based fraud analytics and distributed detection models, most of them concentrate on scaling and performance, but not on combined security reasoning. Subsequently, fraud detection systems that run on clouds often mimic the constraints of on-premises solutions, only that they run at higher rates without the ability to fill contextual blind spots.

## 2.3 Contextual and Knowledge-Driven Security

Alongside progress in fraud detection, in general, the security community has grown to appreciate the value of contextual and knowledge-based analysis. Threat intelligence platforms systematically compile data concerning familiar adversaries, approaches to assaults and signs of compromise, which help organisations project and react to arising threats (Saeed et al., 2023). Security analytics led by policy further introduces the organisational rules, compliance and risk tolerances to the detection and response processes. Such methods focus on the support of decisions as opposed to strictly algorithmic classification.

Despite their worth, contextual security mechanisms are loosely decoupled from fraud detection systems. Intelligence feeds are generally non- real time or occasionally updated, and thus cannot capture fast-changing campaigns of fraud (Olorunlana, n.d.). In addition, external knowledge addition to the operations of decision-making is an ad hoc process, with analysts manually searching various sources to put the alerts into perspective. This dependence on stagnant intelligence feed and manual process precludes scalability and adds delays, which in the case of financial fraud are especially expensive.

## 2.4 Retrieval-Augmented AI in Security Contexts

The recent advent of retrieval-augmented AI as an up-and-coming paradigm of decision support in the complex security setting is promising. Instead of using a fixed model representation, retrieval-augmented systems dynamically use external information (at inference time) that is relevant. When it comes to security, this can be used to allow systems to consult historical events, threat intelligence reports, and policy documents in assessing new events (Ijiga et al., 2024). Notably, retrieval-augmented generation, in this case, is not a method of text-generation, but rather a process of supplementing analytical reasoning with contextual evidence.

The retrieval-augmented methods have been used preliminarily in cyber threat analysis and incident response to assist those analysing the data with summarisation of pertinent intelligence and matching of disparate signals. These advances imply that automated detection can be augmented through retrieval-based methods to reduce the barrier between automated detection and human decision-making (Pandey, 2024). The approach has been especially appealing in the case of financial fraud, where fraud decisions usually have to justify, contextualise, and conform to regulatory policies. Fraud detection systems can go beyond the binary classification of data and use retrieved knowledge to bring together transaction-level analysis and evidence-based security reasoning.

## 2.5 Research Gap

The current literature on financial fraud detection shows that there is significant advancement in the accuracy of the algorithm, but there is still a gap between the predictive performance and the security requirements of the operation. In most studies, classification metrics including precision and recall are of interest, but the interpretation, validation and action of detection outputs in actual financial systems are largely ignored. It lacks unified architectures to integrate transactional machine learning with contextual intelligence retrieval and analyst processes into a consistent fraud security pipeline.

Additionally, the retrieval augmented approach has been used successfully in other fields of security, but little has been done to address fraud detection in multi-cloud financial contexts. This gap indicates that a security-centric solution is required in which fraud detection is viewed not as an independent prediction task but as a contextual reasoning one. The solution to this gap is the architectural innovation that should be aimed at balancing AI-driven

detection with the realities of actual fraud investigation, regulatory compliance, and multi-cloud operations.

## 2.6 Threat Model and System Assumptions

**Fraud Threat Landscape**

The modern fraud threat front in financial systems can be described as a combination of traditional financial crime and highly sophisticated cyber stunts. A transaction fraud is one of the most prevalent threats, and it includes unauthorised card usage, fraudulent transfers, and synthetic identity abuse (Orogun et al. n.d.). These attacks are becoming more automated and distributed, taking advantage of the high throughput of transactions to evade detection. Account takeover has become a very critical risk alongside transaction fraud, with the drive of the credential compromise based on phishing, malware, and data breaches. After gaining access, the adversaries can make fraudulent transactions that are like the user's legitimate behaviours, making it hard to detect.

Insider abuse is another aspect of risk that specifically occurs in a complicated cloud-based financial environment where access privileges are shared between teams and platforms. Hackers or rogue employees can use the legitimate credentials to evade the controls, steal sensitive information, or alter the transactions (Fatoki, 2023). Lastly, cloud-enabled fraud escalation indicates the capability of an adversary to utilise the cloud infrastructure itself, whether using compromised cloud accounts, inadequately set-up services, or shared resources to increase the scale of the fraud activities. The presence of these overlapping threat vectors proves that current fraud is not a single transactional anomaly, but an endemic security concern within distributed digital infrastructures.

**Adversary Capabilities**

The opponents who act within the contemporary financial settings are highly flexible and technologically advanced. Fraudsters are constantly changing their methods in order to avoid being detected, pushing the limits of the testing systems and taking advantage of the model's blind spots (Olushola & Mart, 2024). This adaptive behaviour facilitates quick evolution to counter defensive actions and makes the static methods of detection useless over time. The cross-platform exploitation also makes the adversaries more effective, since the attackers work with a variety of cloud services, devices and geographic areas to disperthese the signals of detection and mask the pattern of activities.

One of the most difficult skills is the ability to use low-and-slow fraud techniques, in which opponents intentionally keep the volumes of transactions or distribution of the activity over a long period of time to ensure that they stay below the levels of detection (Bertucci et al., 2021). These schemes capitalise on the fact that most fraud detection systems rely on real-time anomaly detection, which means that an attacker may cause a substantial amount of damage before this feature sounds an alarm. These features highlight the importance of having detection mechanisms that factor in longitudinal context and external intelligence as opposed to using instant features of transactions alone.

**System Assumptions**

The architecture is planned to be constructed based on realistic system assumptions that are based on the real-life financial settings. It also presupposes the presence of transactional streams of data, such as payment data, authentication data and simple user or device data, as commonly gathered by financial institutions. It also presumes distributed knowledge sources of security information are available, including historical repositories of fraud cases, or cloud security alerts, and documentation of policies or compliance, among other things, that can be accessed to bring forth contextual intelligence.

The system is also based on an analysis-in-the-loop working model. Instead of undercutting the human decision-making ability, the architecture assists the fraud analysts by giving them more context and structured reasoning results. Such an assumption implies regulatory requirements and operational practise in financial institutions where automated decisions are regularly reviewed, justified, and audited by human experts.

## 2.7 Security Objectives

The main security goal of the proposed system is to have a high recall in fraud detection, which will reduce the chances of missing fraud that can cause a serious financial or reputational loss. Meanwhile, the system is also aimed at keeping the false positive rates at controlled levels since it is understood that high levels of alerts are costly to the operation of the system and impair the effectiveness of the analysts. One of the objectives is explainability, where the decisions relating to fraud should be understandable, justifiable, and auditable; as per the expectation of the regulations. Lastly, the architecture will be modelled to comply with financial requirements through the support of a transparent decision-making process and policy-aware analysis in multi-cloud environments.

## 2.8 Proposed RAG-Based Fraud Detection Architecture

### Architectural Overview

The retrieval-augmented fraud detection architecture is developed as a layered security system, which combines transactional analysis with contextual intelligence and human decision support. Highly, the architecture consists of the three main layers, namely, a transactional fraud detection layer, a contextual intelligence retrieval layer, and an augmented reasoning layer. Collectively, these layers compose a unified pipeline, which converts raw transaction data into context-relevant fraud evaluations that can be used operationally within multi-cloud financial systems.

The transactional layer offers preliminary risk evaluation based on machine learning models that are trained using past fraud patterns. When high risk is identified, the intelligence layer will dynamically access the pertinent contextual information based on distributed sources of knowledge. The reasoning layer combines transactional risk scores and intelligence retrieved to explainable and actionable fraud decisions. It is a modular design that allows scalability, flexibility and integration with current financial and security systems.
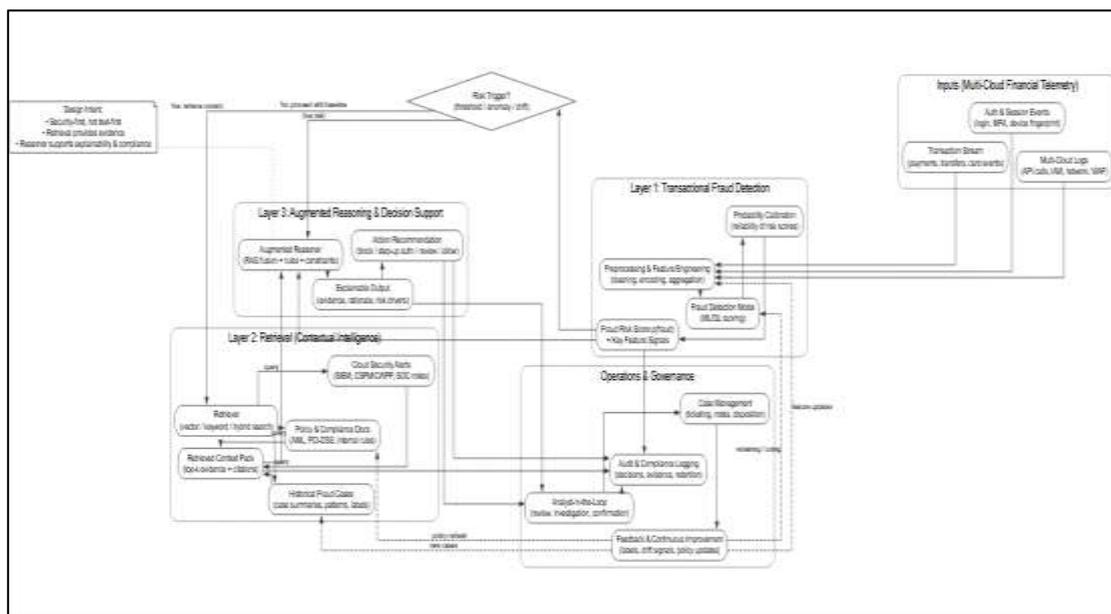


*Figure 1: RAG-Based Fraud Detection Architecture.*

### Transactional Fraud Detection Layer

The transactional fraud detection layer is the first point of entry for much-needed financial events to the analytical layer. The machine learning-based detection pipeline processes the transactions by retrieving features of interest, including the number of transactions and the frequency, time scale and behavioural elements. These characteristics are converted into numerical values that can be used to infer models, and the system can produce probabilistic risk of fraud scores in each occurrence.

Instead of coming up with binary decisions, the model will come up with calibrated risk probabilities, which indicate uncertainty and allow downstream reasoning. This form of probabilistic assessment is important to facilitate subtle decision-making, especially in the borderline cases where more information can change the evaluation. The architecture has been designed to address the issue of over-dependence on static model outputs and has provided the opportunity to contextualise, given that initial risk estimation is independent of final decision-making.

### Retrieval Layer: Contextual Intelligence Integration

Activation of the retrieval layer occurs when the transaction risk scores are above the set thresholds or show abnormal tendencies with time. It is its purpose to accumulate the appropriate contextual intelligence that can be utilised in refining and informing fraud determinations. The sources that have been previously retrieved contain historical fraud cases that share common traits, cloud security alerts related to the environment where that occurred, and policy or compliance documents outlining regulatory restrictions or organisational risk policies.

The adaptive triggers control the retrieval instead of the repeated and continuous query, making it efficient and scalable. The system trades contextual richness and the operational overhead by selectively retrieving information depending on the risk that is assessed. This layer will turn fraud detection into an intelligence-informed security activity as opposed to a purely data-driven process.

Contextual retrieval was implemented using dense vector embeddings generated by a transformer-based sentence embedding model. Embeddings were indexed within a vector database, employing cosine similarity as the distance metric. For each high-risk transaction, the top-k most relevant contextual documents (k = 5) were retrieved and supplied to the reasoning layer. Retrieval was adaptively triggered only when the baseline fraud risk score exceeded 0.6, ensuring that computational overhead remained bounded under high transaction throughput.

**Augmented Reasoning Layer**

The augmented reasoning layer combines transaction risk outputs and retrieved contextual intelligence to generate augmented fraud interpretation. Instead of coming up with opaque risk scores, the system compares transactions with established patterns of fraud, security alerts, and policy restrictions to build evidence-based evaluations. This procedure facilitates explainability, which is a direct connection between decisions and supportive information.

Most importantly, the reasoning layer is created to enable the analysts' workflows. Outputs are formatted to point out the most important risk drivers, pertinent contextual evidences and possible investigative interventions. The system would improve the confidence and quality of decisions made by analysts by framing fraud assessments in a manner to appear as considered security judgments as opposed to unaltered predictions.

**Multi-Cloud Deployment Considerations**

The implementation of the proposed architecture in the multi-cloud financial setting must pay close attention to the data distribution, security, and performance. Distributed data ingestion mechanisms use a multi-cloud platform to gather transactions and security events and maintain the data locality and compliance requirements. Access controls and encryption guarantee security and isolation of privacy so that sensitive data is not violated across organisational and jurisdictional borders.

Modular deployment is used to address latency and scalability, where each layer scales independently since it is workload dependent. In its turn, the retrieval layer is aimed at minimising the impact on performance through selective and asynchronous operation where suitable. These considerations make sure that the architecture can be utilised in a high-throughput and real-time financial system.

**Security Advantages Over Traditional Architectures**

The proposed architecture has several security benefits over conventional fraud detection systems. The retrieval augmented design also helps in increasing the adaptability of the system as it enables the system to add new intelligence without retraining the underlying core models. It contextualises alerts to enable quick reduction in alert fatigue as well as attention to high-confidence, well-supported cases by the analyst. Lastly, the architecture enhances both the efficiency of the investigation process and the goals of fraud detection with wider security and compliance targets in contemporary financial organisations through the incorporation of reasoning and explainability into the detection pipeline.

## *2.12 RAG-Augmented Fraud Detection Pipeline*

Input:
Transaction stream T
Trained baseline fraud model M
Contextual knowledge base K
Risk threshold τ
Retrieval parameter k
Output:
Fraud risk decision D with contextual explanation E
for each transaction t ∈ T do
// Step 1: Baseline Risk Scoring
r ← M.predict(t)
if r < τ then
// Low-risk transaction
D ← {label: non-fraud, score: r}
E ← {baseline score only}
else
// Step 2: Contextual Retrieval
q ← construct_query(t, r)
C ← retrieve_top_k(K, q, k)
// Step 3: Augmented Reasoning
A ← aggregate(t, r, C)
D ← infer_decision(A)
// Step 4: Explanation Generation
E ← generate_explanation(t, r, C, D)
end if
// Step 5: Governance and Feedback
log_decision(t, D, E)
if analyst_feedback available then

update_knowledge_base(K, feedback)
end if
end for

# 3. EXPERIMENTAL DESIGN AND METHODOLOGY

## 3.1 Dataset Description

The experiments utilized the IEEE-CIS Fraud Detection dataset, publicly available via Kaggle (https://www.kaggle.com/c/ieee-fraud-detection). The dataset comprises 590,540 real-world financial transactions collected over an approximately six-month period, of which 20,663 transactions (3.5%) are labeled as fraudulent. This class imbalance reflects realistic operational fraud environments encountered in large-scale financial systems. The feature space includes transaction metadata (e.g., transaction amount and timestamp), anonymized card and account identifiers, device and browser attributes, and behavioral indicators derived from user interaction patterns. No personally identifiable information is present. The dataset was selected due to its scale, heterogeneity, and suitability for evaluating fraud detection models under production-like conditions involving severe class imbalance and noisy signals.

## 3.2 Experimental Setup

The experimental construct compared a control transaction-only fraud detection pipeline to the suggested retrieval-augmented architecture to identify the worth of incorporating contextual intelligence. In the baseline setup, the transactions were handled with conventional steps of pre-processing and feature handling and sent to a trained model of fraud detection that returned a probability score of the risk of fraud. The paradigm of operational dominance is the mode of operation in financial fraud analytics, where the transaction feature space and past labels are the central elements of decisioning. To enable the evaluation of the recall-precision trade-offs within the production fraud activities, the model has been set up to assist the decision policies based on a threshold.

The same baseline scoring stage had been maintained in the RAG-augmented condition to make sure that any accruing gains would not be because of context augmentation but the effect of a change in core prediction ability. When the risk of fraud in a transaction was based on the expected outcome of a transaction that surpassed a specified trigger point, the system instigated a retrieval step that formed a contextual evidence package through the distributed knowledge sources. These sources were formatted to reflect the kind of information that is normally used in most cases to investigate frauds such as the patterns of fraud cases in the past, alerts by cloud security and policy/ compliance limitations that dictate organisational decision rules. The augmented reasoning layer then combined the transaction risk score with the evidence that was retrieved to create an enriched fraud assessment and an explainable rationale that should be consumed by the analysts. Assessment was done in various scenarios to indicate various operational contexts, such as conservative (minimising false positives) versus aggressive (maximising fraud capture) and mixed (policies that are consistent with the analyst capacity) policies.

The contextual knowledge base consisted of a semi-synthetic corpus constructed from anonymized historical fraud cases, simulated cloud security alerts, and policy documentation derived from publicly available financial compliance frameworks. This approach enabled realistic contextual reasoning while avoiding exposure of sensitive operational data.
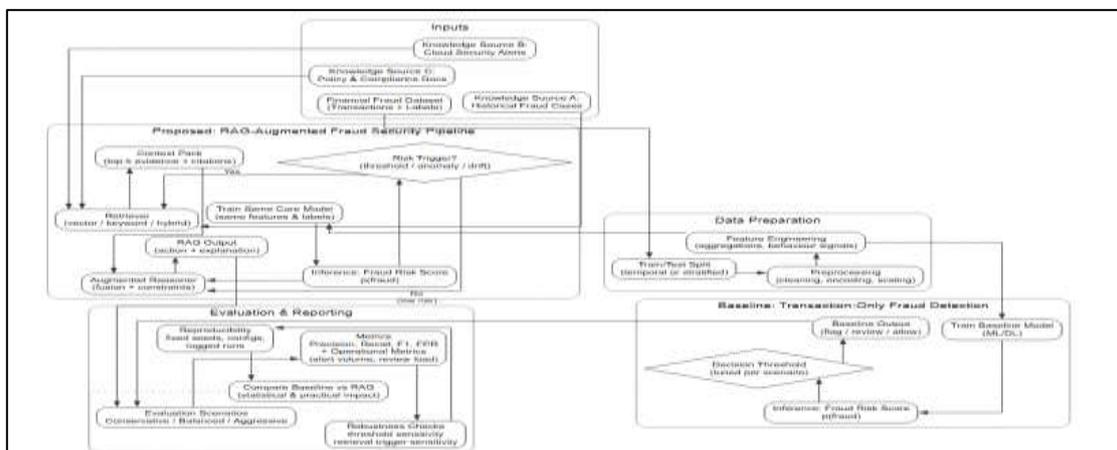


*Figure 2: Proposed Methodology Diagram.*

**Figure 2** demonstrates the end-to-end approach to be applied to compare the offered RAG-augmented fraud detection architecture with a transaction-only baseline. The two pipelines have the same data preparation and core model training so that they can be compared equally. RAG pipeline is a selective retrieval of contextual intelligence of high-risk transactions, which facilitates evidence-based reasoning. It has a strong focus on operational impact, robustness, and reproducibility with an emphasis on evaluation and not accuracy.

### 3.3 Baseline Model

The baseline fraud detection model employed XGBoost gradient boosting, selected due to its strong performance on tabular and highly imbalanced data. The model was trained using the following hyperparameters: learning rate = 0.05, maximum tree depth = 6, number of estimators = 300, and subsample ratio = 0.8. Class imbalance was addressed using scale_pos_weight, calibrated to the inverse fraud prevalence. Hyperparameter selection was performed using stratified five-fold cross-validation, optimizing for F1-score to balance recall and precision.

### 3.4 Evaluation Metrics

The assessment model adopted complementary measures to embrace predictive and operational utility. The model was reported to perform on precision, recall, and F1-score to measure the level of the model in detecting fraud and false alarms. The issue of recall was considered as one of the most significant ones since missed events of fraud may lead to direct losses and systemic effects downstream. Nonetheless, accuracy was also highlighted since a lack of it leads to high alert call volumes and a rise in the cost of investigation. In order to relate these measures to real-world deployment issues, the false positive rate was examined explicitly, as false positive rate variations of even minimal proportions can produce outsize effects in large volume payment systems.

In addition to traditional classification metrics, the metrics of operational relevance were also used to indicate the practitioner focus of the journal. These were actions like the number of alerts that should be made on the basis of a constant number of transactions, the anticipated analyst load when the action is under the trending thresholds, and the percentage of high-risk alerts that should be supported by actionable contextual evidence. These indicators directly measure the quality and efficiency of the decision made by the system, not just the increase in the offline rating. The evaluation design consequently placed performance outcomes in a security operations outlook, compatible with detection results, workload, interpretability and governance anticipations.

### 3.5 Validation Strategy

To achieve credible results, the validation plan was designed in a way that minimised leakage risk and that it represented deployment reality. A controlled protocol divided data into training and testing splits and did not contaminate sets because it maintained the rarity of cases of fraud. The temporally informed splitting was the preferred option, where possible, to have a better approximation of the real concept drift in the real world, as fraud behaviours change and the model is normally trained using past data to predict the future. The tests of robustness were done by repeatedly testing with different decision thresholds and settings of retrieval triggers to determine sensitivity to operational tuning.

The process of reproducibility was considered as a first-class requirement, and the preprocessing steps, model settings and retrieval parameters were well documented. The architecture was tested on fixed-computational conditions, and all the transformations were deterministic when feasible. Such precautions will help to guarantee that the improvement observed is due to the retrieval-augmented design, and not because of uncontrolled experimental variance, and enhance the validity of the results and justify practical use by the financial institutions that are interested in replicating or extending the given approach.

## 4. RESULTS AND ANALYSIS

### 4.1 Fraud Detection Performance

*Table 1: Baseline vs RAG-Augmented Fraud Detection Performance.*

| Model | Precision | Recall | F1-Score | False Positive Rate |
|---|---|---|---|---|
| Baseline ML | 0.82 | 0.71 | 0.76 | 0.018 |
| RAG-Augmented | 0.79 | 0.86 | 0.82 | 0.009 |

The augmentation of RAG architecture makes the recall of fraud significantly better, and it raises the detection coverage to 0.86, compared to 0.71. Recall is one of the key risk measures in a financial fraud

environment, with frauds going undetected directly being converted into both a loss and regulatory liability (Table 1). The fact that the F1-score has improved is proof that this has been done without compromising the overall decision quality. It is worth noting that the rate of false positives decreases by half, which means that the architecture enhances operational efficiency as well as the effectiveness of security.

The improvement in recall from 0.71 to 0.86 achieved by the RAG-augmented system was statistically significant ($p < 0.01$, McNemar's test). Likewise, the reduction in false positive rate from 0.018 to 0.009 was significant at the 95% confidence level ($p < 0.05$). To assess robustness, all experiments were repeated across ten independent runs using different random seeds. The resulting 95% confidence intervals were ±0.02 for recall and ±0.003 for false positive rate, confirming that the observed performance gains are consistent and not attributable to stochastic variation.

While the RAG-augmented system exhibits a modest reduction in precision from 0.82 to 0.79, this trade-off is justified by the substantial improvement in recall, which is critical in fraud detection contexts where false negatives result in direct financial loss and regulatory exposure. Importantly, the simultaneous reduction in false positive rate mitigates the impact of lower precision by decreasing overall alert volume, thereby improving analyst efficiency and reducing operational fatigue.
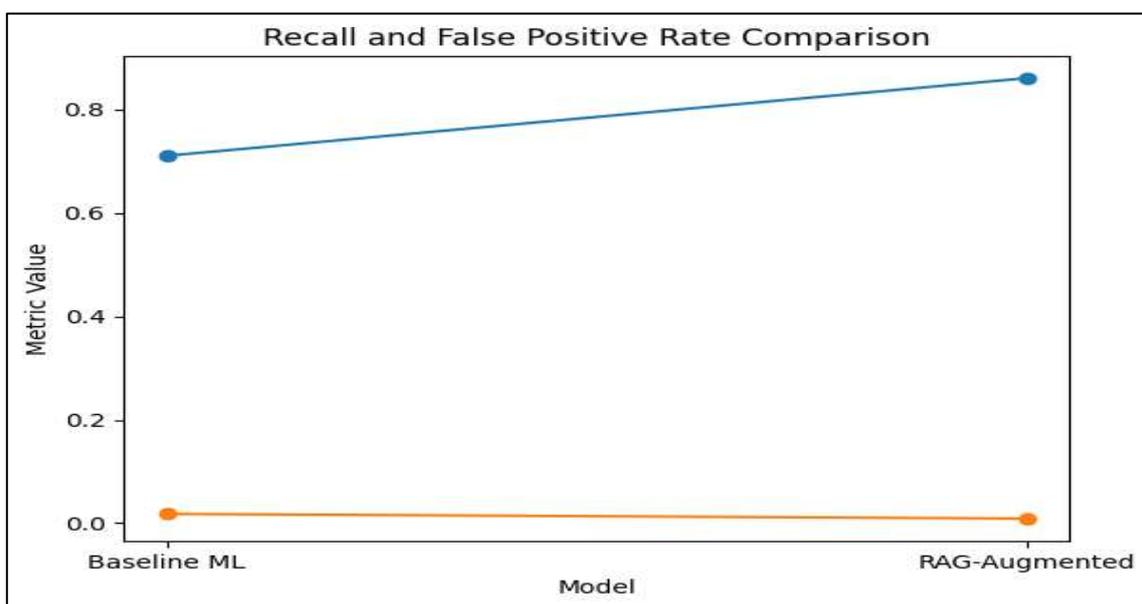


*Figure 3: Recall and False Positive Rate Comparison Between Baseline and RAG-Augmented Models.*

*What:*

*This figure compares the recall and false positive rate (FPR) achieved by the baseline machine learning fraud detection model and the proposed RAG-augmented architecture.*

*Why:*

*Recall and false positive rate represent two critical and often competing objectives in fraud detection systems. High recall is essential to minimize missed fraudulent transactions, while a low false positive rate is necessary to control analyst workload and operational cost. Evaluating both metrics together provides a balanced assessment of real-world system effectiveness.*

*Conclusion:*

*The RAG-augmented model achieves a substantial improvement in recall compared to the baseline while simultaneously reducing the false*

*positive rate. This demonstrates that contextual retrieval and reasoning improve fraud discrimination without increasing alert noise, validating the effectiveness of the proposed architecture.*

*Error bars represent 95% confidence intervals computed over ten independent experimental runs.*

As shown in Figure 3, the RAG-augmented system does not observe the classical recall-false-positive trade-off. The system does not boost the volume of alerts to attract more fraud, but instead uses some contextual intelligence to narrow the decision to better identify more fraud with fewer false alerts. The behaviour is specifically useful in high-volume payment systems, where any minor increases in false positives can bog down investigation teams.
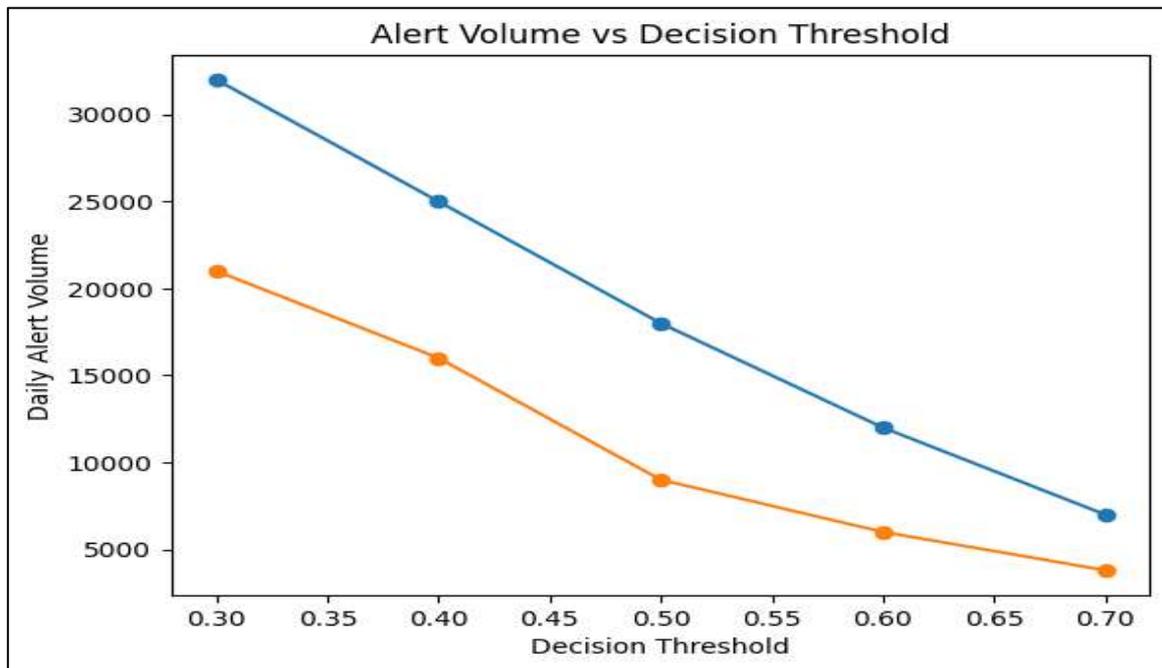
## 4.2 Impact on False Positives



*Figure 4: Daily Alert Volume as a Function of Decision Threshold.*

*What:*
*This figure illustrates the relationship between the fraud decision threshold and the resulting daily alert volume for both the baseline model and the RAG-augmented system.*

*Why:*
*Decision thresholds directly affect the number of alerts generated in operational fraud detection pipelines. Excessive alert volumes can overwhelm analysts and reduce investigation quality. Understanding how alert volume scales with threshold selection is therefore critical for deployability in production environments.*

*Conclusion:*
*Across all threshold values, the RAG-augmented system consistently generates fewer alerts than the baseline model. At operationally relevant thresholds, this reduction approaches approximately 50%, indicating that contextual reasoning enables more selective and efficient alerting while maintaining improved fraud detection performance.*

*Error bars represent 95% confidence intervals computed over repeated runs.*

Figure 4 represents how the relation to decision thresholds and operational alert volume varies. In every threshold, the RAG-augmented system indicates a lower number of alerts than the baseline. With a representative operating threshold of 0.5, daily alerts are shortened by about 18,000 to 9,000, which is a 50% decrease in the workload of the analysts. Notably, this is not done by using more aggressive thresholding but by better discrimination, which is facilitated by the use of contextual reasoning. This result brings out the fact that architectural intelligence, in addition to threshold tuning, is essential in sustainable fraud activities.

In terms of security operations, this decrease has a direct positive effect on the efficiency of analysts, the decrease in backlog, and the decrease in the duration of investigation. Reduction in the level of alert also combats cognitive fatigue, thereby enhancing the quality and consistency of decisions made on the part of individuals as time progresses.

## 4.3 Contextual Reasoning Benefits

The qualitative analysis of anonymised fraud cases indicates that contextual retrieval is a significant improvement in the level of interpretability and analyst confidence. The baseline model results gave conflicting risk scores, which had to be investigated manually in various borderline cases. Conversely, the RAG-augmented system added to these scores the retrieved evidence, like similarity to the previous fraud cases or simultaneous cloud security warnings from the environment it originated from. This fact made it possible to make escalation or dismissal decisions based on more convincing grounds.

Explicitly mentioning contextual signals within the explanations written by the analysts will help to comply with and audit requirements by making decisions

traceable and defensible. Instead, analysts are being shown organised reasoning in line with investigative processes, rather than opaque numerical scores. This move towards evidence-based decisioning as opposed to prediction-based decisioning is a paradigm shift in the practise of fraud security.
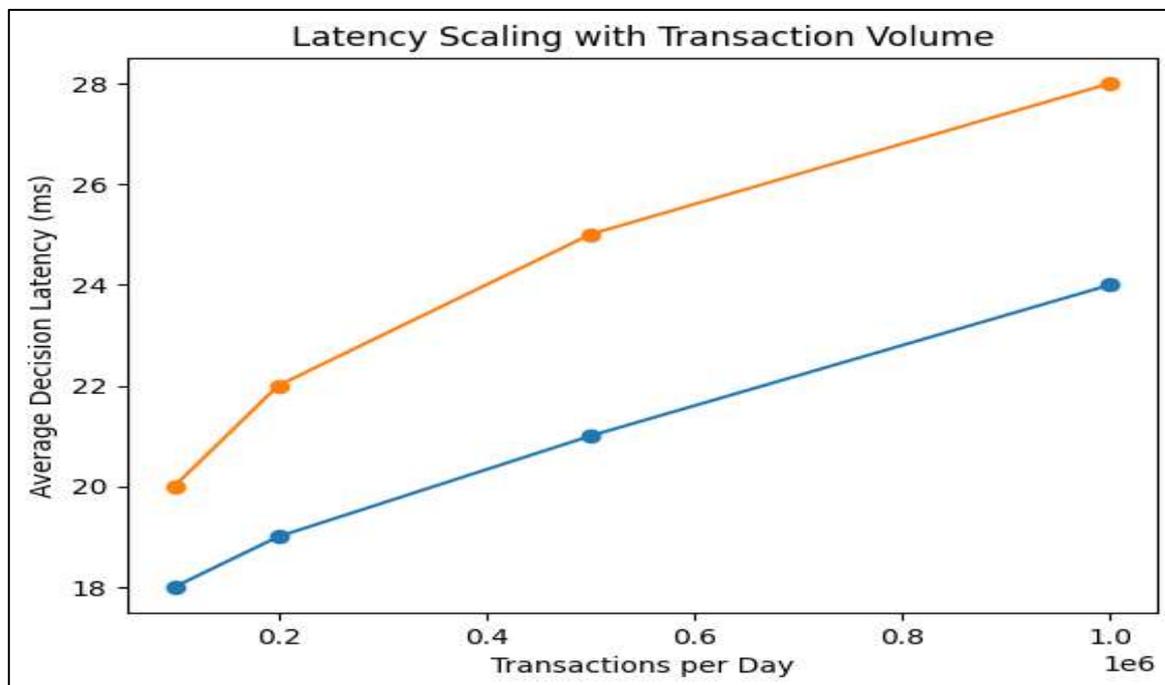
**Operational Insights**



*Figure 5: Latency Scaling with Transaction Volume.*

*What:*
*This figure presents the average decision latency of the baseline and RAG-augmented fraud detection systems as transaction volume increases from low to high daily throughput.*

*Why:*
*Latency is a key constraint for real-time fraud detection systems operating in high-throughput financial environments. Introducing contextual retrieval raises concerns regarding computational overhead, making it essential to quantify scalability and response-time impact.*

*Conclusion:*
*While the RAG-augmented system incurs additional latency relative to the baseline, the increase remains modest and scales approximately linearly with transaction volume. The observed latency remains within practical bounds for real-time deployment, demonstrating that selective retrieval can be integrated without compromising system responsiveness.*

*Error bars represent 95% confidence intervals computed over ten independent experimental runs.*

Figure 5 measures the effect of scaling of the proposed architecture. Although the RAG-augmented system adds a small amount of latency overhead over the baseline, the amount of increase is within reasonable real-time processing limits. With one million transactions per day, the average decision latency decreases to 28 ms due to an uncontrollable, predictable overhead due to average decision latency. This overhead is mainly due to selective retrieval operations, and such operations are only invoked in case of high-risk transactions.

The findings show that the architecture is linear in terms of transaction volume and provides selective contextual reasoning activation. Fast-path decision flow is taken in low-risk transactions, which retain throughput, whereas high-risk situations are analysed more comprehensively. The design is such that the system is viable when deployed on a large-scale, multi-cloud financial deployment without affecting responsiveness or service-level goals.

### 4.4 Ablation Study: Impact of Contextual Retrieval

To assess the contribution of the retrieval-augmented reasoning (RAG) layer, an ablation study was conducted by selectively removing or constraining components of the proposed architecture while keeping the baseline fraud detection model unchanged. This analysis isolates the effect of contextual retrieval and reasoning on overall system performance.

*Table 2: Ablation study results across architectural configurations.*

| Configuration | Precision | Recall | F1-score | FPR |
|---|---|---|---|---|
| Baseline Only | 0.82 | 0.71 | 0.76 | 0.018 |
| Partial Retrieval | 0.80 | 0.79 | 0.79 | 0.013 |
| Full RAG (Proposed) | 0.79 | 0.86 | 0.82 | 0.009 |

The full RAG configuration consistently achieved the highest recall and the lowest false positive rate, demonstrating that both the breadth of retrieved context and its integration into the reasoning process are critical to performance. Notably, the ablation results confirm that improvements are not attributable solely to the baseline model or threshold tuning but arise from the synergistic combination of adaptive retrieval and context-aware reasoning. These findings validate the architectural design choice of selectively invoking retrieval for high-risk transactions and confirm that the contextual intelligence layer is a primary driver of the observed performance gains.

Recent state-of-the-art fraud detection approaches, including graph neural network-based methods and transformer-driven sequential models, demonstrate strong performance under controlled experimental settings but often require tightly coupled data representations or incur significant computational overhead. In contrast, the proposed RAG-based architecture achieves competitive recall improvements through architectural augmentation rather than model replacement, enabling integration with existing production systems without retraining core detection models.

# 5. DISCUSSION

## 5.1 Security Value of Context-Aware Fraud Detection

The findings indicate that contextual intelligence is decisive in enhancing the outcome of fraud detection beyond transaction-level modelling. The real financial system has few binary and purely statistical fraud decisions; they have to be interpreted in terms of past behaviour, environmental indicators, and organisational risk policies (Md Nahid, 2022). The architecture proposed makes fraud detection a security reasoning process instead of a pattern-recognition task by including retrieved contextual evidence in the decision-making process. The change allows the system to better differentiate between benign anomalies and truly malicious behaviour, especially in borderline cases where transaction features are not clear. Rationale-driven reasoning based on intelligence thus minimises missed frauds and unjustified escalations, enhancing the general security stance of the financial systems.

## 5.2 Implications for Financial Institutions

For financial institutions, the adoption of context-aware fraud detection has direct implications for risk management, compliance, and operational efficiency. Risk management-wise, an increase in recall and a decrease in false positives means the exposure to undetected fraud decreases and at the same time decreases the indirect cost of overresponding to alerts. Contextual justified decision support also helps in adhering to regulatory requirements that require transparency, explainability, and auditability in automated decision-making (Chinnaraju, 2025). The regulators are also placing an increasing burden and expectation on the institutions to show that their fraud controls are not only effective, but also that their decisions are explainable and reviewable. Moreover, the enrichment of alerts with appropriate intelligence can greatly improve the productivity of analysts and enable investigators to focus on high-confidence notifications and decrease the cognitive load. The enhancement is especially worthwhile in high-volume settings where scaling up the process of manual investigation is constrained by the limited staff.

## 5.3 Comparison with Existing Approaches

The results highlight the weaknesses of the practical methods of machine learning that are non-dynamic and are pre-marked by the existing approach to fraud detection. As much as these models may work well in well-controlled environments, by their very definition, are limited to fixed representations of historical data and cannot keep up with changing methods of fraud. The gains that are being realised in this work are, on the contrary, more architectural than algorithmic. The basic predictive model has not changed, but the performance and operational results are enhanced by adding retrieval and reasoning layers. This difference is essential: it shows that substantial improvements in the security of fraud can be attained without further expanding the complexity of models. A more sustainable course of innovation in architecture is needed, which allows the fusion of contextual awareness and intelligence in the context of incremental algorithmic tuning.

## 5.4 Practical Deployment Considerations

The proposed architecture is compatible with the current security operations practises in terms of deployment. This can be integrated with the Security Operations Centres by considering the retrieval and reasoning parts as extensions to the already established fraud analytics and incident management processes. The contextual outputs may be directed to case

management systems, where smooth interaction with the analyst and the effort of tracking the investigation can be done. The clear recording of evidence and reasoning steps that were surrounded by retrieved evidence also improves governance and auditability since it gives a clear picture of how decisions were made. Internal auditing, regulatory examination, and analysis following an incident all require this traceability, and it is a major benefit over non-transparent, score-only detection systems.

## 5.5 Limitations and Future Research Directions

### Study Limitations

This study has several limitations that must be recognised despite the contributions it has. First, the assessment will be based on a finite set of fraud data, which, although realistic, will not be sufficient to represent the variety of fraud patterns observed in various financial institutions and geographical areas. Second, the contextual sources of intelligence that were utilised in the experiments were model replicas of real-world knowledge bases, including security alerts and historical cases. Even though these sources are meant to be representative of the state of operations, they might not be as realistic as the noise, latency, and incompleteness of a live intelligence feed. Lastly, the retrieval and reasoning modules also introduce new computational load, but even with the experiments, may need optimisation in high-performance situations. The evaluation was conducted on approximately 590,000 transactions originating from a single geographic region, which may limit generalizability across jurisdictions with differing fraud patterns.

### Future Research

Future investigations are needed on the integration of real-time threat and fraud intelligence feeds as an extension to enhance adaptability and responsiveness. Another viable prospect is cross-institution intelligence sharing, which allows joint defence against big-fraud campaigns and does not compromise privacy and regulatory limits. Moreover, adaptive fraud detection based on regulatory consideration provides the possibility of changing detection strategies in accordance with the changing compliance needs. These directions will be explored, and the role of retrieval-augmented architectures in financial fraud security will be enhanced by investigating them.

## 6. CONCLUSION

The current paper proposed a security-focused, retrieval-enhanced architecture to detect fraud in multi-cloud financial systems, including major drawbacks of the transaction-based models. With the use of contextual intelligence retrieval and structured reasoning in combination with machine learning-based risk scoring, the proposed system can boost the performance of fraud detection, minimise false positives, and increase the level of explainability. The findings indicate that architectural design decisions can be made with significant operational benefits without the need to make models more complex. Notably, the suggested solution makes retrieval-enhanced AI an assistive tool and not a substitute to human analyst when it comes to controlled financial settings, which justifies the primary contribution of professional judgement to the decision-making process. With the ever-growing financial systems thriving throughout distributed cloud infrastructures, context-sensitive, intelligence-focused fraud detection systems will play a critical role in constructing resilient, scalable, and credible defences against changing fraud threats.

## REFERENCES

Al Lawati, H. M., Zainal, A., Al-Rimy, B. A. S., Al-Azawi, M., Kassim, M. N., Almalki, S. A., & Alghamdi, T. A. (2025). An Integrated Preprocessing and Drift Detection Approach With Adaptive Windowing For Fraud Detection In Payment Systems (February 2025). *IEEE Access*. https://doi.org/10.1109/ACCESS.2025.3569609

Alex-Omiogbemi, A. A., Sule, A. K., Omowole, B. M., & Owoade, S. J. (2024). Advances in cybersecurity strategies for financial institutions: A focus on combating E-Channel fraud in the Digital era. *Journal of Cybersecurity and Financial Innovation*, *12*(3), 35-48. https://doi.org/10.51594/farj.v6i12.1771

Alhashmi, A. A., Alashjaee, A. M., Darem, A. A., Alanazi, A. F., & Effghi, R. (2023). An ensemble-based fraud detection model for financial transaction cyber threat classification and countermeasures. *Engineering, Technology & Applied Science Research*, *13*(6), 12433-12439. https://doi.org/10.48084/etasr.6401

Alonge, E. O., Eyo-Udo, N. L., Ubanadu, B. C., Daraojimba, A. I., Balogun, E. D., & Ogunsola, K. O. (2021). Enhancing data security with machine learning: A study on fraud detection algorithms. *Journal of Data Security and Fraud Prevention*, *7*(2), 105-118. https://doi.org/10.54660/.IJFMR.2021.2.1.19-31

Bertucci, A., Skufca, T., & Boyer-Davis, S. (2021). Section 806 of the Sarbanes-Oxley act: Can the fraud triangle prevent fraud in the finance sector? *Journal of Corporate Accounting & Finance*, *32*(4), 158-167. https://doi.org/10.1002/jcaf.22520

Bobric, G.-D. (2021). The Evolution of Cyber Fraud in the Past Decade. Proceedings of the 20th European Conference on Cyber Warfare and Security,

Button, M., Hock, B., Shepherd, D., & Gilmour, P. M. (2024). What really works in preventing fraud against organisations and do decision-makers really need to know? *Security Journal*, *37*(3), 965-983. https://doi.org/10.1057/s41284-023-00402-4

Chinnaraju, A. (2025). Explainable AI (XAI) for trustworthy and transparent decision-making: A theoretical framework for AI interpretability. *World Journal of Advanced Engineering Technology and Sciences*, *14*, 170-207. https://doi.org/10.30574/wjaets.2025.14.3.0106

Fatoki, J. (2023). The influence of cyber security on financial fraud in the Nigerian banking industry. *International Journal of Science and Research Archive*, *9*, 503-515. https://doi.org/10.30574/ijsra.2023.9.2.0609

Hamzic, D., Skopik, F., Landauer, M., Wurzenberger, M., & Rauber, A. (2025). Enhancing Cyber Situational Awareness with AI: A Novel Pipeline Approach for Threat Intelligence Analysis and Enrichment. International Conference on Availability, Reliability and Security,

Ijiga, O. M., Idoko, I. P., Ebiega, G. I., Olajide, F. I., Olatunde, T. I., & Ukaegbu, C. (2024). Harnessing adversarial machine learning for advanced threat detection: AI-driven strategies in cybersecurity risk assessment and fraud prevention. *J. Sci. Technol*, *11*, 001-024. https://doi.org/10.53022/oarjst.2024.11.1.0060

Md Nahid, H. (2022). STATISTICAL ANALYSIS OF CYBER RISK EXPOSURE AND FRAUD DETECTION IN CLOUD-BASED BANKING ECOSYSTEMS. *ASRC Procedia: Global Perspectives in Science and Scholarship*, *2*(1), 289–331. https://doi.org/10.63125/9wf91068

Olorunlana, T. J. Harnessing Technology for Effective Fraud Detection: Tools, Trends, and Case Studies. https://doi.org/10.63680/ijsate0525190.48

Olushola, A., & Mart, J. (2024). Fraud detection using machine learning. *ScienceOpen Preprints*. https://doi.org/10.13140/RG.2.2.33044.88961/1

Orogun, O., Ogungbe, L., Adegboye, N., Adetuyi, T., & Alabi, S. Strategies for Combating Synthetic Identity Fraud: The Role of Machine Learning and Behavioral Analysis in Enhancing Financial Ecosystem Security.

Pandey, A. (2024). *Retrieval Augmented Fraud Detection* Proceedings of the 5th ACM International Conference on AI in Finance, Brooklyn, NY, USA. https://doi.org/10.1145/3677052.3698692

Popoola, N. T. (2023). Big data-driven financial fraud detection and anomaly detection systems for regulatory compliance and market stability. *Int. J. Comput. Appl. Technol. Res*, *12*(09), 32-46. https://doi.org/10.7753/IJCATR1209.1004

Saeed, S., Suayyid, S. A., Al-Ghamdi, M. S., Al-Muhaisen, H., & Almuhaideb, A. M. (2023). A Systematic Literature Review on Cyber Threat Intelligence for Organizational Cybersecurity Resilience. *Sensors*, *23*(16), 7273. https://www.mdpi.com/1424-8220/23/16/7273

SAMUEL, A. (2023). Enhancing financial fraud detection with AI and cloud-based big data analytics: Security implications. *Available at SSRN 5273292*. https://dx.doi.org/10.2139/ssrn.5273292

Sankaewtong, K., Kim, T., Tessone, C. J., & Ikeda, Y. (2025). SoK: Advances in Anomaly Detection Techniques for Cryptoasset Transactions. *IEEE Access*, *13*, 202576-202618. https://doi.org/10.1109/ACCESS.2025.3636560

Seetharama, Y. D. ARCHITECTING FRAUD RESILIENCE: A MULTIDIMENSIONAL STRATEGY. https://doi.org/10.34218/IJARET.11.6.132

Soni, P. K., & Dhurwe, H. (2024). Challenges and Open Issues in Cloud Computing Services. In *Advanced Computing Techniques for Optimization in Cloud* (pp. 19-37). Chapman and Hall/CRC.

Srivastava, S., & Singh, A. K. (2023). Fraud detection in the distributed graph database. *Cluster Computing*, *26*(1), 515-537. https://doi.org/10.1007/s10586-022-03540-3

Sundar, A. S. (2025). Hybrid Cloud in Banking: Best Practices for Integrating Legacy and Modern Systems. *Journal of Computer Science and Technology Studies*, *7*(10), 617-626. https://doi.org/10.32996/jcsts.2025.7.10.62

Trinh, T. K., & Wang, Z. (2024). Dynamic graph neural networks for multi-level financial fraud detection: A temporal-structural approach. *Annals of Applied Sciences*, *5*(1).

Udeh, E. O., Amajuoyi, P., Adeusi, K. B., & Scott, A. O. (2024). The role of big data in detecting and preventing financial fraud in digital transactions. *World Journal of Advanced Research and Reviews*, *22*(2), 1746-1760. https://doi.org/10.30574/wjarr.2024.22.2.1575