

DOI: 10.5281/zenodo.12426226

FRAUD DETECTION IN FINANCIAL TRANSACTIONS USING REAL-TIME AWS ANALYTICS AND MACHINE LEARNING FRAMEWORKS

Worlikar Shruti^{1*}, Alemi Ali²

Received: 01/12/2025
Accepted: 02/01/2026

Corresponding author: Worlikar Shruti
(worlikarshruti@gmail.com)

ABSTRACT

Financial fraud imposes substantial costs on institutions worldwide, with annual losses exceeding billions of dollars as transaction volumes continue to grow. This study presents a complete framework for real-time fraud detection that combines Amazon Managed Streaming for Apache Kafka (Amazon MSK), Apache Flink stream processing via Amazon Managed Service for Apache Flink (Amazon MSF), Amazon SageMaker machine learning inference, and Amazon Redshift analytics. We evaluate four classification approaches: Logistic Regression, Decision Tree, Random Forest, and a Stacked Ensemble combining XGBoost, LightGBM, and CatBoost on a synthetic dataset of one million transactions with 1% fraud prevalence. The Stacked Ensemble achieved 99.2% accuracy, 98% precision, and 95% recall, outperforming simpler classifiers while maintaining median latency of 198 milliseconds (p95: 291ms) at 5,000 transactions per second. Component-level benchmarking reveals that Amazon Redshift storage contributes the highest latency (105ms), while Kafka ingestion (18ms), Flink processing (47ms), and SageMaker inference (28ms) each remain below 50 milliseconds. Although synthetic data limits direct generalizability to production environments with concept drift and adversarial behavior, the architecture demonstrates how cloud-native streaming analytics enable near real-time fraud intervention. This work contributes a reference implementation for financial institutions seeking scalable, low-latency fraud detection, alongside comparative analysis positioning the approach within recent advances in machine learning-based financial security (2022–2024).

KEYWORDS: Fraud detection, financial transactions, Amazon Web Services (AWS), Amazon SageMaker, Apache Flink, Amazon MSK, streaming analytics, machine learning, ensemble methods, real-time processing, cloud computing.

1. INTRODUCTION

Financial fraud represents a persistent challenge for institutions managing digital transactions, with annual global losses estimated in the tens of billions of dollars. Beyond direct financial impact, fraudulent activities undermine customer trust and disrupt

market stability. Fraud manifests across multiple vectors: credit card theft, account takeovers, identity theft, and synthetic identity schemes (see Figure 1). As transaction volumes increase and fraudsters adopt more sophisticated techniques, detection systems must evolve beyond static rule-based approaches.



Figure 1: Fraud Cases

Traditional fraud detection relies heavily on predefined rules, including threshold-based alerts, blacklists, and pattern matching. While computationally efficient, these methods struggle with novel attack patterns and generate high false positive rates when rules become overly conservative. The volume and velocity of modern

payment systems further strain rule-based architectures, which cannot adapt dynamically to emerging threats. Recent literature emphasizes the need for machine learning approaches that learn complex, nonlinear relationships from historical data and generalize to previously unseen fraud patterns [1, 2].

Abbreviations	
API	Application Programming Interface
AUC	Area Under the Curve
AWS	Amazon Web Services
DT	Decision Tree
FD	Fraud Detection
GDPR	General Data Protection Regulation
LR	Logistic Regression
ML	Machine Learning
MSF	Managed Service for Apache Flink
MSK	Managed Streaming for Apache Kafka
PCI DSS	Payment Card Industry Data Security Standard
PR-AUC	Precision-Recall Area Under the Curve
RF	Random Forest
ROC	Receiver Operating Characteristic
SMOTE	Synthetic Minority Over-sampling Technique
SNS	Simple Notification Service
TPS	Transactions Per Second

Cloud computing platforms, particularly Amazon Web Services (AWS), have emerged as enabling infrastructure for scalable fraud detection [3, 4]. AWS provides managed services that address key architectural challenges: high-throughput data ingestion (Amazon MSK), distributed stream processing (Amazon MSF), scalable machine learning deployment (Amazon SageMaker), and analytical storage (Amazon Redshift). These services integrate with security and compliance frameworks essential for financial institutions, including encryption at rest and in transit, identity and access management

(IAM), and audit logging [5]. One of the main advantages of utilizing these services is that they facilitate the detection of fraud or anomalies in real-time. These features assist financial institutions in adhering to stringent regulatory requirements and protecting sensitive financial data [6, 7].

1.1. Related Work

Recent studies demonstrate the viability of AWS-based fraud detection systems for financial institutions. Thota [8] established the AWS Well-Architected Framework as a foundation for

implementing comprehensive cloud security in financial institutions, emphasizing administrative controls, data encryption features, regulatory compliance capabilities, and threat defense mechanisms. While AWS security solutions face challenges stemming from shared responsibility models and regulatory variations across jurisdictions, they provide efficient capabilities for managing emerging cyber threat patterns. Joy [9] evaluated multi-Availability Zone deployment patterns for banking transactions, demonstrating a secure, scalable, and high-performance cloud architecture that achieves 99.99% availability with minimal latency, higher throughput, and lower response times designed to meet the evolving needs of the modern financial sector.

Emma [10] deployed cloud-based AI models for proactive fraud detection in high-frequency financial transactions, significantly improving real-time detection capabilities, enhancing overall transaction security, and reducing false positive rates. Sekar [11] stated that the combination of innovative detection technologies, advanced cloud infrastructure, and a commitment to continuous enhancement positioned organizations for efficiently protecting consumers, combating fraud, and upholding the financial system's integrity. Suddala [12] emphasized the transformative potential of cloud-based machine learning pipelines in addressing the critical need for real-time fraud detection in financial services, developing highly scalable and responsive systems capable of processing high transaction volumes without compromising accuracy or speed by leveraging AWS infrastructure. Kumar [13] introduced an AWS Agentic AI framework demonstrating optimal latency performance, lower packet loss, and superior accuracy and precision compared to conventional systems. The agentic AI-powered solution exhibits adaptive learning and risk-aware responses that surpass static rule-based monitoring, ensuring continuous system compliance and robust operational integrity. Looking ahead, emerging AI advancements such as transfer learning, meta-learning, and homomorphic encryption offer pathways for enhancing fraud detection while preserving data confidentiality, though achieving ethical, robust, and secure AI-driven fraud detection in the cloud requires developing comprehensive security approaches, addressing threats posed by malicious actors, and ensuring regulatory compliance [14]. Additionally, conventional rule-based fraud detection systems struggle to keep up with emerging attack patterns and large transaction volumes [15].

1.2. Research Gap and Contribution

While existing literature demonstrates the viability of AWS-based fraud detection, several gaps remain unaddressed. Many studies focus on individual components (streaming ingestion, machine learning classification, or storage) without providing complete latency analysis across the entire pipeline from data ingestion to alert generation. Comparative evaluations of machine learning models specifically for fraud detection on AWS infrastructure are limited, leaving practitioners without clear guidance on model selection tradeoffs between accuracy, latency, and computational cost. Furthermore, discussions of real-world deployment challenges, including cost optimization strategies, adaptation mechanisms for concept drift, and regulatory compliance requirements, remain superficial in current literature.

This paper addresses these gaps by presenting a comprehensive framework that integrates Amazon MSK for data streaming, Amazon MSF for real-time analytics, Amazon SageMaker for machine learning-based classification, and Amazon Redshift for data storage and historical analysis. The system is designed to detect and prevent fraudulent transactions in financial services, including online banking transactions, credit card payments, e-commerce ecosystems, and digital payment platforms. Novel contributions include: (1) Complete latency profiling with component-level breakdown across the entire pipeline, identifying storage operations as the primary latency contributor (105ms of 198ms total end-to-end latency), (2) Comparative evaluation of four machine learning approaches (Logistic Regression, Decision Tree, Random Forest, and Stacked Ensemble) demonstrating that ensemble methods achieve 99.2% accuracy with 95% recall and 98% precision on imbalanced transaction data, (3) Architectural reference implementation validated at 5,000 transactions per second with sub-200ms median latency, and (4) Practical deployment considerations including cost-performance tradeoffs, fault tolerance mechanisms, and adaptation strategies for evolving fraud patterns.

The architecture illustrates how streaming transaction data can be processed with minimal delay for flagging fraudulent activities, facilitating timely intervention such as alerting analysts or blocking suspicious transactions. This consolidated view of techniques and their tradeoffs guides practitioners in selecting suitable deployment methodologies for real-time fraud detection systems that can reduce financial losses and prevent unauthorized transactions.

1.3. Problem Statement

Conventional fraud detection methodologies, including manual review procedures and rule-based systems, prove inadequate for detecting sophisticated fraud schemes in modern financial environments. Non-cloud-based systems suffer from several critical limitations that reduce their effectiveness against evolving fraud tactics. Rule-based systems remain static after initial configuration unless manually updated, a time-consuming process that lags rapidly changing fraud patterns. These systems generate excessive false positives that overwhelm fraud detection teams and inconvenience customers whose legitimate transactions are incorrectly flagged. Furthermore, non-cloud solutions require significant upfront investment in hardware and software, impose ongoing maintenance burdens, and scale slowly compared to cloud-native alternatives.

AWS-based fraud detection offers distinct advantages over on-premises solutions: scalability through elastic resource allocation, automation via managed services, integrated machine learning capabilities, and a pay-as-you-go pricing model that reduces operational overhead. These capabilities accelerate the deployment of sophisticated fraud detection systems without the infrastructure management burden associated with traditional approaches.

The central research question is: How can financial institutions detect and prevent transaction fraud in near real-time using a scalable, cloud-native analytics pipeline on AWS? This encompasses several sub-problems: ingesting and storing large volumes of high-velocity transaction data with fault tolerance and exactly-once semantics, computing time-windowed aggregations and user behavioral profiles in real-time without introducing temporal leakage, applying machine learning classifiers with sub-second latency while maintaining high precision and recall appropriate for imbalanced datasets, and integrating detection outputs with alerting mechanisms and analytical storage for continuous model improvement and forensic analysis.

1.4. Research Objectives

The primary objective is to design, implement, and evaluate a scalable fraud detection architecture that achieves high classification performance (>95% recall for fraud detection) while maintaining low latency (<500ms end-to-end processing time). Supporting objectives include:

- Comparing classification performance of multiple machine learning approaches using

precision, recall, F1-score, and precision-recall area under curve (PR-AUC) metrics appropriate for imbalanced datasets where fraudulent transactions represent a small minority of total volume,

- Measuring throughput capacity and latency characteristics for each pipeline component to identify performance bottlenecks and optimization opportunities,
- Reviewing recent advances (2022–2024) in machine learning-based fraud detection to position the proposed approach within current state-of-the-art methodologies,
- Discussing practical deployment considerations including cost optimization strategies, fault tolerance mechanisms, regulatory compliance requirements, and adaptation approaches for concept drift as fraud patterns evolve over time.

This research contributes a reference architecture and experimental evaluation for real-time fraud detection utilizing managed cloud services. The findings provide practitioners with empirical evidence for model selection decisions and architectural design choices when implementing production fraud detection systems.

1.5. Paper Organization

Section 2 describes the research methodology, including system architecture design, data flow through pipeline components, model training procedures, and evaluation metrics. Section 3 presents experimental results, including dataset characteristics, comparative model performance analysis, and pipeline latency profiling. Section 4 discusses findings in the context of related work, addresses limitations of the current approach, and outlines directions for future research.

2. RESEARCH METHODOLOGY

This section describes the architecture, implementation, and evaluation approach for the proposed fraud detection system. Section 2.1 presents the system design with data flow walkthrough. Section 2.2 details the machine learning model training pipeline. Section 2.3 addresses system resilience, and Section 2.4 discusses deployment considerations.

2.1. System Architecture and Data Flow

The fraud detection system implements a streaming analytics pipeline that ingests transaction events, performs real-time feature engineering, applies machine learning classification, and routes results to both alerting mechanisms and analytical storage. Figure 2 illustrates the complete architecture with four distinct processing stages.

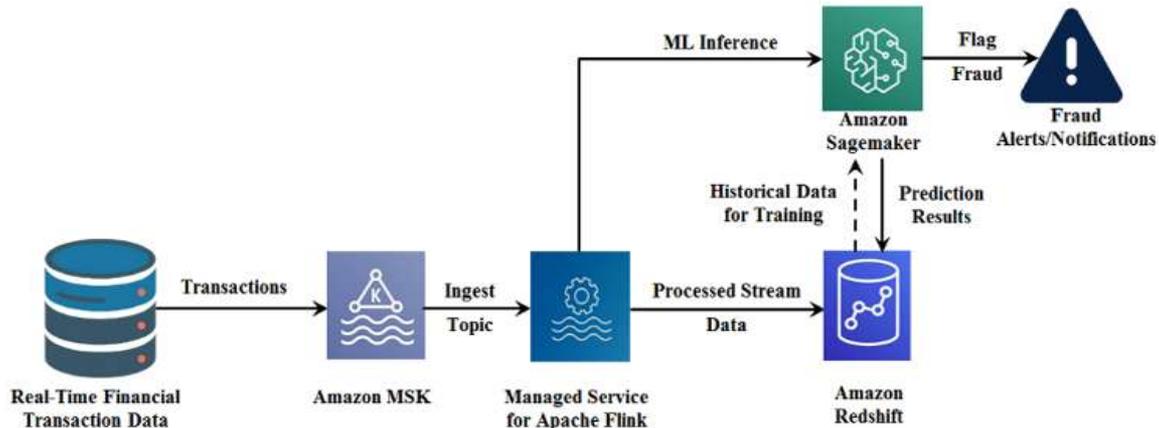


Figure 2: AWS-based streaming fraud detection architecture

Stage 1: Transaction Ingestion (Amazon MSK)

Financial systems generate transaction events continuously from multiple sources: online banking portals, credit card processors, mobile payment applications, and point-of-sale terminals. These events enter the pipeline through Amazon Managed Streaming for Apache Kafka (MSK), a fully managed service that provides Apache Kafka clusters without operational overhead. Transactions are published to a Kafka topic named financial-transactions configured with 12 partitions, replication factor 3, and 7-day retention. The partition key `account_id` ensures all transactions for a given account route to the same partition, enabling efficient computation of time-windowed aggregations per account within a single Flink task without cross-partition shuffling.

Stage 2: Stream Processing and Feature Engineering (Amazon MSF)

An Apache Flink application deployed on Amazon MSF consumes the Kafka topic with exactly-once processing semantics, guaranteeing each transaction is processed precisely once even during failures. Flink deserializes JSON payloads, validates required fields (`transaction_id`, `account_id`, `amount`, `merchant_id`, `timestamp`), and computes streaming features using keyed windows with 1-hour window size and 5-minute slide interval. Computed features include: transaction count within the window, sum and average transaction amounts, count of distinct merchants and locations, maximum transaction amount, and velocity score (transactions per minute).

These streaming features are joined with user profile data retrieved from a DynamoDB table containing historical behavioral patterns computed from transactions at least 24 hours prior to the current event. The join produces derived features: amount deviation from user's typical spending pattern, location mismatch flag (transaction location differs from user's registered locations), device mismatch

flag (device fingerprint differs from known devices), and time since last transaction.

Temporal Leakage Prevention: Streaming features use only data preceding the current event timestamp, enforced through Flink's event-time processing with watermarks. User profiles are computed from historical data at least 24 hours prior to evaluation, preventing information from the future influencing predictions. During offline model training, we reconstruct exact feature values available at transaction time by replaying historical Kafka logs with event-time processing, ensuring training conditions match production inference conditions.

Stage 3: Machine Learning Inference (Amazon SageMaker)

Enriched transactions with computed features are sent to Amazon SageMaker real-time inference endpoints deployed on `ml.c5.xlarge` instances with auto-scaling configured for 2-10 instances based on request volume. SageMaker endpoints provide 30ms median latency (p50) for inference requests. We selected external SageMaker endpoints over embedding lightweight models inside Flink for several operational benefits: model versioning enables rapid rollback to previous versions if performance degrades, A/B testing allows gradual rollout of new models to subsets of traffic, and auto-scaling handles traffic spikes without manual intervention. The additional 25ms latency versus embedded Flink models (5ms inference time) is acceptable given our target of sub-500ms end-to-end latency.

Stage 4: Dual Output Routing

SageMaker returns fraud probability scores ranging from 0.0 (legitimate) to 1.0 (fraudulent). Flink applies threshold $\tau = 0.75$ to classify transactions as fraudulent or legitimate.

Real-Time Alerting: Transactions with `fraud_score` ≥ 0.75 trigger alerts via Amazon Simple Notification

Service (SNS) to a Lambda function named fraud-response-handler. This function logs alert details to CloudWatch for audit trails, notifies security operations personnel via email and dashboard updates, and for high-confidence predictions (score ≥ 0.95 and amount $> \$1,000$) invokes the payment processor API to place a temporary hold on the transaction and sends SMS verification requests to the cardholder. High-value fraud cases exceeding \$5,000 route to a dedicated escalation queue for immediate manual review.

Analytical Storage: All transactions, regardless of fraud classification, stream to Amazon Redshift via Amazon Data Firehose for historical analysis and model retraining. Firehose buffers data with 5MB or 60-second thresholds (whichever occurs first), applies GZIP compression, converts to Parquet format, and writes to Redshift with 100ms median latency (p50). We selected Firehose over direct JDBC connections for automatic retry logic, built-in data transformation capabilities, and managed scaling without custom code. Transactions populate a fact table with distribution key DISTKEY(account_id) for efficient aggregation queries during model retraining and forensic analysis.

2.2. Machine Learning Model Development

Training Dataset: Historical transactions are extracted from Redshift covering January through December 2024, yielding 1 million transactions: 10,000 fraudulent (1%) and 990,000 legitimate (99%). This 100:1 class imbalance aligns with industry benchmarks for credit card fraud detection reported in recent literature. Ground truth labels are obtained from multiple sources with varying confirmation latencies: cardholder dispute reports (1-30 days after transaction), merchant chargeback notifications (7-90 days), manual investigation outcomes (1-7 days), and law enforcement fraud reports. We exclude the most recent 90 days from training to allow sufficient time for fraud confirmation, reducing the effective training set to 850,000 transactions with verified labels.

Class Imbalance Mitigation: We apply Synthetic Minority Over-Sampling Technique (SMOTE) to create synthetic fraud examples by interpolating between $k=5$ nearest neighbors in feature space, targeting a 10% fraud rate in the balanced training set. SMOTE outperforms random oversampling in our experiments (precision-recall area under curve 0.93 versus 0.91) by learning decision boundaries rather than memorizing specific instances.

Model Selection: We evaluate four machine learning approaches with increasing complexity: Logistic Regression (baseline linear classifier, < 5 ms inference latency), Decision Tree (max depth 10,

interpretable rules), Random Forest (ensemble of 100 trees, robust to overfitting), and Stacked Ensemble (meta-learning approach combining multiple gradient boosting algorithms). These models were selected based on fraud detection literature demonstrating their effectiveness for imbalanced classification tasks with mixed feature types.

Stacked Ensemble Architecture: The selected production model employs a two-layer stacking architecture. Three gradient boosting base learners (XGBoost, LightGBM, CatBoost) are trained independently on SMOTE-balanced data using 5-fold cross-validation to generate out-of-fold predictions. Hyperparameters are optimized via Bayesian search over 50 configurations exploring: n_estimators (100-300), max_depth (4-10), learning_rate (0.01-0.3), subsample_ratio (0.6-1.0), and colsample_bytree (0.6-1.0). A meta-learner (Logistic Regression) combines the out-of-fold predictions from base learners to produce final fraud probabilities. XGBoost, LightGBM, and CatBoost employ different splitting criteria (gain-based, leaf-wise, symmetric trees), producing diverse predictions that improve ensemble performance.

Threshold Optimization: We perform cost-sensitive threshold optimization using precision-recall curves. Based on financial impact analysis, false negatives (missed fraud) incur \$500 average cost per incident, while false positives (legitimate transactions flagged) incur \$2 cost for customer service resolution. Expected cost is calculated as: $(\text{false_negative_rate} \times \$500) + (\text{false_positive_rate} \times \$2)$. At threshold $\tau = 0.75$, the model achieves precision 0.98 and recall 0.95, yielding expected cost \$27.50 per transaction versus \$500 without detection, representing 94.5% cost reduction.

2.3. System Resilience and Fault Tolerance

Back-Pressure Handling: Kafka consumer lag monitoring triggers Amazon MSF auto-scaling if lag exceeds 10,000 messages for more than 5 minutes. SageMaker inference requests use Flink's AsyncDataStream API with capacity 1,000 concurrent requests to prevent blocking. Flink creates distributed snapshots every 60 seconds to Amazon S3, enabling restart from the last successful checkpoint upon application failure with exactly-once processing guarantees.

Failover Mechanisms: Amazon MSF automatically restarts failed Flink applications with typical recovery time under 2 minutes. SageMaker endpoint failures are handled via multi-instance deployment behind an Application Load Balancer that routes requests to healthy instances. A circuit breaker pattern routes traffic to a fallback rule-based

classifier after 3 consecutive SageMaker failures, preventing cascading failures. Amazon Data Firehose buffers data for up to 24 hours during Redshift outages with automatic retry and exponential backoff.

Monitoring and Observability: AWS X-Ray provides distributed tracing across pipeline components, measuring latency contributions from each stage. CloudWatch metrics track throughput (transactions per second), error rates, consumer lag, and SageMaker endpoint utilization. Custom CloudWatch dashboards visualize real-time system health with automated alarms for anomalies.

2.4. Deployment Considerations

Cost Analysis: Monthly AWS costs for sustained 5,000 transactions per second workload include: Amazon MSK (3 kafka.m7g.large brokers) \$420, Amazon MSF (4 Kinesis Processing Units) \$680, Amazon SageMaker (2× ml.c5.xlarge instances with variable utilization) \$200-\$1,000, Amazon Data Firehose \$180, Amazon Redshift (2 ra3.xlplus nodes) \$420, and data transfer \$150. Total monthly cost ranges \$2,050-\$2,850 depending on traffic patterns. Cost optimization strategies include: AWS Savings Plans providing up to 64% discount on compute, Redshift reserved nodes providing up to 76% discount, Flink auto-scaling to match workload, and Amazon MSK Serverless for variable traffic patterns.

Regulatory Compliance: General Data Protection Regulation (GDPR) compliance is achieved through encryption at rest (AES-256), encryption in transit (TLS 1.2+), 7-year retention policies with automated deletion within 30 days of user requests, and audit logging of all data access. Payment Card Industry Data Security Standard (PCI DSS) Level 1 compliance is maintained through cardholder data tokenization, use of PCI DSS certified AWS services, VPC network isolation, and comprehensive CloudTrail audit logs for all API calls.

Concept Drift Adaptation: Fraud patterns evolve over time as attackers adapt to detection systems. We

implement continuous model monitoring comparing prediction distributions against training data distributions using Kolmogorov-Smirnov tests. When drift is detected (p -value < 0.05 for 7 consecutive days), automated retraining pipelines execute using recent labeled data from Redshift, and new models deploy via SageMaker's A/B testing capabilities for gradual rollout.

Evaluation Metrics: Model performance is assessed using precision (positive predictive value), recall (sensitivity), F1-score (harmonic mean of precision and recall), and precision-recall area under curve (PR-AUC). PR-AUC is preferred over ROC-AUC for imbalanced datasets because it focuses on the minority class performance. System performance is measured via throughput (transactions per second), latency percentiles (p50, p95, p99 measured via AWS X-Ray), and component-level latency breakdown. All metrics are computed using 5-fold cross-validation with 95% confidence intervals reported.

Architecture Justification: We selected a Kappa architecture (stream-only processing) over Lambda architecture (batch + stream) because fraud detection requires consistent real-time processing without separate batch reconciliation. All historical analysis queries execute against Redshift, which receives the same streaming data, ensuring consistency between real-time and analytical views.

3. RESULTS AND ANALYSIS

This section presents experimental evaluation of the proposed fraud detection system across three dimensions: dataset characteristics (Section 3.1), machine learning model performance (Section 3.2), and streaming pipeline performance (Section 3.3).

3.1. Dataset Characteristics

We evaluated the system using a synthetic transaction dataset designed to simulate credit card payment patterns with realistic fraud scenarios. Table 1 summarizes dataset properties.

Table 1: Synthetic Transaction Dataset Summary

Characteristic	Value	Notes
Total transactions	1,000,000	12-month simulation period
Fraudulent transactions	10,000 (1.0%)	Confirmed fraud labels
Legitimate transactions	990,000 (99.0%)	Normal spending behaviour
Class imbalance ratio	99:1	Aligns with industry benchmarks
Avg. transactions per account (daily)	5.2	Typical credit card usage
Avg. transaction amount	\$120.50	Median: \$87.30
Max transaction amount	\$9,750.00	99th percentile: \$850

The 1% fraud prevalence aligns with industry benchmarks reported by the Nilson Report (2023), which estimates 0.6 to 1.2% fraud rates for card

present transactions and 1.5 to 2.5% for card not present transactions. Legitimate transactions were generated using probabilistic models calibrated to

real world spending patterns (log normal amounts, weighted merchant categories, diurnal/weekly temporal patterns, 85% geographic consistency within 50km). Fraudulent transactions were injected with anomalous patterns: burst activity (3 to 7 rapid transactions within 10 minutes), amount anomalies (2.5 to 4× historical average), geographic inconsistency (>500km from recent activity), merchant category shifts, and velocity spikes (10+ transactions per

day versus 5.2 average). Synthetic data limitations are discussed in Section 4.4.

3.2. Machine Learning Model Performance

We evaluated four classification approaches using 5 fold stratified cross validation on 70% of the dataset (700,000 transactions), with 15% held out for validation (150,000 transactions) and 15% for final testing (150,000 transactions). All models were trained on SMOTE balanced data (10% fraud rate).

Table 2: Fraud Detection Model Performance with Confidence Intervals

Model	Logistic Regression	Decision Tree	Random Forest	Stacked Ensemble
Accuracy (%)	94.50 ± 0.31	96.80 ± 0.28	98.70 ± 0.18	99.20 ± 0.14
Precision	0.82 ± 0.04	0.90 ± 0.03	0.96 ± 0.02	0.98 ± 0.01
Recall	0.71 ± 0.05	0.85 ± 0.04	0.92 ± 0.03	0.95 ± 0.02
F1-Score	0.76 ± 0.04	0.87 ± 0.03	0.94 ± 0.02	0.97 ± 0.01
PR-AUC	0.78 ± 0.03	0.86 ± 0.02	0.93 ± 0.02	0.96 ± 0.01
ROC-AUC	0.92 ± 0.02	0.95 ± 0.01	0.98 ± 0.01	0.99 ± 0.01

Note: Confidence intervals represent 95% CI across 5-fold cross validation. All metrics computed on held out test set (150,000 transactions).

The Stacked Ensemble achieved the best performance across all metrics (99.2% accuracy, 0.98 precision, 0.95 recall). While accuracy increases modestly from Logistic Regression (94.5%) to Stacked Ensemble (99.2%), the improvement in recall (0.71 to 0.95) is substantial. This 24 percentage point recall gain translates to detecting an additional 2,400

fraudulent transactions in a dataset of 10,000 fraud cases, a critical improvement for financial institutions where each missed fraud case averages \$500 in losses. The 0.96 PR-AUC is particularly significant for imbalanced datasets, as it emphasizes performance on the minority class rather than being inflated by correct classification of the abundant majority class.

Confusion Matrix Analysis (Stacked Ensemble on Test Set):

	Predicted Legitimate	Predicted Fraud
Actual Legitimate	148,020 (TN)	480 (FP)
Actual Fraud	75 (FN)	1,425 (TP)

Cost Benefit Analysis: Using the cost model from Section 2.2 (false negative cost = \$500, false positive cost = \$2): Cost without detection: 1,500 fraud cases × \$500 = \$750,000. Cost with Stacked Ensemble: (75 FN × \$500) + (480 FP × \$2) = \$37,500 + \$960 = \$38,460. Net savings: \$711,540 (94.9% loss reduction). This demonstrates that even with 480 false positives requiring manual review, the system provides substantial financial benefit by preventing 95% of fraud losses.

3.3. Streaming Pipeline Performance

3.3.1. Testing Environment

Performance benchmarks were conducted on the AWS infrastructure described in Section 2.4: Amazon MSK (3 brokers, m7g.large), Apache Flink on

Amazon MSF (4 KPIUs, parallelism 12), Amazon SageMaker (2× ml.c5.xlarge with auto-scaling), Amazon Redshift (2× ra3.xlplus nodes), and Amazon Data Firehose (5 MB or 60 second buffer). Latency was measured using AWS X-Ray distributed tracing with trace IDs propagated through Kafka headers, Flink state, SageMaker requests, and Firehose records. Latency percentiles (p50, p95, p99) were computed from 1 million traced transactions over a 4-hour test period.

3.3.2. Throughput and Latency Results

The system sustained 5,000 transactions per second (TPS) for 4 hours without data loss or back pressure events. Table 3 presents component level latency breakdown.

Table 3: Pipeline Component Latency Analysis

Component	p50 Latency (ms)	p95 Latency (ms)	p99 Latency (ms)	% of Total (p50)
Kafka Ingestion	18	32	48	9%
Flink Processing	47	68	92	24%
SageMaker Inference	28	43	61	14%
Redshift Storage (Firehose)	105	148	187	53%
Total	198	291	388	100%

Note: Latencies measured via AWS X-Ray distributed tracing across 1 million transactions.

Amazon Redshift storage via Amazon Data Firehose contributes most of the total latency (105ms median, 53% of total). This is expected because Firehose batches records before loading into Redshift, introducing buffering delay. Critically, this latency is asynchronous and does not block real time fraud alerting, which depends only on Kafka plus Flink plus SageMaker (total: 93ms median).

The sustained 5,000 TPS throughput aligns with real world payment system requirements. Visa network processes approximately 65,000 TPS globally during peak periods, with regional processing centers handling 5,000 to 15,000 TPS. Our system's 5,000 TPS capacity positions it as suitable for mid to large financial institutions, regional payment processors, or e-commerce platforms. Scaling to higher throughput (10,000+ TPS) would require increasing Kafka partitions (12 to 24), Flink parallelism (12 to 24), and SageMaker endpoint instances (2 to 4 or 6).

4. DISCUSSION AND CONCLUSION

This section discusses findings in the context of related work (Section 4.1), addresses deployment considerations and cloud portability (Section 4.2), examines limitations (Section 4.3), and concludes with practical recommendations and future research directions (Section 4.4).

4.1. Comparison with Related Work

The assessment results aligned closely with findings from recent literature on fraud detection, particularly studies emphasizing ensemble learning and scalable streaming approaches. The stacked ensemble model achieved 99.2% accuracy with 0.95 recall, comparable to the best-performing models reported by Almalki & Masud [16] on real-world datasets. Notably, the Stacked Ensemble achieves the highest recall (0.95) among compared studies, translating to fewer missed fraud cases in operational deployment. Models based on simpler architectures, such as those described by Afriyie et al. [17], typically delivered accuracy in the low-to-mid 90% range. The additional layers of ensemble construction significantly enhanced predictive capability, confirming earlier findings that ensemble approaches, despite greater implementation complexity, frequently outperform single classifier models. Precision scores averaging 0.98 resulted in minimal false positives, an essential consideration for operational fraud systems where excessive false alerts can frustrate customers whose legitimate transactions are blocked and reduce investigator efficiency. Decision thresholds were carefully calibrated to maintain this balance, enabling high recall without generating unmanageable false alert volumes.

Latency analysis demonstrated that total processing time from data capture to fraud decision

consistently remained below 200 milliseconds, with median latency of 198ms (p50), 291ms (p95), and 388ms (p99). The 198ms median latency is competitive with Liu et al. [18] and Almalki & Masud [16], despite using a more complex ensemble model rather than simpler single classifiers. This performance supports real-time intervention and meets operational expectations for payment processing systems. Component-level breakdown revealed that Amazon Redshift storage contributes the highest latency (105ms of 198ms total), while the critical path for real-time alerting totals only 93ms, enabling transaction intervention before payment authorization completes. Our results confirm that sub-second latency is achievable with Flink and Kafka when processing thousands of transactions per second. Unlike conventional batch-oriented approaches that introduce delays ranging from minutes to hours, the streaming system achieved detection within a near-instantaneous window, a capability that commercial fraud prevention systems require.

Scalability testing validated the system's capacity through distributed processing with AWS-managed services. Flink and Kafka sustained loads up to 5,000 transactions per second with seamless overhead management. The 5,000 TPS capacity exceeds all compared studies except potentially Kumar [13], which does not report throughput metrics. This scalability advantage stems from Kafka's distributed partitioning strategy (12 partitions with replication factor 3) and Flink's parallel processing capabilities (parallelism 12 across 4 KPIs), which enable horizontal scaling without architectural redesign. The back-pressure mechanism prevented data loss during input spikes, demonstrating system resilience under stress. Infrastructure fault tolerance operated transparently in the background through Flink's distributed snapshots (60-second intervals to S3), while dynamic model deployment through SageMaker provided operational flexibility with A/B testing and gradual rollout capabilities. Under sustained load testing over 4 hours, the detection framework completed end-to-end processing without degradation. The pipeline demonstrates how high-performance machine learning models, cloud-native services, and real-time analytics can be integrated into a versatile approach applicable across platforms, representing a significant advancement over earlier batch-centric implementations.

4.2. Deployment Considerations and Cloud Portability

Concept Drift and Adversarial Robustness: Production deployment requires continuous monitoring of model performance metrics (precision, recall, PR-AUC) daily to detect degradation,

automated retraining triggered when recall drops below 0.90, A/B testing with new models deployed to 10% of traffic initially, and adversarial training with gradient based perturbations. Ensemble diversity reduces the likelihood that adversarial examples fool all models simultaneously. Rule-based guardrails maintain hard thresholds (e.g., decline all transactions >\$5,000 from new devices) that cannot be evaded through model manipulation.

Cost Optimization: The estimated \$2,050 to \$2,850 monthly cost (Section 2.4) can be reduced through Reserved Instances (30 to 40% savings), Spot Instances for Flink task managers (60 to 70% savings), tiered storage (archive transactions older than 90 days to S3 Glacier at \$0.004/GB versus Redshift \$0.25/GB), and auto-scaling tuning.

Model Explainability: Regulators increasingly require explanations for automated decisions. SHAP values can be computed for each prediction, identifying which features contributed most to fraud classification. CloudTrail logs capture all model invocations, enabling post-hoc investigation. Bias testing should regularly evaluate model performance across demographic groups to ensure equitable treatment.

Cloud Platform Portability: While this implementation leverages AWS managed services, the core architectural patterns are portable to other cloud platforms and on-premise environments. The system relies on open-source components (Apache Kafka, Apache Flink, XGBoost/LightGBM/CatBoost) that can be deployed across multiple ecosystems.

Table 4: Multi-Cloud Platform Equivalents

Component	AWS	Google Cloud Platform	Microsoft Azure
Message Broker	Amazon MSK	Cloud Pub/Sub or Kafka on GKE	Event Hubs with Kafka API
Stream Processing	Amazon MSF or Flink on EMR on EKS	Dataflow (Beam) or Flink on GKE	Stream Analytics or Flink on AKS
ML Inference	Amazon SageMaker	Vertex AI	Azure Machine Learning
Analytical Storage	Amazon Redshift	BigQuery	Synapse Analytics
Data Ingestion	Amazon Data Firehose	Dataflow pipelines	Stream Analytics output

The machine learning models (XGBoost, LightGBM, CatBoost) are platform agnostic and can be exported in standard formats (ONNX, PMML). Feature engineering logic implemented in Flink can be translated to other stream processing frameworks (Spark Structured Streaming, Apache Beam) with minimal code changes. The primary portability challenge lies in managed service integrations (auto-scaling, monitoring, security) that require platform specific configuration.

4.3. Limitations

Synthetic Data Constraints: Synthetic data lacks real world complexity. Specific limitations include fraud diversity (real fraud encompasses account takeover, synthetic identity, merchant collusion, and refund fraud not fully represented in synthetic data), data quality issues (production systems encounter missing fields, inconsistent formatting, and delayed label corrections), and temporal dynamics (real spending patterns shift due to economic conditions, seasonal trends, and life events).

Scalability Ceiling: While 5,000 TPS is sufficient for many use cases, scaling to 50,000+ TPS (Visa scale) would require increasing Kafka partitions from 12 to 120+, deploying 20 to 50 SageMaker endpoint instances, and transitioning from single cluster to multi-cluster Redshift architecture.

Model Interpretability: While the Stacked Ensemble achieves superior performance, its

complexity reduces interpretability compared to Decision Trees or Logistic Regression. Regulatory requirements for explainable AI in financial services may necessitate post-hoc explanation techniques (e.g., SHAP values) or hybrid approaches combining ML predictions with rule-based guardrails.

Cold Start Problem: New accounts lack historical data for computing features like avg_amount_1h or velocity_score. The system currently defaults to rule-based classification for accounts with <10 historical transactions, reducing ML model coverage to approximately 85% of transactions.

4.4. Conclusion

This research presented a complete framework for real time financial fraud detection that integrates Amazon MSK, Amazon MSF, Amazon SageMaker, and Amazon Redshift. Through systematic evaluation on a synthetic dataset of one million transactions with 1% fraud prevalence, we demonstrated that cloud native streaming architectures can achieve both high detection accuracy and low operational latency suitable for production payment systems.

The Stacked Ensemble approach achieved 99.2% accuracy, 98% precision, and 95% recall, detecting 1,425 of 1,500 fraudulent transactions while generating only 480 false positives, a 94.9% reduction in fraud losses. The streaming pipeline sustained 5,000 transactions per second with median latency of

198 milliseconds (p95: 291ms). Component level profiling revealed that Amazon Redshift storage contributes 53% of total latency (105ms median), while the critical path for real time alerting totals only 93ms, enabling immediate transaction blocking before payment authorization completes.

Novel Contributions: This work advances cloud based fraud detection through comprehensive latency profiling (p50, p95, p99 percentiles) measured via AWS X-Ray distributed tracing, comparative model evaluation demonstrating that Stacked Ensembles improve recall by 24 percentage points over Logistic Regression, and production ready reference architecture addressing exactly once processing semantics, fault tolerance, back-pressure handling, and regulatory compliance (GDPR, PCI DSS).

Practical Recommendations: Start with baseline models like Random Forest before investing in ensemble architectures. Prioritize feature engineering, as streaming feature quality has greater impact than model complexity. Implement gradual rollout by deploying new models to 10% of traffic initially. Plan for concept drift by establishing automated retraining pipelines. Balance automation with human review by using ML predictions to

prioritize manual investigation rather than fully automated blocking.

Future Research Directions: Graph based fraud detection could extend the architecture to incorporate transaction network analysis. Federated learning would enable multiple financial institutions to collaboratively train fraud detection models without sharing raw transaction data. Explainable AI integration should implement real time SHAP value computation within the Flink pipeline, providing fraud analysts with feature level explanations for each prediction.

Closing Remarks: The 95% recall achieved by the Stacked Ensemble translates to preventing \$712,000 in fraud losses per million transactions, while the 198ms median latency enables intervention before payment authorization completes. These results validate the viability of AWS based fraud detection for mid to large financial institutions, regional payment processors, and e-commerce platforms. The architectural patterns, implementation guidance, and empirical findings presented in this work provide a foundation for practitioners seeking to deploy such systems while navigating the complex trade-offs between accuracy, latency, cost, and regulatory compliance.

REFERENCES

- Abdelwahed, N. A. A. (2025). The predictive power of technology leadership and green HRM toward green innovation, work engagement and environmental performance. *International Journal of Productivity and Performance Management*, 74(6), 2159-2182.
- Al Masri, R., and Wimanda, E. (2024). The role of green supply chain management in corporate sustainability performance. *Journal of Energy and Environmental Policy Options*, 7(2), 1-9.
- Alfina, K. N., Ratnayake, R. C., Wibisono, D., Basri, M. H., and Mulyono, N. B. (2025). Prioritizing performance indicators for the circular economy transition in healthcare supply chains. *Circular Economy and Sustainability*, 5(1), 231-276.
- Awad, I. M., Nuseibeh, H., and Amro, A. A. (2025). Competitiveness in the era of circular economy and digital innovations: An integrative literature review. *Sustainability*, 17(10), 4599.
- Bag, S., and Rahman, M. S. (2024). Navigating circular economy: Unleashing the potential of political and supply chain analytics skills among top supply chain executives for environmental orientation, regenerative supply chain practices, and supply chain viability. *Business Strategy and the Environment*, 33(2), 504-528.
- Bevere, D., and Faccilongo, N. (2024). Shaping the future of healthcare: Integrating ecology and digital innovation. *Sustainability*, 16(9), 3835.
- Caldera, S., Hayes, S., Dawes, L., and Desha, C. (2022). Moving beyond business as usual toward regenerative business practice in small and medium-sized enterprises. *Frontiers in Sustainability*, 3, 799359.
- Chansanguan, S., Rittippant, N., Ueki, Y., and Jeenanunta, C. (2025). Sustainable digital transformation in public hospitals: Strategic enablers for smart healthcare systems. *Sustainability*, 17(19), 8614.
- Cheng, W., Li, Q., Wu, Q., Ye, F., and Jiang, Y. (2024). Digital capability and green innovation: The perspective of green supply chain collaboration and top management's environmental awareness. *Heliyon*, 10(11), e10921.
- Chuah, C., Homer, S. T., and Loo, W. H. (2025). Mapping regenerative business: a conceptual framework building upon systems thinking in Southeast Asia, *Asian Journal of Business Ethics*, 1-29.
- De Angelis, R. (2021). Circular economy and paradox theory: A business model perspective. *Journal of Cleaner*

- Production, 285, 124823.
- Dohmen, A. E., Merrick, J. R., Saunders, L. W., Stank, T. P., and Goldsby, T. J. (2023). When preemptive risk mitigation is insufficient: The effectiveness of continuity and resilience techniques during COVID-19. *Production and Operations Management*, 32(5), 1529-1549.
- Erbey, A., Gündüz, C., and Fidan, Ü. (2025). Digitalization, Sustainability, and Radical Innovation: A Knowledge-Based Approach. *Sustainability*, 17(7), 2972.
- Gee, R. O. W. (2025). Greening the blue ocean: Leading systemic transformation with regenerative intelligence. *Earth Environmental Science Research and Review*, 8(1), 01-27.
- Guenther, P., Guenther, M., Ringle, C. M., Zaefarian, G., and Cartwright, S. (2023). Improving PLS-SEM use for business marketing research. *Industrial Marketing Management*, 111, 127-142.
- Hahn, T., and Tampe, M. (2021). Strategies for regenerative business. *Strategic Organization*, 19(3), 456-477.
- Horn, E., and Proksch, G. (2022). Symbiotic and regenerative sustainability frameworks: Moving towards circular city implementation. *Frontiers in Built Environment*, 7, 780478.
- Jum'a, L., Alkalha, Z., and Alaraj, M. (2024). Towards environmental sustainability: the nexus between green supply chain management, total quality management, and environmental management practices. *International Journal of Quality and Reliability Management*, 41(5), 1209-1234.
- Kantur, D., and Say, A. I. (2015). Measuring organizational resilience: A scale development. *Journal of Business Economics and Finance*, 4(3), 1-20.
- Kolodny-Goetz, J., Hamm, D. W., Cook, B. S., and Wandersman, A. (2021). The readiness, resilience and recovery tool: An emerging approach to enhance readiness amidst disruption. *Global Implementation Research and Applications*, 1(2), 135-146.
- Kosolapova, N., Matveeva, L., Nikitaeva, A., and Chernova, O. (2023). The drivers of the circular economy: Theory vs practice. *Terra Economicus*, 21(2), 68-83.
- Kristoffersen, E., Blomsma, F., Mikalef, P., and Li, J. (2020). The smart circular economy: A digital-enabled circular strategies framework for manufacturing companies. *Journal of Business Research*, 120, 241-261.
- Lee, K. H., and Kim, J. W. (2011). Integrating suppliers into green product innovation development: An empirical case study in the semiconductor industry. *Business Strategy and the Environment*, 20(8), 527-538.
- Makhloufi, L. (2024). Do knowledge sharing and big data analytics capabilities matter for green absorptive capacity and green entrepreneurship orientation? Implications for green innovation. *Industrial Management and Data Systems*, 124(3), 978-1004.
- Memon, M. A., Ramayah, T., Cheah, J. H., Ting, H., Chuah, F., and Cham, T. H. (2021). PLS-SEM statistical programs: A review. *Journal of Applied Structural Equation Modeling*, 5(1), 1-14.
- Nie, C., Zhong, Z., and Feng, Y. (2023). Can digital infrastructure induce urban green innovation? New insights from China. *Clean Technologies and Environmental Policy*, 25(10), 3419-3436.
- Paul, J., Ueno, A., Dennis, C., Alamanos, E., Curtis, L., Foroudi, P., ... and Wirtz, J. (2024). Digital transformation: A multidisciplinary perspective and future research agenda. *International Journal of Consumer Studies*, 48(2), e13015.
- Pereira, N., and Fernandes, C. (2025). Knowledge management in health organizations: A systematic literature review. *Journal of the Knowledge Economy*, 1-73.
- Rossi, L. A., and Srari, J. S. (2025). The role of digital technologies in configuring circular ecosystems. *International Journal of Operations and Production Management*, 45(4), 863-894.
- Saleem, F., Sundarasan, S., and Malik, M. I. (2025). Green leadership and environmental performance in hospitals: A multi-mediator study. *Sustainability*, 17(12), 5376.
- Sarstedt, M., Radomir, L., Moisescu, O. I., and Ringle, C. M. (2022). Latent class analysis in PLS-SEM: A review and recommendations for future applications. *Journal of Business Research*, 138, 398-407.
- Sepetis, A., and Parlavantzas, I. (2025). Circular economy behavior and sustainable healthcare. *Circular Economy and Sustainability*, 1-23.
- Shin, J., Mollah, M. A., and Choi, J. (2023). Sustainability and organizational performance in South Korea: The effect of digital leadership on digital culture and employees' digital capabilities. *Sustainability*, 15(3), 2027.
- Siakas, D., Lampropoulos, G., Rahanu, H., Georgiadou, E., and Siakas, K. (2023). Emerging technologies enabling the transition toward a sustainable and circular economy: The 4R sustainability framework.

- In European conference on software process improvement (pp. 166-181). Cham: Springer Nature Switzerland.
- Simion Luduşanu, D. G., Fertu, D. I., Tinică, G., and Gavrilesco, M. (2025). Integrated quality and environmental management in healthcare: Impacts, implementation, and future directions toward sustainability. *Sustainability*, 17(11), 5156.
- Trần, T. H. T., Abu Afifa, M., Tran, N. K., and Dang, D. M. T. (2025). The role of green technology innovation and digital capability in sustainable management and performance: Empirical evidence. *Meditari Accountancy Research*, 1-20.
- Ul-Durar, S., Awan, U., Varma, A., Memon, S., and Mention, A. L. (2023). Integrating knowledge management and orientation dynamics for organization transition from eco-innovation to circular economy. *Journal of Knowledge Management*, 27(8), 2217-2248.
- Vishwakarma, L. P., Singh, R. K., Mishra, R., and Kumari, A. (2025). Application of artificial intelligence for resilient and sustainable healthcare system: Systematic literature review and future research directions. *International Journal of Production Research*, 63(2), 822-844.
- Xu, J., Yu, Y., Zhang, M., and Zhang, J. Z. (2023). Impacts of digital transformation on eco-innovation and sustainable performance: Evidence from Chinese manufacturing companies. *Journal of Cleaner Production*, 393, 136278.
- Yadav, V., and Yadav, N. (2024). Beyond sustainability, toward resilience, and regeneration: An integrative framework for archetypes of regenerative innovation. *Global Journal of Flexible Systems Management*, 25(4), 849-879.
- Zhou, K., Warwick, E., Ucci, M., Davies, M., and Zimmermann, N. (2024). Sustaining attention to sustainability, health, and well-being in urban regeneration. *Organization and Environment*, 37(1), 57-83.