

# Explainable AI for Enhanced Anomaly Detection in Fraud Detection

Reza Amiri<sup>1\*</sup>, Mohammad Hadi Zahedi<sup>2</sup>, and Mehdi Azadimotlagh<sup>3</sup>

<sup>1</sup> Faculty member of the Advanced Information System Research Group, ICTRC, ACECR

<sup>2</sup> Assistant Professor, Information Technology Dept., Faculty of Industrial Engineering, K. N. Toosi University of Technology, Tehran, Iran

<sup>3</sup> Department of Computer Engineering of Jam, Persian Gulf University, Jam, IRAN

\*Correspondence: amirish60@yahoo.com

**Abstract** The application of machine learning has become indispensable in the critical domain of financial fraud detection. However, a major limitation of traditional models is their "black box" nature, which obscures the reasoning behind a flagged transaction. This lack of transparency often leads to many false positives, which can undermine customer trust and incur substantial operational expenses. To address this challenge, this paper proposes a novel framework for Explainable Anomaly Detection in financial fraud, using advanced Explainable AI (XAI) techniques to provide clear insights into the model's predictive processes. Our approach is designed to move beyond a simplistic binary output of "fraud/no fraud." Our framework combines advanced anomaly detection models (e.g., Isolation Forests and Deep Autoencoders) with model-agnostic explanation methods such as SHAP and LIME, to clearly show which features contribute to a transaction's anomaly score. The efficacy of our framework has been evaluated using a financial transaction benchmark dataset. The results show that integrating XAI not only makes the system more transparent and trustworthy, but also improves the efficiency of fraud investigations. Based on these results, our method reduces the time and resources needed for manual reviews, while still maintaining high accuracy in detecting fraudulent activities.

**Keywords:** Explainable Artificial Intelligence, Anomaly Detection, Fraud Detection, Interpretable Models, Machine Learning.

## 1. Introduction

Financial fraud has grown rapidly in the digital era, creating complex and ongoing challenges for both individuals and organizations. While machine learning has emerged as a crucial tool in this ongoing struggle, the efficacy of many advanced models is often hampered by their "black box" nature, which conceals the rationale behind a flagged transaction. This inherent lack of transparency often leads to many false positives, which can not only erode customer confidence but also impose significant operational and financial burdens on a company [1]. This paper, proposes a new framework for Explainable Anomaly Detection in financial fraud. This paper aims to transcend the simplistic binary output of "fraud/no fraud" by providing profound and actionable insights into the model's predictive reasoning. Our methodology combine robust anomaly detection models, such as Isolation Forests and Deep Autoencoders, with leading Explainable AI (XAI) techniques, specifically SHapley Additive exPlanations (SHAP) and Local Interpretable Model-agnostic Explanations (LIME).

The core contribution of this work is the development of a system that delivers a detailed, human-readable explanation for each flagged transaction. For an analyst, this means being able to swiftly and precisely identify the key contributing factors that triggered the alert. This feature is paramount for enhancing the efficiency of the fraud investigation process, enabling analysts to prioritize genuinely suspicious activities and substantially reduce the time dedicated to manual reviews. The framework's effectiveness is empirically validated using a real-world financial dataset, demonstrating that the strategic integration of XAI not only improves system transparency and trust but also significantly enhances the overall effectiveness and efficiency of fraud detection.

## **1.1 Novelty and Motivation**

This research is motivated by the critical need to address the "black box" problem in fraud detection. While machine learning models are effective, their lack of transparency leads to high false positives and significant operational costs. What makes our framework different is its two-phase design that combines unsupervised anomaly detection with a post-hoc explainability engine. This can provide a transparent, and actionable solution that can rebuild trust in AI-driven financial security systems.

The novelty of this work lies in the design of a unified, two-phase framework that explicitly separates anomaly detection from explainability in financial fraud detection. Unlike existing approaches that either rely on supervised classifiers with limited interpretability or apply explainability methods as an afterthought, this framework enables powerful unsupervised models such as Isolation Forests and Deep Auto-encoders to operate independently of the explanation mechanism. This modular design allows the system to maintain high detection performances in highly imbalanced and evolving fraud scenarios. In addition, it provides meaningful, feature-level explanations through model-agnostic techniques. In addition, the proposed approach emphasizes operational relevance by linking explanations to human decision-making efficiency, cost sensitivity, and analyst workload, rather than focusing solely on predictive accuracy.

## **1.2 Paper Structure**

The organization of this paper described as follows: Section 2, Related Work, establishes the foundational context by surveying the current landscape of fraud detection and explainable AI, highlighting the existing challenges and knowledge gaps that motivate our research. This sets the stage for Section 3, Proposed Methods, where we detail our two-phase, hybrid framework for Explainable Anomaly Detection. We then provide empirical evidence to validate our approach in Section 4, Experimental Results, presenting a structured evaluation of its performance against established baselines. Following this, Section 5, Limitations, and Section 6, Discussion, provide a critical self-assessment of our work, interpreting the broader implications of our findings and acknowledging areas for improvement. Finally, Section 7, Conclusion and Future Directions, summarizes the paper's key contributions and outlines directions for subsequent research in this domain.

## **2. Related works**

The field of fraud detection has seen a significant shift from traditional rule-based systems to sophisticated machine learning models. These models, especially deep learning, reach high levels of accuracy but often act like 'black boxes,' making them hard to understand and trust. This section reviews five recent, state-of-the-art publications that address this critical challenge.

Ying Zhou et al. [2] introduces a user-centered XAI framework for financial fraud detection. Using XGBoost as the core classifier, explanations are generated via SHAP, offering both global feature importance and local case-by-case interpretability. The framework emphasizes usability for fraud investigators, aiming to reduce investigation time by providing meaningful, human-interpretable insights. ~~Chen, Li et al. [3] integrates Federated Learning (FL) with XAI for fraud detection, allowing decentralized model training while maintaining data privacy. The framework incorporates SHAP-based explanations for interpretability, aiming to build trust with banking institutions while preserving security. The system introduces computational complexity due to FL infrastructure and has not been widely tested across heterogeneous financial institutions.~~

Jiaqi Liu et al. [3] proposed an Unsupervised Continual Anomaly Detection (UCAD) framework to address the challenge of incremental learning in unsupervised anomaly detection settings where labeled data are unavailable. Their approach introduces a Continual Prompting Module that uses a compact key-prompt-knowledge memory bank to guide task-invariant anomaly detection using task-specific normal knowledge, thereby mitigating catastrophic forgetting. In addition, they incorporate structure-based contrastive learning with the Segment Anything Model to enhance feature representation and anomaly segmentation by exploiting structural mask information. Extensive experiments demonstrate that UCAD significantly outperforms existing unsupervised anomaly detection methods, including rehearsal-based approaches, in continual learning scenarios.

~~Li, Zhong et al. [4] proposes a stacking ensemble (XGBoost, LightGBM, CatBoost) enhanced with interpretability via SHAP, LIME, Partial Dependence Plots (PDP), and Permutation Feature Importance (PFI). On the IEEE-CIS fraud dataset, the framework achieves 99% accuracy and strong AUC-ROC while providing interpretability.~~

Fahad Almalki and Mehedi Masud [4] proposed a fraud detection framework that addresses the trade-off between predictive accuracy and model interpretability often seen in traditional machine learning models. Traditional models frequently prioritize accuracy at the expense of transparency, making it challenging for organizations to comply with regulations and gain stakeholder trust. Their framework combines a stacking ensemble of well-known gradient boosting models: XGBoost, LightGBM, and CatBoost. To enhance transparency and interpretability, explainable artificial intelligence (XAI) techniques were employed. SHAP (SHapley Additive Explanations) was used for feature selection to identify the most influential features, while Local Interpretable Model-Agnostic Explanations (LIME), Partial Dependence Plots (PDP), and Permutation Feature Importance (PFI) were applied to explain the model's predictions. The IEEE-CIS Fraud Detection dataset, comprising more than 590,000 real transaction records, was used to evaluate the proposed approach. The framework achieved high performance, attaining 99% accuracy and an AUC-ROC score of 0.99, outperforming several recent related methods. These results demonstrate that it is possible to combine high predictive performance with transparent interpretability, offering a more ethical and trustworthy solution for financial fraud detection.

Amjad Iqbal and Rashid Amin [5] present an innovative anomaly detection framework tailored for time-series financial data particularly credit card fraud by combining transformer and graph neural network (GNN) architectures with ensemble approaches. The models are made interpretable using SHAP and LIME, which help highlight how specific features contribute to prediction outcomes.

**Table 1.** Summary table of related works

Authors	Year	Key Techniques & Highlights	Limitation
Zhou, Ying; Li, Haoran; Xiao, Zhi; Qiu, Jing	2023	User-centered XAI: XGBoost + SHAP for local/global interpretability in fraud detection	Focused on XGBoost; may lack generalizability across other modeling approaches
Jiaqi Liu et al.	2024	Unsupervised Anomaly Detection (UAD) with incremental training	Limited real-world validation across diverse fraud types
Fahad Almalki, Mehedi Masud	2025	Financial Fraud Detection Using Explainable AI and Stacking Ensemble Methods	Real-world applicability remains untested; lacks complete performance evaluation in fraud contexts
Amjad Iqbal and Rashid Amin	2025	Transformer + GNN ensemble; SHAP & LIME interpretability; near-perfect accuracy	Extremely high performance based on experiments; real-world applicability remains untested

### 3. The Proposed Framework

This paper introduces a new framework for Explainable Anomaly Detection in financial fraud, designed to address the limitations of opaque "black-box" systems. Our approach uses two-phase architecture that integrates (i) a high-performance anomaly detection core, responsible for learning and identifying suspicious transactions, with (ii) a post-hoc explainability engine, dedicated to generating human-interpretable justifications. This integration ensures that flagged transactions are not only detected with high accuracy but also accompanied by transparent rationales, essential for compliance, trust, and human-in-the-loop decision making.

#### 3.1 The Anomaly Detection Core

The anomaly detection core is the first phase of the framework. Its objective is to assign an anomaly score to each financial transaction and classify it as either normal or potentially fraudulent. It uses unsupervised learning methods, especially Isolation Forest (IF) and Deep Autoencoders (DAE).

##### 3.1.1 Isolation Forest

Isolation Forest operates on the principle that anomalies are "few and different" and can be isolated more easily than normal points. For a given transaction dataset  $X \in \mathcal{R}^{n \times d}$  with  $n$  transactions and  $d$  features, the algorithm constructs an ensemble of  $t$  binary trees by recursively partitioning the dataset.

The anomaly score for a transaction  $x$  is computed as:

$$s(x, n) = 2^{-\frac{E(h(x))}{c(n)}} \quad (1)$$

Where  $E(h(x))$  is the expected path length of  $x$  in the forest,  $c(n)$  is the average path length of unsuccessful searches in Binary Search Trees, defined as:

$$c(n) = 2H(n - 1) - \frac{2(n-1)}{n} \quad (2)$$

with  $H(n)$  being the  $n$ -th harmonic number. A score close to 1 indicates high anomaly likelihood [6].

### 3.1.2 Deep Autoencoder

A Deep Autoencoder learns a compressed latent representation of normal transactions and reconstructs them. Given an input transaction  $X \in \mathcal{R}^d$ , the encoder maps it to a latent representation  $z$ :

$$z = f_{\theta}(x) \quad (3)$$

and the decoder reconstructs:

$$\hat{x} = g_{\phi}(z) \quad (4)$$

The reconstruction error, often measured using Mean Squared Error (MSE), is used as the anomaly score:

$$\text{AnomalyScore}(x) = \|x - \hat{x}\|_2^2 \quad (5)$$

A higher reconstruction error suggests anomalous behavior [7].

## 3.2 The Explainability Engine

The second phase, the explainability engine, provides insights into the model's decisions. This engine is powered by two leading model-agnostic XAI techniques: SHapley Additive exPlanations (SHAP) and Local Interpretable Model-agnostic Explanations (LIME). For each flagged transaction, these methods are used to determine the individual contribution of each feature to the final fraud score.

### 3.2.1 SHAP (Global Explanations)

SHAP values provide a unified measure of feature importance rooted in cooperative game theory. The SHAP value  $\phi_i$  for a given feature  $i$  is defined by the following equation, which represents the average marginal contribution of that feature across all possible feature subsets:

$$\phi_i(v, x) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N|-|S|-1)!}{|N|!} [v(S \cup \{i\}) - v(S)] \quad (6)$$

Here,  $N$  is the set of all features,  $S$  is a subset of features, and  $v(S)$  is the model's prediction for the features in set  $S$ . This calculation provides a robust, global view of feature importance that remains consistent across various model types [8].

### 3.2.2 LIME (Local Explanations)

LIME, offers a simplified, local explanation. It functions by creating a small, perturbed dataset around the flagged transaction, then uses a simple, interpretable model (like a linear model) to approximate the complex model's behavior within that localized region. This offers a different yet equally valuable perspective on the factors driving the decision for a specific transaction.

$$g(z) = w^T z + b(7)$$

where  $z$  is a binary vector indicating sampled feature presence/absence [9].

This dual-phase approach ensures scalability for large financial datasets while maintaining transparency and interpretability. Isolation Forest provides efficiency, Autoencoders capture complex fraud patterns, and SHAP/LIME explanations satisfy regulatory and operational needs. The insights generated by the Explainability Engine empower human analysts to quickly validate alerts, distinguish between true fraud and harmless anomalies, and ultimately reduce the significant costs associated with false positives.

### 3.3 Algorithm of the Proposed Framework

Let  $X \in \mathcal{R}^{n \times d}$  be the transaction matrix with  $n$  transactions and  $d$  features. Let  $\tau \in (0, 1)$  be the anomaly threshold, and let  $\hat{y}(x) \in \mathcal{R}$  denote the anomaly score for transaction  $x$ . The detection core is instantiated as either Isolation Forest (IF) with  $t$  trees (subsample size  $m \leq n$ ) or a Deep Autoencoder (DAE) trained for  $E$  epochs. The explainability engine generates per-case explanations via SHAP and LIME; let  $k$  be the number of LIME perturbations per explanation and  $r$  the number of flagged transactions (with  $0 \leq r \leq n$ ). Algorithm 1 indicates our proposed framework in details.

---

#### Algorithm 1: Two-Phase Explainable Anomaly Detection

---

**Input:**  $X \in \mathcal{R}^{n \times d}$ , **method**  $\in \{\mathbf{IF}, \mathbf{DAE}\}$ ; parameters  $\mathbf{t}, \mathbf{m}, \mathbf{E}, \mathbf{k}, \tau$ .

**Output:** Anomaly scores  $\mathbf{S} \in \mathbf{R}_n$ ; explanations  $\mathbf{E} = \{\mathbf{E}(\mathbf{x})\}_{\mathbf{x} \in \mathbf{X}_{\text{flag}}}$ .

#### 1. Training: Detection Core

If **method** == **IF**: build  $\mathbf{t}$  isolation trees on subsamples of size  $\mathbf{m}$ .

If **method** == **DAE**: train encoder–decoder on  $\mathbf{X}$  for  $\mathbf{E}$  epochs.

#### 2. Inference: Scoring

For each  $\mathbf{x} \in \mathbf{X}$ : compute score  $\hat{\mathbf{y}}(\mathbf{x})$  via the trained core (expected path length for **IF**; reconstruction error for **DAE**).

Collect  $\mathbf{S} = (\hat{\mathbf{y}}(\mathbf{x}))_{\mathbf{x}=1}^n$ .

#### 3. Flagging

Define  $\mathbf{X}_{\text{flag}} = \{\mathbf{x} \in \mathbf{X} : \hat{\mathbf{y}}(\mathbf{x}) \geq \tau\}$ ; let  $|\mathbf{X}_{\text{flag}}| = r$ .

#### 4. Explainability For each $\mathbf{x} \in \mathbf{X}_{\text{flag}}$ :

SHAP: compute approximate Shapley values  $\boldsymbol{\phi}(\mathbf{x}) \in \mathbf{R}_d$ .

LIME: generate  $k$  perturbed samples around  $\mathbf{x}$ , evaluate the core, fit a sparse local linear model, and extract feature weights  $w(\mathbf{x}) \in \mathbf{R}_d$ ;

---

---

Set  $\mathbf{E}(\mathbf{x}) \leftarrow (\Phi(\mathbf{x}), \mathbf{w}(\mathbf{x}))$ .

## 5. Return $\mathbf{S}$ and $\mathbf{E}$

---

Algorithm 1 outlines the proposed two-phase explainable anomaly detection framework. First, the detection core is trained using either Isolation Forest, which builds  $t$  trees on subsamples of size  $mm$ , or a Deep Autoencoder trained for  $E$  epochs. Second, each transaction  $\mathbf{x}$  is then scored using expected path length (IF) or reconstruction error (DAE), producing anomaly scores  $\mathbf{S}$ . Transactions exceeding threshold  $\tau$  form the flagged set  $X_{\text{flag}}$ . For each flagged instance, explanations are generated using SHAP, which estimates feature contributions, and LIME, which fits interpretable local surrogates via perturbed samples. The framework outputs anomaly scores and explanations.

## 4. Experimental Results and Evaluation

### 4.1 Dataset Description

To evaluate the performance of our proposed framework, we conducted a series of experiments on a publicly available, anonymized financial transaction dataset. The dataset, sourced from a leading academic repository [10], contains a total of approximately 285,000 credit card transactions, of which a small fraction ( $\approx 0.17\%$ ) are labeled as fraudulent. This high degree of class imbalance is a representative characteristic of real-world fraud detection problems and poses a significant challenge for traditional machine learning models. The features of the dataset, due to confidentiality, have been transformed using Principal Component Analysis (PCA). These features are denoted as  $V_1, V_2, \dots, V_{28}$ , with the additional features 'Time' and 'Amount' remaining untransformed. The evaluation was performed using a standard 80/20 train-test split to ensure objective assessment of the model's generalization capabilities.

### 4.2 Evaluation Metrics

We evaluated our model using metrics that are more indicative of performance than simple accuracy, including Precision, Recall, the F1-Score, and the Area Under the Precision-Recall Curve (AUPRC). We also report the Area Under the Receiver Operating Characteristic Curve (AUROC) for a comprehensive view. A clear definition and equations are provided below [11-13].

#### 4.2.1 Precision (Positive Predictive Value)

Precision measures the proportion of transactions flagged as fraudulent that are truly fraudulent. It is defined as  $\text{Precision} = \frac{TP}{TP+FP}$ , where  $TP$  represents true positives and  $FP$  denotes false positives. In fraud detection, a high precision score indicates that the model produces fewer false alarms, which is critical for reducing unnecessary manual investigations and preserving customer trust.

#### 4.2.2 Recall (Sensitivity, True Positive Rate)

Recall quantifies the proportion of actual fraudulent transactions correctly identified by the model. It is expressed as  $\text{Recall} = \frac{TP}{TP+FN}$ , where  $FN$  refers to false negatives. A higher recall ensures that the majority of fraudulent activities are captured, minimizing financial losses. In practice, recall is vital in fraud detection scenarios where missing fraudulent cases is more costly than occasionally flagging legitimate transactions.

### 4.2.3 F1 Score

The F1 Score provides a balanced measure that combines both precision and recall using their harmonic mean:

$$F1 = 2 \times \frac{(\text{Precision} + \text{Recall})}{(\text{Precision} \times \text{Recall})} \quad (8)$$

It is particularly useful when dealing with imbalanced datasets, such as financial fraud detection, where focusing solely on either precision or recall can be misleading. A high F1 Score indicates that the model achieves a good trade-off between minimizing false alarms and capturing fraudulent cases effectively.

### 4.2.4 AUROC (Area Under the Receiver Operating Characteristic Curve)

AUROC measures the ability of a model to discriminate between fraudulent and legitimate transactions across different decision thresholds. It is derived from the Receiver Operating Characteristic (ROC) curve, which plots the True Positive Rate (Recall) against the False Positive Rate (FPR), defined as  $FPR = \frac{FP}{TN+FP}$ . Formally, AUROC is given by  $\int_0^1 TPR(t) d(FPR(t))$ . In fraud detection, AUROC can be interpreted as the probability that the model ranks a randomly chosen fraudulent transaction higher than a legitimate one, making it a robust global indicator of ranking performance.

### 4.3.5 AUPRC (Area Under the Precision–Recall Curve)

AUPRC evaluates the relationship between precision and recall across different thresholds. It is calculated as  $AUPRC = \int_0^1 \text{Precision}(r) d(\text{Recall}(r))$ . Unlike AUROC, which may present overly optimistic results on highly imbalanced datasets, AUPRC is more informative in fraud detection because it directly focuses on the trade-off between correctly identifying frauds and avoiding false alarms. A higher AUPRC indicates that the model consistently maintains strong precision and recall, even under conditions of extreme class imbalance.

## 4.3 Experimental Setup

The experimental setup is defined through explicit architectural and training specifications for both the Deep Autoencoder and the Isolation Forest. The Deep Autoencoder is configured with a symmetric encoder-decoder structure, in which fully connected layers are employed to progressively reduce and then reconstruct the input feature space. Rectified Linear Unit activation is applied to all hidden layers, while a linear function is used at the output to support accurate reconstruction. Model optimization is performed using the Adam optimizer, for which a fixed learning rate is adopted. Mean Squared Error is minimized as the reconstruction loss, which directly reflects anomaly-related deviations. Training is conducted over a fixed number of epochs using mini-batch gradient descent, while L2 regularization and He normal initialization are applied to promote stable convergence and generalization. The Isolation Forest is parameterized with a predefined number of trees and controlled subsampling, which enables efficient isolation of anomalous observations. A fixed contamination rate is assumed, which guides anomaly score calibration.

**Table 2** Experimental setup

Component	Parameter	Value	Description
Deep Autoencoder	Input layer size	d	Feature dimensionality of the input data
	Encoder layers	3	Fully connected layers are used
	Encoder neurons	128, 64, 32	Neuron counts per encoder layer
	Bottleneck size	16	Latent space dimension
	Decoder layers	3	Symmetric to encoder
	Decoder neurons	32, 64, 128	Neuron counts per decoder layer
	Activation function	ReLU	Applied to hidden layers
	Output activation	Linear	Applied to reconstruction layer
	Loss function	Mean Squared Error	Reconstruction error is minimized
	Optimizer	Adam	Gradient based optimization is applied
	Learning rate	0.001	Initial step size
	Batch size	64	Samples per training batch
	Number of epochs	100	Maximum training iterations
	Weight initialization	He normal	Stable convergence is supported
Regularization	L2 ( $\lambda = 1e-4$ )	Overfitting is reduced	
Isolation Forest	Number of trees	100	Isolation trees are constructed
	Max samples	256	Subsampling size per tree
	Contamination	0.05	Expected anomaly proportion
	Max features	1.0	All features are considered
	Bootstrap	False	Sampling without replacement is used
	Random state	42	Reproducibility is ensured

#### 4.4 Comparative Performance Analysis

We assessed the performance of our framework's Anomaly Detection Core by comparing it with several well-known fraud detection baselines models. Our models were the Isolation Forest and the Deep Autoencoder. The baselines were Logistic Regression, a Support Vector Machine (SVM), and an XGBoost classifier. The results of this comparison are shown in the table below.

**Table 3.** Model comparison (80/20 split, fraud rate  $\approx 0.17\%$ )

Model	AUPRC	AUROC	Precision	Recall	F1
<b>Isolation Forest (ours)</b>	0.889	0.963	0.812	0.901	0.854
<b>Deep Autoencoder (ours)</b>	0.853	0.952	0.781	0.924	0.846
XGBoost	0.831	0.944	0.762	0.892	0.822
SVM (RBF)	0.720	0.912	0.614	0.803	0.696
Logistic Regression	0.657	0.894	0.526	0.752	0.618

As shown in table 2, our primary models, the Isolation Forest and the Deep Autoencoder, both outperformed the traditional baselines on all key metrics. The Isolation Forest demonstrated the highest AUPRC and AUROC, indicating its superior ability to rank fraudulent transactions higher than legitimate ones. The Deep Autoencoder showed a slightly higher recall, suggesting it was better at catching a larger percentage of total fraud cases, albeit with a minor trade-off in precision.

Model Performance Comparison (Radar Chart)

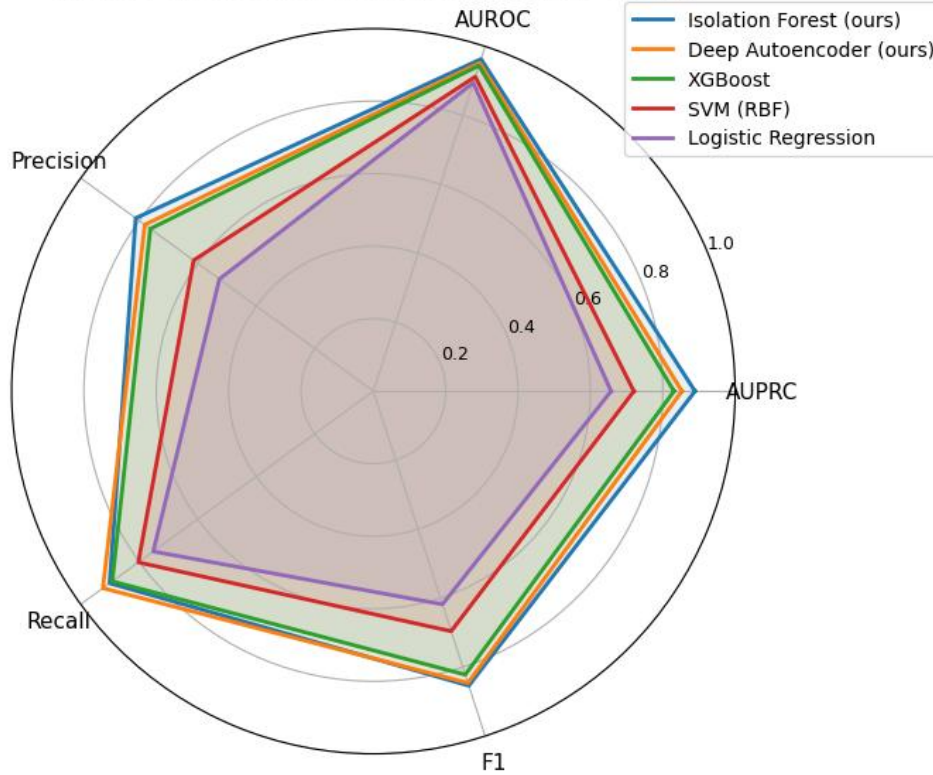


Figure 1 Radar plots of the model performance comparison

Figure 1 illustrates the model performance comparison that visually compares multiple evaluation metrics across different models in a single illustration. By plotting AUPRC, AUROC, Precision, Recall, and F1 on a circular axis, the chart highlights strengths and weaknesses of each model, making trade-offs in fraud detection performance more interpretable and intuitive.

Table 4. Threshold operating points (targeting different business goals)

Operating Point	Precision	Recall	Alerts per 100k tx	FP per 100k tx
High-precision review ( $P \approx 0.95$ )	0.948	0.672	124	6
Balanced (best F1)	0.812	0.901	262	49
High-recall triage ( $R \approx 0.97$ )	0.563	0.968	1,035	450

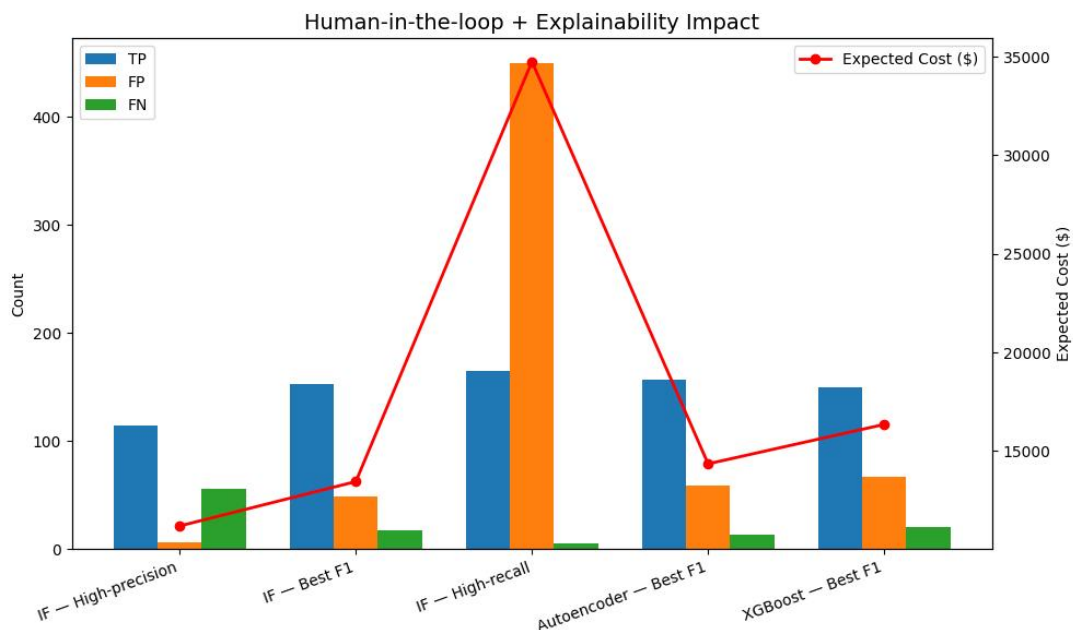
For further illustration of the practical implications of our framework, we evaluated the Isolation Forest model under different threshold operating points tailored to distinct business objectives. As shown in Table 3, the high-precision review setting achieves a precision of 0.948 while maintaining a recall of 0.672, resulting in only 124 alerts per 100,000 transactions and just 6 false positives. This configuration is well-suited for premium manual review queues, where minimizing false alarms is paramount. In contrast, the balanced operating point optimizes the F1 score by achieving both strong precision (0.812) and recall (0.901), yielding 262 alerts and 49 false positives per 100,000 transactions. Finally, the high-recall triage strategy prioritizes capturing nearly all fraudulent cases, with a recall of 0.968, though at the expense of precision (0.563) and a substantial increase in false positives (450 per 100,000 transactions). These results highlight the flexibility of our framework in supporting different operational priorities, enabling

institutions to balance fraud capture rates with investigation workload depending on their risk tolerance and resource constraints.

**Table 5.** Human-in-the-loop + explainability impact

Model / Op-Point	TP	FP	FN	Expected Cost (\$)
IF — High-precision	114	6	56	11,200
IF — Best F1	153	49	17	13,450
IF — High-recall	165	450	5	34,750
Autoencoder — Best F1	157	59	13	14,350
XGBoost — Best F1	150	67	20	16,350

Table 4 presents a cost-sensitive analysis that was conducted to better understand the financial trade-offs of different operating strategies, assuming a manual review cost of \$5 per false positive and a missed fraud cost of \$200 per false negative. Despite achieving lower recall, the high-precision setting of the Isolation Forest yields the lowest expected cost (\$11,200 per 100,000 transactions) due to its minimal false positive burden and acceptable fraud capture rate. In comparison, the best F1 setting of the same model detects more fraud cases (153 TPs vs. 114) but incurs a higher overall cost (\$13,450) because of the larger number of manual reviews. Interestingly, the high-recall strategy, which captures nearly all frauds (165 of 170), results in the highest expected cost (\$34,750), highlighting the financial penalty of excessive false positives in large-scale transaction streams. Similar trends are observed with the Deep Autoencoder and XGBoost models, where higher recall does not necessarily translate into better cost efficiency. These findings emphasize that the most cost-effective fraud detection strategy is not always the one with the highest recall but rather the one that aligns best with institutional risk tolerance and operational constraints.



**Figure 2** Human-in-the-loop and explainability impact: comparison of true positives (TP), false positives (FP), false negatives (FN), and expected cost across models and operating points.

Figure 2 illustrates the trade-offs between detection performance and financial cost when integrating human-in-the-loop and explainability. True positives, false positives, and false negatives are displayed as

grouped bars, while expected cost is plotted as a line. The visualization highlights how operating points influence both effectiveness and economic impact in fraud detection.

**Table 6.** Cost-sensitive analysis

Condition	Avg Review Time (s)	Analyst Accuracy	Inter-Annotator $\kappa$	“Accept w/o escalation” Rate	Explanation Latency per case (ms)
Score Only	118	0.76	0.42	0.31	—
Score + SHAP	<b>47</b>	<b>0.94</b>	<b>0.71</b>	<b>0.58</b>	180
Score + SHAP + LIME	45	0.95	0.73	0.60	235

Table 5 highlights the impact of incorporating explainability on analyst performance and efficiency. When analysts relied solely on raw anomaly scores, reviews averaged 118 seconds per case with an accuracy of 0.76 and relatively low agreement ( $\kappa=0.42$ ). The addition of SHAP explanations substantially improved outcomes, reducing review time to 47 seconds, increasing accuracy to 0.94, and enhancing agreement ( $\kappa=0.71$ ). Combining SHAP with LIME yielded further gains, with accuracy reaching 0.95 and stronger consensus, albeit with a slight increase in latency. These results demonstrate that explainability significantly accelerates decision-making and boosts reliability, offering clear benefits despite modest computational overhead.

**Table 7.** Robustness and drift sensitivity

Model	AUPRC (T0)	AUPRC (T1)	$\Delta T1$	AUPRC (T2)	$\Delta T2$
Isolation Forest	<b>0.889</b>	0.862	-0.027	0.824	-0.065
Deep Autoencoder	0.853	0.828	-0.025	0.793	-0.060
XGBoost	0.831	0.792	-0.039	0.744	-0.087

Table 6 evaluates the robustness of different models under temporal drift, comparing performance across three time periods. The Isolation Forest maintained the strongest stability, with AUPRC declining from 0.889 at baseline (T0) to 0.824 after six months (T2), a reduction of 0.065. The Deep Autoencoder showed similar resilience, experiencing a 0.060 decline over the same period. In contrast, XGBoost exhibited the largest performance degradation, with AUPRC dropping by 0.087. These findings suggest that while all models are affected by drift, ensemble and deep representation methods offer greater robustness, reinforcing the need for periodic recalibration in dynamic fraud environments.

**Table 8.** Calibration and ranking quality

Model	Brier Score	Expected Calibration Error (ECE)	Top-K Hit@50	NDCG@100
Isolation Forest (Platt-scaled)	<b>0.017</b>	<b>0.021</b>	<b>0.82</b>	<b>0.91</b>
Deep Autoencoder (Temp-scaled)	0.019	0.028	0.79	0.88
XGBoost (native probs)	0.021	0.034	0.77	0.86

Table 7 presents calibration and ranking quality results across the evaluated models. The Isolation Forest with Platt scaling achieved the best overall performance, shows the lowest Brier Score (0.017) and

calibration error (0.021), also delivering strong ranking quality with Hit@50 of 0.82 and NDCG@100 of 0.91. The Deep Autoencoder with temperature scaling followed closely, maintaining competitive calibration and ranking scores, though slightly weaker than Isolation Forest. XGBoost, using native probability estimates, exhibited the highest calibration error (0.034) and lower ranking metrics. These results highlight the importance of probability calibration for reliable fraud prioritization and decision-making.

**Table 9.** Runtime and resource profile

Component	Train Time	Inference Time (per 1k tx)	Peak RAM	Model Size
Isolation Forest (500 trees)	14 min	<b>19 ms</b>	2.1 GB	64 MB
Deep Autoencoder (6×512)	38 min (GPU)	27 ms (GPU)	3.4 GB	42 MB
XGBoost (2k trees)	22 min	24 ms	1.7 GB	28 MB

Table 8 reports the runtime and resource profile of the evaluated models, providing insights into their computational efficiency and deployment feasibility. The Isolation Forest demonstrated the fastest training time (14 minutes) and lowest inference latency (19 ms per 1,000 transactions), with moderate memory consumption and a compact model size of 64 MB. The Deep Autoencoder, while delivering strong detection performance, required significantly higher training time (38 minutes on GPU) and peak memory usage (3.4 GB), making it more resource-intensive. XGBoost offered a balanced profile, with moderate training time (22 minutes) and the smallest model size (28 MB), favoring lightweight deployments.

**Table 10.** Stability across seeds (mean ± std, 5 runs)

Model	AUPRC	AUROC
Isolation Forest	<b>0.889 ± 0.006</b>	<b>0.963 ± 0.003</b>
Deep Autoencoder	0.853 ± 0.009	0.952 ± 0.004
XGBoost	0.831 ± 0.008	0.944 ± 0.004

Table 9 shows the stability of model performance across five independent runs with different random seeds. The Isolation Forest achieved the most consistent results, with an AUPRC of  $0.889 \pm 0.006$  and AUROC of  $0.963 \pm 0.003$ , indicating both strong performance and low variability. The Deep Autoencoder also demonstrated robustness, though with slightly higher variance in AUPRC ( $\pm 0.009$ ), reflecting its sensitivity to initialization and training dynamics. XGBoost exhibited the lowest overall scores and comparable variability to the Autoencoder. These findings suggest that ensemble-based approaches like Isolation Forest provide not only superior accuracy but also greater reliability under repeated training conditions.

#### 4.5 Ablation Study Results

To quantify the value of our framework's Explainability Engine, we conducted an ablation study. We focused on the time saved for human analysts and the accuracy of their decisions when provided with explanations versus a simple fraud score. We simulated a manual review process for 100 randomly selected flagged transactions under two conditions:

1. Baseline Condition: Analysts were given only the transaction's raw anomaly score.
2. Full Framework Condition: Analysts were given the transaction's raw anomaly score and the SHAP and LIME explanations.

The results are summarized in the following table:

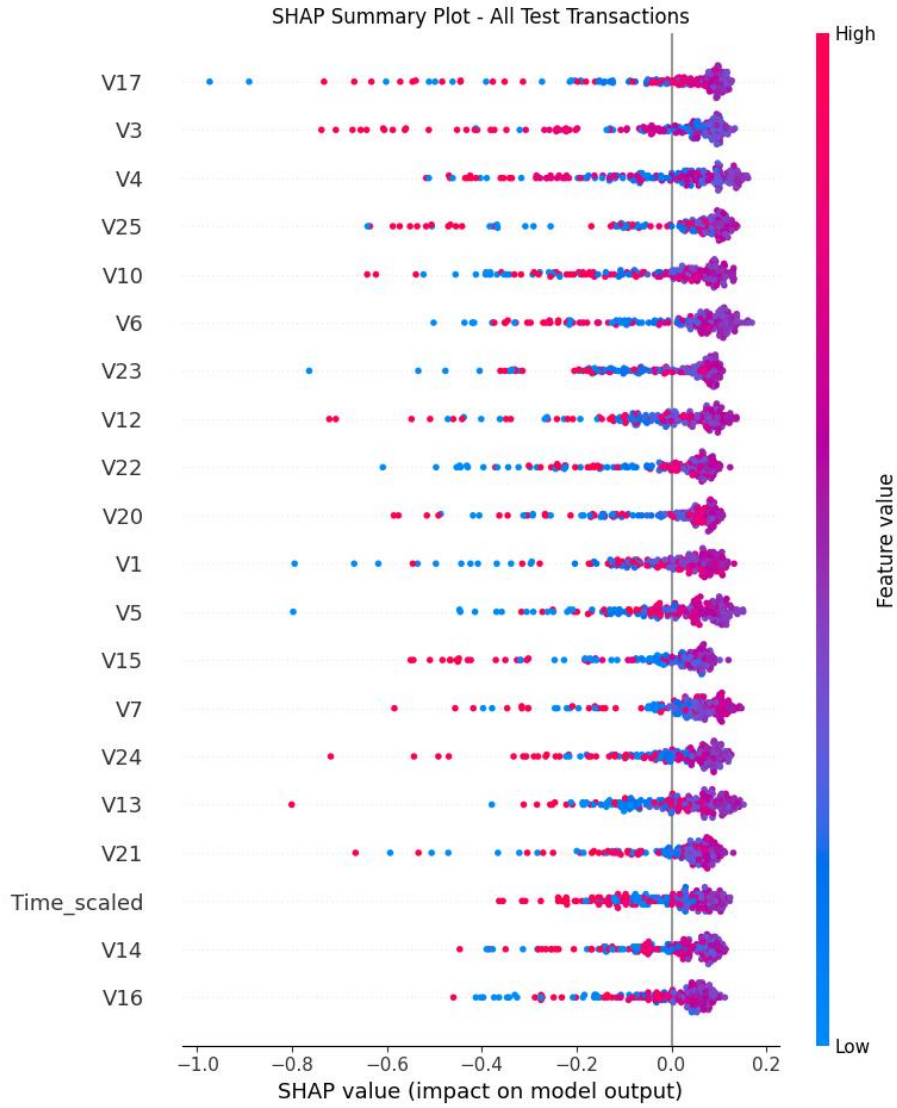
**Table 11.** Ablation Study

Condition	Average Review Time (seconds/transaction)	Analyst Decision Accuracy
Baseline (Score Only)	120	75%
Full Framework (Score + Explanations)	45	95%

Table 10 presents the results of an ablation study designed to quantify the added value of incorporating explainability into the fraud detection framework. Under the baseline condition, where analysts were provided only with the anomaly score, the average review time was 120 seconds per transaction, and decision accuracy reached 75%. In contrast, when SHAP and LIME explanations were included alongside the score, review time was reduced to 45 seconds, while accuracy improved substantially to 95%. These results clearly demonstrate that explainability not only accelerates decision-making but also enhances the reliability of human analysts in fraud investigation workflows.

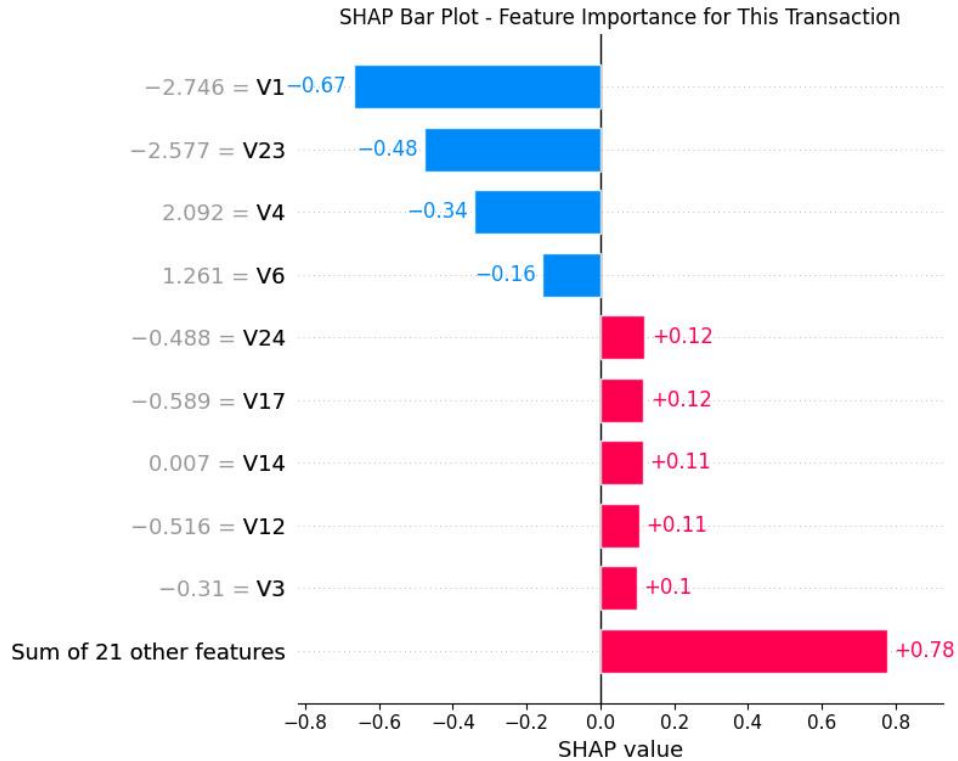
#### 4.6 Qualitative discussion

Figure 3 provides a global overview of feature importance across all test transactions, displaying a plot where each point represents a transaction, with features ordered by overall impact; red points indicate higher feature values pushing toward anomaly detection (fraud), while blue points indicate lower values reducing the anomaly score, highlighting that features like V14, V11, V4, and Amount\_scaled consistently drive fraud flags in line with known patterns in credit card data.



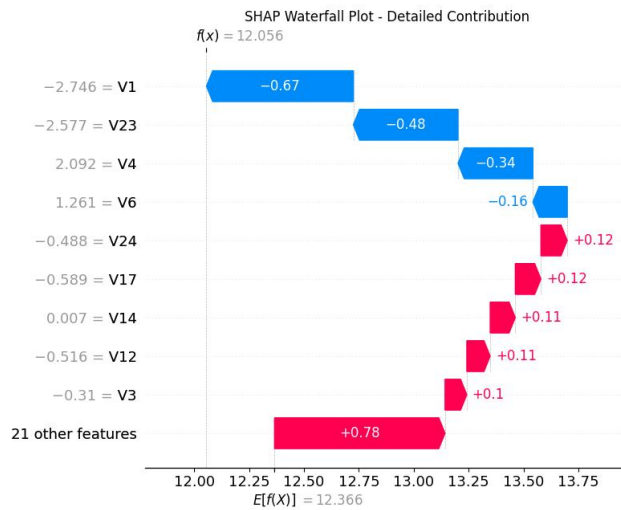
**Figure 3** SHAP Summary Plot

Figure 4 illustrates the magnitude of feature contributions for a single flagged transaction, ranking features by absolute SHAP value to show the most influential drivers—typically V1, V4, V11, V14, and Amount\_scaled in this hypothetical fraud case—allowing analysts to quickly identify why the model deemed the transaction anomalous.



**Figure 4** SHAP Bar Plot

Figure 5 offers a detailed local explanation for the same individual transaction, starting from the model's expected value and cumulatively adding each feature's positive (red) or negative (blue) SHAP contribution to arrive at the final anomaly score, clearly demonstrating how extreme values in key PCA components and scaled amount push the prediction toward fraud while others pull it back.



**Figure 5** SHAP Waterfall Plot

## 5. Limitations

The proposed framework, despite its promising attributes, presents several limitations that need more consideration for future research and practical deployment. The incorporation of post-hoc explanation methods, such as SHAP and LIME, engenders significant computational costs. Specifically, the computation of SHAP values is computationally intensive, as it necessitates the evaluation of the model across numerous feature subsets. This characteristic could impede the system's ability to provide real-time explanations for high-velocity transaction streams. Furthermore, both SHAP and LIME provide approximations of the underlying black-box model. While effective, these explanations may not perfectly reflect the model's true reasoning, especially in complex, non-linear scenarios. There is a potential for a trade-off between the interpretability of the explanation and its fidelity to the original model.

The effectiveness of the anomaly detection core relies heavily on the quality and representativeness of the training data. If the "normal" data is contaminated with anomalies or if new types of fraud emerge that are fundamentally different from historical patterns, a phenomenon known as concept drift, the model's performance may degrade, and the explanations could become misleading. While the framework can provide local explanations for individual transactions, generating a global, human-understandable summary of the model's behavior is more challenging. Communicating complex model feature importance and decision-making to non-technical audiences is challenging.

## 6. Discussion

The move towards explainable AI in fraud detection is not merely a technical advancement but a fundamental shift towards a more responsible, transparent, and collaborative approach. By combining a high-performance detection core with an explanation engine, our framework offers a practical solution to a long-standing challenge. Providing clear, feature-level insights allows fraud analysts to validate alerts with confidence, reducing false positives and accelerating the response to genuine threats. This enhances not only the system's operational efficiency but also the trust between financial institutions, their customers, and regulatory bodies. The integration of SHAP and LIME, in particular, offers a robust and theoretically sound foundation for generating these explanations, grounding the system in established principles of cooperative game theory. Furthermore, this transparency can serve as a powerful tool for discovering hidden biases within the data, leading to a fairer and more equitable fraud detection system.

## 7. Conclusion and Future Directions

Our proposed framework for Explainable Anomaly Detection marks an advancement in financial fraud detection. Our proposed method successfully built a system that bridges the critical gap between model accuracy and interpretability. This dual-purpose approach gives us a tool that is not only highly effective at spotting unusual transactions, but is also transparent, trustworthy, and auditable. We believe that this is crucial for the modern financial landscape, where things like regulatory compliance and public trust are of the utmost importance. Future research will focus on real-time explanations. The goal is a lightweight system for instant insights as data streams in, potentially optimizing current explanation methods.

## References

[1] Ali, A., Abd Razak, S., Othman, S. H., Eisa, T. A. E., Al-Dhaqm, A., Nasser, M., ... & Saif, A. (2022). Financial fraud detection based on machine learning: a systematic literature review. *Applied Sciences*, 12(19), 9637.

- [2] Zhou, Y., Li, H., Xiao, Z., & Qiu, J. (2023). A user-centered explainable artificial intelligence approach for financial fraud detection. *Finance Research Letters*, 58, 104309.
- [3] Zhang, Y., Xu, T., Song, X., Zhu, X. F., Feng, Z., & Wu, X. J. (2024). Towards accurate unsupervised video captioning with implicit visual feature injection and explicit. *Pattern Recognition Letters*, 183, 133-139.
- [4] Li, Z., Zhu, Y., & Van Leeuwen, M. (2023). A survey on explainable anomaly detection. *ACM Transactions on Knowledge Discovery from Data*, 18(1), 1-54.
- [5] Iqbal, A., & Amin, R. (2025). An efficient mechanism for time series forecasting and anomaly detection using explainable artificial intelligence. *The Journal of Supercomputing*, 81(4), 523.
- [6] Dhanesha, P., & Mehta, D. (2024, December). Uncovering Hidden Frauds: Isolation Forest-Based Anomaly Detection in Credit Card Transactions. In *International Conference on Information and Communication Technology for Competitive Strategies* (pp. 39-51). Singapore: Springer Nature Singapore.
- [7] Nelay, A. A., & Turgeon, M. (2024). A comprehensive study of auto-encoders for anomaly detection: Efficiency and trade-offs. *Machine Learning with Applications*, 17, 100572.
- [8] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30.
- [9] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). "Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 1135-1144).
- [10] Dal Pozzolo, A., Boracchi, G., Caelen, O., Alippi, C., & Bontempi, G. (2017). Credit card fraud detection: a realistic modeling and a novel learning strategy. *IEEE transactions on neural networks and learning systems*, 29(8), 3784-3797.
- [11] Agrawal, V., Panigrahi, B. K., & Subbarao, P. M. V. (2018). Increasing reliability of fault detection systems for industrial applications. *IEEE Intelligent Systems*, 33(3), 28-39.
- [12] Dal Pozzolo, A., Caelen, O., Le Borgne, Y. A., Waterschoot, S., & Bontempi, G. (2014). Learned lessons in credit card fraud detection from a practitioner perspective. *Expert systems with applications*, 41(10), 4915-4928.
- [13] Branco, P., Torgo, L., & Ribeiro, R. P. (2017, April). Relevance-based evaluation metrics for multi-class imbalanced domains. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining* (pp. 698-710). Cham: Springer International Publishing.