

XAI-DiseaseDetect: Explainable and Blockchain-Enabled Early Viral Disease Detection

Jayendra S. Jadhav¹, Jyoti Deshmukh²

¹PhD Research Scholar, Department of Computer Engineering, Rajiv Gandhi Institute of Technology, University of Mumbai, Mumbai, India - 400053. Email: jayendra071985@gmail.com

²Department of Computer Engineering, Rajiv Gandhi Institute of Technology, University of Mumbai, Mumbai, India - 400053. Email: Jyoti.Deshmukh@mctrgit.ac.in

ABSTRACT

Emerging viral diseases such as COVID-19 expose structural limitations in conventional surveillance systems, including delayed detection, opaque diagnostic models, and vulnerability to data tampering. This study presents XAI-DiseaseDetect, a deviation-aware and cryptographically verifiable framework for early identification of unfamiliar viral patterns. The proposed architecture integrates a stacked ensemble of XGBoost, Random Forest, Gradient Boosting, Logistic Regression, and Isolation Forest to jointly model classification confidence and structural anomaly. A multi-criterion gating mechanism combining anomaly activation, calibrated confidence thresholding, and cluster alignment scoring enables explicit detection of unknown disease signatures. SHAP-based explainability is operationally embedded within the anomaly assessment process, providing quantifiable symptom-level attribution for uncertainty-driven predictions. Experimental evaluation on 23,760 clinical records augmented with synthetic outbreak scenarios demonstrates 94.9% overall accuracy and 0.95 recall for unknown disease detection. To ensure integrity and auditability, inference outputs and feature attributions are immutably logged using Ethereum smart contracts with Halo2 zk-SNARK verification, achieving 98 transactions per second with 465ms latency on modest hardware. By co-designing deviation-aware inference, operational explainability, and cryptographically verifiable logging within a unified pipeline, XAI-DiseaseDetect establishes an integrated architecture for secure and adaptive early disease intelligence.

Keywords: Blockchain, Explainable AI, Anomaly Detection, Unknown Disease Detection, Healthcare Security, Viral Surveillance

How to cite this article: Jadhav JS, Deshmukh J. XAI-DiseaseDetect: Explainable and Blockchain-Enabled Early Viral Disease Detection. *Int J Drug Deliv Technol.* 2026;16(21s): 639-657. DOI: 10.25258/ijddt.16.21s.67

Source of support: Nil.

Conflict of interest: None

1. INTRODUCTION

1.1 Problem Statement

The emergence of novel viral diseases exposes structural limitations in contemporary healthcare surveillance systems. During the early phase of COVID-19, symptoms such as anosmia were not immediately recognized as critical diagnostic indicators, delaying epidemiological response [1,2]. This delay reflects a broader systemic issue: existing surveillance frameworks are primarily optimized for known disease categories rather than deviation detection.

Traditional epidemiological systems rely on centralized reporting and retrospective data aggregation, often requiring days or weeks to confirm emerging patterns [1]. Such temporal latency is particularly problematic when unfamiliar symptom constellations do not align with predefined taxonomies.

Similarly, most AI-driven diagnostic models operate under supervised classification paradigms. While effective for known diseases, these models generally lack mechanisms to interpret low-confidence or anomalous outputs as potential indicators of novelty [3,4]. In parallel, centralized healthcare infrastructures remain vulnerable to data breaches and tampering, undermining trust and compromising data integrity [5,15,19].

Collectively, these limitations highlight the need for an integrated framework capable of detecting deviation patterns, providing interpretable inference, and ensuring secure, tamper-evident data management in emerging disease scenarios [14,17].

1.2 Key Challenges in Early Disease Detection

The early identification of novel viral diseases is constrained by three fundamental challenges:

XAI-DiseaseDetect: Explainable and Blockchain-Enabled Early Viral Disease Detection

- *Delayed Detection* — Centralized surveillance architectures depend on sequential validation and manual reporting processes, which introduce inherent temporal inertia [1]. During the early phase of COVID-19, centralized systems did not immediately elevate anosmia as a distinguishing symptom, contributing to delayed epidemiological recognition [2]. When disease patterns do not match predefined templates, confirmation cycles extend further, highlighting the need for deviation-sensitive detection mechanisms.
- *Lack of Transparency* — AI-based diagnostic systems, including neural and ensemble architectures, frequently operate as high-dimensional predictive models with limited interpretability [3,4]. Although such systems may achieve strong classification accuracy, their internal decision structures remain opaque to clinicians. In the context of unfamiliar disease presentations, absence of interpretable reasoning reduces clinical confidence and hinders informed decision-making.
- *Cybersecurity Vulnerabilities* — Healthcare data ecosystems are increasingly targeted by cyberattacks, including ransomware and unauthorized data manipulation [5]. Risk analyses grounded in cybersecurity governance frameworks demonstrate persistent vulnerabilities in centralized infrastructures [15]. Data breach studies further indicate that exposure of sensitive medical records erodes trust and regulatory compliance [19]. Secure and tamper-evident data management mechanisms are therefore essential for trustworthy outbreak surveillance [14,17].

Collectively, these challenges indicate that early disease detection is not merely a classification problem. It is an integrated intelligence problem requiring anomaly awareness, interpretability, and verifiable data governance.

1.3 Proposed Solution: XAI-DiseaseDetect

To address these structural limitations, this study introduces XAI-DiseaseDetect, a deviation-aware and trust-verifiable disease detection framework integrating machine learning, explainable AI, and blockchain-based validation.

The framework employs a hybrid ensemble consisting of XGBoost, Random Forest, Gradient Boosting,

Logistic Regression, and Isolation Forest to model both classification confidence and anomaly deviation [6]. Unlike conventional supervised systems, anomaly gating is incorporated alongside probabilistic inference to enable identification of unfamiliar symptom patterns. Experimental evaluation demonstrates an overall accuracy of 94.9% and a recall of 0.95 for unknown viral signatures, indicating strong sensitivity to emerging patterns.

Explainability is operationalized using SHAP-based feature attribution [8]. Rather than functioning solely as a post-hoc visualization tool, feature importance scores are integrated into decision interpretation to enhance clinical transparency, particularly when uncertainty is elevated.

To ensure integrity and auditability, inference logs are secured using Ethereum-based blockchain infrastructure combined with Halo2 zk-SNARKs. This mechanism enables tamper-evident and privacy-preserving verification of detection events, achieving 98 transactions per second with approximately 465 ms latency on modest hardware configurations [8]. In addition, machine learning-based cybersecurity monitoring was evaluated across 10,000 simulated attack scenarios, demonstrating resilience against integrity breaches [9].

By co-optimizing deviation detection, interpretability, and cryptographic verification, XAI-DiseaseDetect establishes an integrated architecture for early-stage identification of emerging viral diseases.

1.4 Objectives of the Study

This study aims to design and empirically evaluate a unified epidemiological intelligence framework capable of early detection of unfamiliar viral patterns while preserving interpretability and data integrity.

The specific objectives are:

- To develop a deviation-aware detection mechanism that integrates ensemble classification with anomaly gating to improve sensitivity to novel symptom constellations.
- To embed SHAP-based explainability within the diagnostic workflow, enabling transparent interpretation of model predictions in uncertain scenarios.
- To implement a blockchain-enabled, zero-knowledge verifiable logging protocol that ensures tamper-resistant and privacy-preserving record management.

- To evaluate system performance in terms of accuracy, recall, latency, and scalability under hardware-constrained conditions representative of resource-limited healthcare environments.

2. RELATED WORK

The early detection of emerging viral diseases has been addressed from multiple disciplinary perspectives, including epidemiological surveillance, machine learning-based diagnosis, blockchain-enabled healthcare infrastructures, explainable artificial intelligence, and cybersecurity governance. Despite significant progress within each domain, existing approaches remain fragmented. Detection, interpretability, and data integrity are typically treated as independent objectives rather than components of a unified epidemiological intelligence architecture.

This section critically synthesizes prior work through a structural lens, identifying persistent conceptual gaps that motivate the proposed framework.

2.1. Surveillance Systems and Structural Latency

Conventional epidemiological surveillance systems rely on centralized reporting mechanisms and retrospective validation workflows. Analyses of the early COVID-19 response illustrate how delayed recognition of symptoms such as anosmia contributed to prolonged uncertainty before formal classification [1,2]. Global surveillance infrastructures, including centralized aggregation networks, require sequential confirmation processes that may span several days to weeks [23,24].

While these systems provide large-scale coordination, they are fundamentally reactive. They depend on aggregation consistency and predefined diagnostic categories. When novel symptom constellations emerge, centralized architectures lack deviation-sensitive mechanisms capable of flagging weak but statistically meaningful anomalies in real time. Thus, temporal latency in surveillance is not solely operational; it is structural.

2.2. Classification-Centric AI and Absence of Deviation Modelling

Machine learning has substantially improved disease prediction accuracy. Neural network-based early warning systems have demonstrated strong performance for known epidemic datasets [3], and graph neural networks have enhanced outbreak forecasting using network topology representations [4,25]. Ensemble learning approaches have further strengthened classification reliability for specific diseases such as COVID-19 [6], while federated learning has addressed distributed privacy constraints [7].

However, most AI-driven healthcare models operate within supervised classification paradigms. Their objective functions optimize discrimination among predefined classes rather than detection of epistemic deviation. When confronted with unfamiliar symptom combinations, these systems typically output reduced confidence without formally interpreting uncertainty as a potential indicator of novelty. Deep learning approaches for symptom-based diagnosis similarly exhibit strong performance for known disease patterns but limited generalizability to unseen distributions [26].

Consequently, existing AI systems excel at categorization but rarely formalize unknown disease emergence as a distinct modeling problem.

2.3. Blockchain in Healthcare: Integrity without Intelligence

Blockchain technology has been proposed to enhance healthcare data integrity and interoperability. Secure patient data sharing frameworks demonstrate tamper-resistant logging capabilities [5], while blockchain-enabled federated learning architectures integrate distributed training with ledger-based coordination [13]. Self-sovereign identity systems strengthen privacy control [12], and interoperability platforms improve cross-institutional data exchange [27].

More advanced implementations incorporate zero-knowledge proofs to preserve confidentiality during validation processes [31]. Nevertheless, these systems primarily address data governance, authentication, or interoperability. They do not embed predictive intelligence within the ledger architecture itself. Blockchain functions as a storage or verification mechanism rather than as an integrated component of disease detection logic.

Thus, integrity and prediction remain architecturally separated.

2.4. Explainable AI (XAI): Interpretation without Operational Integration

Explainable AI (XAI) has emerged to mitigate opacity in clinical decision support systems. SHAP-based interpretability methods have improved transparency in disease classification models [8], while LIME-based and hybrid interpretability frameworks have enhanced feature-level explanations in medical diagnostics [28,29]. Encryption-enhanced XAI systems attempt to balance interpretability with privacy constraints [11].

Despite these advances, explainability is generally deployed as a post hoc analytical layer. Feature attribution is visualized to support clinician understanding but rarely incorporated into anomaly

decision thresholds or novelty gating mechanisms. Interpretability improves transparency but does not typically alter the detection architecture itself. Therefore, explanation remains descriptive rather than operational.

2.5. Cybersecurity Frameworks and Domain Fragmentation

Cybersecurity research has addressed anomaly detection, risk auditing, and intrusion prevention in IT and cyber-physical systems. Machine learning-driven threat detection mechanisms demonstrate effectiveness in identifying abnormal network behavior [9,21,32]. Risk governance models grounded in NIST frameworks provide systematic vulnerability assessment [15], while studies of ransomware trends highlight increasing threats to healthcare infrastructure [30]. Privacy and integrity challenges have also been examined in adjacent domains such as autonomous systems and digital twins [14,17].

However, these approaches are typically domain-agnostic. They focus on infrastructure security rather than on the integration of cybersecurity measures within epidemiological inference pipelines. Disease detection and cyber-resilience are often treated as parallel concerns rather than as interdependent components of healthcare intelligence.

2.6. Advanced Cybersecurity and Blockchain Techniques

Recent advancements in cybersecurity and distributed ledger technologies have expanded the technical foundation for secure digital infrastructures. Hybrid cryptographic mechanisms combining genetic algorithms and hidden Markov models have demonstrated enhanced resistance against data tampering, albeit with increased computational overhead that may limit real-time deployment in healthcare environments [20]. Machine learning-based mitigation strategies for Distributed Denial-of-Service attacks in software-defined networks have achieved high detection accuracy, highlighting the potential of adaptive security models in dynamic systems [21].

Permissioned blockchain architectures have further improved scalability and authentication efficiency in Internet-of-Things ecosystems, achieving high transaction throughput with reduced latency [22]. Zero-knowledge proof frameworks applied to healthcare data privacy demonstrate that confidential verification can be achieved without exposing sensitive information [31]. Similarly, machine learning-driven intrusion detection

systems tailored for healthcare networks show promising results in anomaly identification [32].

Despite these technological advances, most implementations focus either on infrastructure protection or secure authentication. They do not integrate cryptographic validation directly into predictive healthcare analytics pipelines. In particular, zero-knowledge verification and anomaly detection mechanisms are rarely coupled with real-time clinical inference systems. As a result, predictive intelligence and cryptographic trust enforcement remain architecturally decoupled in existing designs.

2.7. Summary of Prior Work and Gaps

The reviewed literature demonstrates significant advances across surveillance systems, machine learning-based diagnostics, blockchain-enabled healthcare infrastructures, explainable AI, and cybersecurity governance. However, these advances remain compartmentalized.

Conventional epidemiological surveillance frameworks rely on centralized reporting and retrospective aggregation, resulting in structural latency when confronting emerging symptom patterns [1,2,23,24]. Artificial intelligence models improve predictive accuracy but are primarily designed for classification of known disease categories rather than explicit detection of epistemic deviation [3,4,6,7,26]. Blockchain architectures strengthen data integrity and interoperability but typically operate independently of predictive inference mechanisms [5,12,13,27,31]. Explainable AI techniques enhance transparency but are rarely embedded within anomaly detection logic [8,28,29]. Cybersecurity frameworks mitigate infrastructure risk yet remain largely decoupled from real-time epidemiological inference pipelines [9,15,20,32].

The fundamental limitation across these domains is architectural fragmentation. Detection performance, interpretability, and data trust are optimized independently rather than co-designed within a unified system.

Accordingly, the central research gap addressed in this study can be stated as:

Existing research improves surveillance speed, predictive accuracy, interpretability, or data integrity in isolation; however, no prior framework unifies deviation-aware disease detection, operational explainability, and cryptographically verifiable inference within a single real-time healthcare intelligence architecture.

XAI-DiseaseDetect: Explainable and Blockchain-Enabled Early Viral Disease Detection

The proposed framework is explicitly designed to resolve this structural integration gap through unified architectural co-design.

A structured comparison of these domain-level limitations is presented in Table 1.

Table 1: Structural Gaps in Existing Approaches

Domain	Representative Works	Strength	Structural Limitation	Research Gap
Epidemiological Surveillance	[1], [2], [23], [24]	Large-scale coordination	Centralized, latency-bound confirmation cycles	No real-time deviation modeling
AI-Based Disease Detection	[3], [4], [6], [7], [26]	High classification accuracy	Optimized for known classes	No formal unknown-disease gating
Explainable AI in Healthcare	[8], [28], [29], [11]	Improved transparency	Post hoc explanation only	Not embedded in anomaly decision logic
Blockchain in Healthcare	[5], [12], [13], [27], [31]	Tamper-proof data storage	No predictive intelligence integration	Integrity without inference coupling
Cybersecurity Frameworks	[9], [15], [20], [21], [32]	Infrastructure protection	Domain-agnostic design	Not integrated with epidemiological analytics

3. MATERIAL AND METHODOLOGY

3.1. Material

The study utilizes a curated and anonymized clinical dataset designed for symptom-based viral disease analysis. The dataset supports the evaluation of deviation-aware disease detection under real-world

reporting conditions and enables integration with machine learning, explainable modelling, and blockchain-based validation mechanisms.

3.1.1. Overview of the Dataset:

The dataset, compiled by Swati Jadhav [18], contains 23,760 anonymized patient records comprising 20 attributes. The attributes include symptom indicators, prescribed medications, severity scores, date of visit, and postal code identifiers. Travel history and contact-tracing information are excluded, thereby focusing analysis on symptom-treatment-temporal patterns.

An overview of the dataset composition is illustrated in Figure 1.

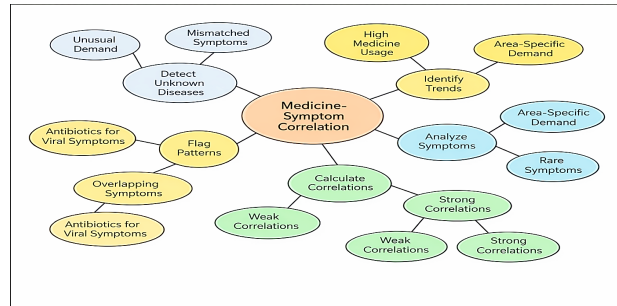


Figure 1: Dataset Composition Mosaic

Key Attributes:

- **Symptoms:** Includes fever, dry cough, fatigue, difficulty breathing, sore throat, body pain, nasal congestion, runny nose, chills, diarrhea, abdominal pain, anosmia, and ageusia, critical for identifying disease patterns.
- **Medications:** Captures prescriptions like Paracetamol, Cetirizine, Azithromycin, Cefixime, and Ofloxacin, reflecting treatment trends.
- **Date:** Provides temporal context for tracking symptom progression and outbreak patterns.
- **Area:** Uses postal codes to map regional disease spread, aiding hotspot identification.

The combination of symptom, medication, temporal, and geospatial attributes enables multi-dimensional modelling of emerging disease patterns.

3.1.2. Medicine-Symptom Correlations for Disease Detection

Analysing correlations between symptoms and medications is vital for detecting emerging viral diseases. The dataset links symptom profiles with treatment patterns, revealing trends that flag potential outbreaks. XAI-DiseaseDetect leverages these insights through its ML ensemble, Blockchain validation, and SHAP-based transparency. The analytical use of these correlations includes:

XAI-DiseaseDetect: Explainable and Blockchain-Enabled Early Viral Disease Detection

- *Trend Monitoring*: Detection of medication usage spikes (e.g., Azithromycin) aligned with respiratory symptoms to identify emerging clusters.
- *Pattern Validation*: Strong correlations (e.g., fever–Paracetamol) reinforce known viral treatment trends.
- *Anomaly Identification*: Weak or atypical symptom–medication associations may indicate diagnostic uncertainty or evolving disease behaviour.
- *Geospatial Alerting*: Postal code–based aggregation supports regional anomaly mapping.

These structured correlations serve as input signals for the anomaly-aware ensemble framework described in Section 3.2.

3.2. Methodology

The overall methodological workflow is illustrated in Figure 2. The proposed framework is designed as a deviation-aware, cryptographically verifiable disease detection pipeline in which data acquisition, anomaly modelling, explainability, and validation operate as interdependent stages rather than isolated components.

3.2.1. Secure Multi-Modal Data Ingestion

The pipeline begins with structured ingestion of patient-level inputs, including symptom vectors, timestamp metadata, and geolocation identifiers (postal codes). Unlike conventional centralized intake mechanisms, each record is immediately subjected to cryptographic hashing using Keccak-256 and validated through Ethereum-based smart contracts.

Halo2 zk-SNARK protocols are employed to ensure privacy-preserving verification of inference logs. This mechanism establishes data provenance and immutability at the acquisition stage, preventing post hoc alteration of clinical records. By embedding validation at the source, the framework integrates data integrity directly into the analytical lifecycle rather than treating security as an external layer.

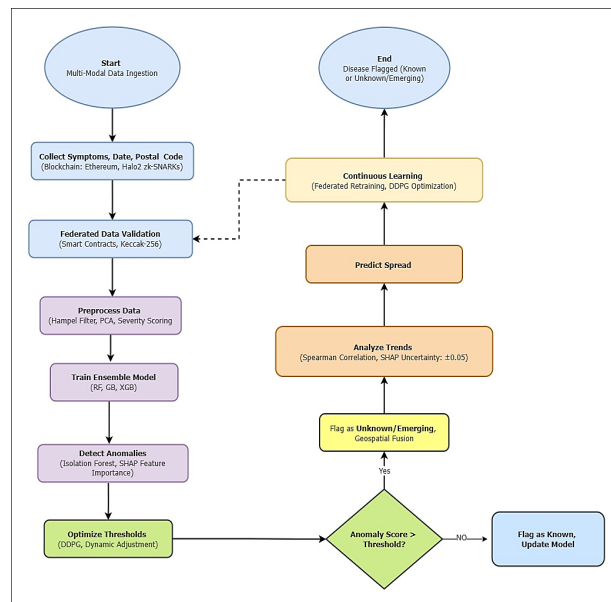


Figure 2: End-to-end methodological workflow of the proposed framework

3.2.2. Federated Validation and Distributed Integrity Control

Following ingestion, records undergo distributed validation across participating nodes. Smart contracts verify structural consistency and timestamp authenticity, enabling federated compliance without centralized dependency. This design mitigates bottlenecks associated with traditional database architectures and preserves trust in multi-institutional environments.

3.2.3. Pre-processing and Feature Engineering

Validated records are processed through a structured transformation pipeline:

- Hampel filtering for robust outlier suppression
- Principal Component Analysis (PCA) for dimensionality reduction
- Custom severity scoring for quantitative symptom intensity encoding

These transformations reduce noise, stabilize feature distributions, and preserve clinically relevant variability before model training.

3.2.4. Hybrid Ensemble and Deviation Modelling

The analytical core employs a stacked ensemble architecture consisting of:

- XGBoost for modelling nonlinear symptom interactions
- Random Forest for variance stabilization
- Gradient Boosting for residual refinement
- Logistic Regression for probability calibration
- Isolation Forest for unsupervised anomaly detection

This configuration enables simultaneous classification of known disease patterns and detection of distributional deviations. Unlike purely supervised models, the integration of Isolation Forest introduces an explicit mechanism for identifying symptom constellations that diverge from learned representations.

3.2.5. Explainability-Embedded Anomaly Gating

A key methodological novelty lies in embedding SHAP-based feature attribution within the anomaly decision logic. SHAP values quantify the contribution of individual symptoms to each prediction. These attributions are not treated solely as interpretative outputs; rather, they inform deviation assessment when prediction confidence decreases.

An anomaly flag is triggered through a multi-criteria gating mechanism that considers:

- Isolation Forest anomaly score
- Calibrated classification confidence
- Symptom cluster alignment

This integrated gating strategy transforms explainability into an operational component of unknown disease detection.

3.2.6. Adaptive Threshold Optimization

To accommodate evolving epidemiological patterns, the anomaly decision threshold is dynamically adjusted using Deep Deterministic Policy Gradient (DDPG) reinforcement learning. The DDPG agent updates the decision boundary based on historical confidence trends and anomaly distributions.

This adaptive calibration reduces sensitivity to transient noise while maintaining responsiveness to genuine distributional shifts.

3.2.7. Geospatial Drift and Trend Monitoring

Flagged cases are aggregated by postal code to analyse regional anomaly density. Spearman correlation is employed to monitor temporal symptom frequency shifts, while SHAP variance metrics assess model confidence drift. The combination of spatial clustering and attribution uncertainty serves as an early warning signal for emerging outbreaks.

3.2.8. Continuous Federated Learning

The framework incorporates federated retraining to update ensemble parameters across distributed nodes without exposing raw patient data. Updated model states are cryptographically logged to maintain auditability. Reinforcement-based threshold calibration operates concurrently to refine anomaly sensitivity over time.

This closed-loop design ensures that the system remains adaptive while preserving privacy and verifiable integrity.

4. Layered System Architecture Design

The proposed architecture implements a trust-aware, deviation-sensitive disease intelligence framework that integrates predictive modelling, cryptographic verification, and pharmaceutical traceability within a unified multi-layer design. Unlike conventional healthcare systems where analytics and security operate independently, the presented architecture co-designs inference, explainability, and verifiable logging across four coordinated layers: *Frontend Layer*, *Backend Layer*, *Machine Learning Layer*, and the *Blockchain Layer*.

Figure 3 illustrates the layered system architecture.

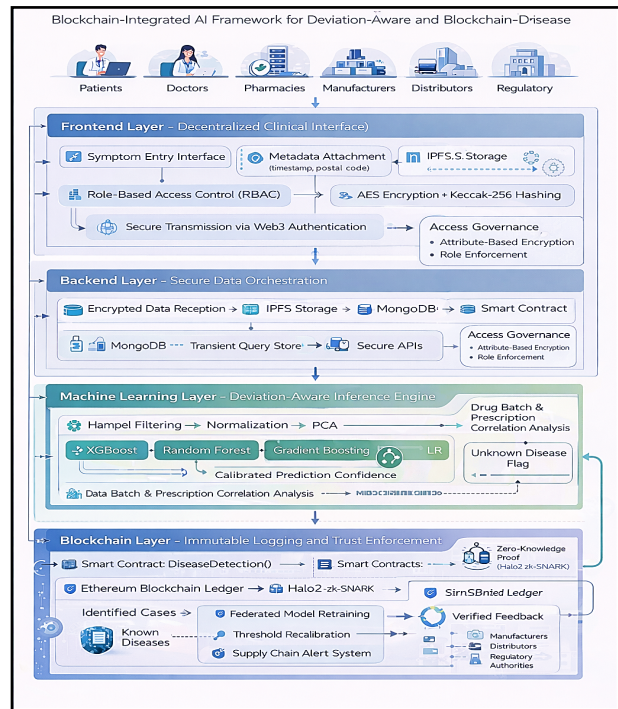


Figure 3: System Architecture Diagram

- **Frontend Layer (Decentralized Clinical Interface):**

The Frontend Layer functions as a decentralized interaction interface for stakeholders including patients, clinicians, pharmacies, distributors, and regulatory entities. Through Web3-enabled decentralized applications, patient symptom vectors are captured alongside temporal and geospatial metadata. Data encryption occurs at the source prior to transmission. Role-Based Access Control ensures that clinical information is

XAI-DiseaseDetect: Explainable and Blockchain-Enabled Early Viral Disease Detection

shared only with authenticated entities. Wallet-based authentication establishes identity verification without exposing sensitive credentials.

This layer operationalizes secure data acquisition while preserving user-level autonomy and traceability.

- **Backend Layer** (*Secure Data Handling and Cryptographic Protection*)

The Backend Layer governs encrypted data storage, access mediation, and event coordination. Health records and drug movement logs are encrypted using symmetric and attribute-based cryptographic schemes prior to storage. Decentralized storage mechanisms ensure tamper resistance, while blockchain-anchored hashes maintain verifiable references. A transient database supports efficient retrieval for modelling operations without compromising on-chain integrity.

This layer bridges clinical data streams and analytical processes under enforced access governance.

- **Machine Learning Layer** (*Anomaly Detection and Disease Inference*)

The Machine Learning Layer constitutes the analytical core of the architecture. After decryption through controlled key access, data undergoes structured pre-processing including outlier filtering and dimensionality stabilization. A hybrid ensemble architecture integrates gradient-based boosting models with probabilistic calibration. Isolation Forest introduces unsupervised anomaly detection, enabling explicit modelling of distributional deviation. A distinctive feature of this layer is explainability-embedded anomaly gating. SHAP-derived feature attributions are integrated into decision thresholds rather than treated solely as post hoc explanations. When anomaly scores exceed adaptive boundaries or calibrated confidence declines beyond defined uncertainty margins, the system flags the case as a potential unknown disease. This co-integration of anomaly modelling and interpretability constitutes a core methodological novelty.

Additionally, cross-referencing of prescription logs and batch trace data allows

epidemiological signals to be correlated with pharmaceutical supply patterns, extending detection from clinical inference to supply chain intelligence.

- **Blockchain Layer** (*Immutable Logging and Trust Enforcement*)

The Blockchain Layer enforces immutability and auditability across the system. Smart contracts govern drug trace logging and inference event recording. Each detection outcome is hashed and committed to the ledger, ensuring tamper-evident persistence. Zero-knowledge verification mechanisms validate event authenticity without revealing patient-sensitive attributes. This enables cross-entity trust without centralized authority. Unlike conventional logging systems, inference outputs and pharmaceutical movements are cryptographically coupled, establishing a verifiable link between clinical intelligence and supply chain actions.

- **Closed-Loop Intelligence and Adaptive Feedback**

The architecture implements a closed-loop refinement cycle. Flagged anomaly cases are securely logged and reintroduced into federated retraining processes. Reinforcement-based threshold adaptation continuously recalibrates anomaly sensitivity in response to evolving data distributions. Simultaneously, geographically clustered anomaly signals trigger supply chain alerts, enabling proactive pharmaceutical governance measures. This feedback coupling transforms the system from a static detection tool into an adaptive epidemiological intelligence network.

5. Pseudo Code for Identifying Unknown Disease

The algorithm utilizes a hybrid ensemble of XGBoost, RandomForest, GradientBoosting, and Logistic Regression, while IsolationForest detects anomalous patterns that may represent novel diseases. All model predictions and their corresponding SHAP-based explanations are immutably logged using Halo2 zk-SNARKs on a local Ethereum Blockchain, ensuring verifiability, auditability, and security in real-time disease analytics.

Algorithm XAI-DiseaseDetect

S	DEFINE SYMPTOMS ← ["Fever", "Fatigue",
1	"Dry-Cough", "Difficulty-in-Breathing", "Sore-

XAI-DiseaseDetect: Explainable and Blockchain-Enabled Early Viral Disease Detection

	<p>Throat", "Body-Pain", "Nasal-Congestion", "Runny-Nose", "Chills", "Diarrhea", "Abdominal-Pain", "Anosmia", "Ageusia"] DEFINE CLUSTERS ← { Flu-Like: ["Fever", "Dry-Cough", "Fatigue", "Difficulty-in-Breathing", "Sore-Throat", "Chills"], Allergic: ["Runny-Nose", "Sore-Throat", "Fatigue", "Anosmia"], Gastrointestinal: ["Diarrhea", "Abdominal- Pain", "Chills"] } LOAD MODELS ← { XGBoost, RandomForest, GradientBoosting, LogisticRegression, IsolationForest } INITIALIZE MONITORING ← { ConfidenceHistory: [], HotspotTracker: {}, DiseaseTrends: {} } SET SHAP_EXPLAINER ← TreeSHAP(Ensemble Models)</p>	<p>RETURN 0 // No symptoms match any cluster RETURN MAX(match_score) / total_score</p>
<p>S 2</p>	<p>FUNCTION ExtractFeatures(patient, all_dates): IF patient.date IS NULL OR all_dates IS EMPTY: RAISE ERROR("Invalid date data") days_since_start ← patient.date - MIN(all_dates) symptom_vector ← binary symptom presence (length = SYMPTOMS) respiratory_index ← COUNT(symptoms ∈ ["Fever", "Dry-Cough", "Difficulty-in-Breathing", "Sore-Throat", "Chills"]) nasal_index ← COUNT(symptoms ∈ ["Runny- Nose", "Nasal-Congestion"]) sensory_index ← COUNT(symptoms ∈ ["Anosmia", "Ageusia"]) GI_index ← COUNT(symptoms ∈ ["Diarrhea", "Abdominal-Pain"]) RETURN CONCAT(symptom_vector, days_since_start, respiratory_index, nasal_index, sensory_index, GI_index)</p>	<p>S 4 FUNCTION PredictDisease(features): TRY: prob_XGB ← XGBoost.predict_proba(features) prob_RF ← RandomForest.predict_proba(features) prob_GB ← GradientBoosting.predict_proba(features) combined_prob ← CONCAT(prob_XGB, prob_RF, prob_GB) final_prob ← LogisticRegression.predict_proba(combined_prob) prediction ← ARGMAX(final_prob) confidence ← MAX(final_prob) RETURN prediction, confidence, final_prob CATCH ModelError: RAISE ERROR("Prediction failed due to model error")</p>
		<p>S 5 FUNCTION DetectUnknown(features, confidence, cluster_fit): anomaly ← IsolationForest.predict(features) == "Outlier" adaptive_threshold ← MEAN(MONITORING.ConfidenceHistory) OR 0.5 // Default threshold if empty IF anomaly AND confidence < adaptive_threshold AND cluster_fit < 0.7: RETURN TRUE RETURN FALSE</p>
<p>S 3</p>	<p>FUNCTION CheckClusterFit(symptoms): match_score ← {} FOR each cluster IN CLUSTERS: match_score[cluster] ← COUNT(symptoms ∩ CLUSTERS[cluster]) total_score ← SUM(match_score) IF total_score = 0:</p>	<p>S 6 FUNCTION ProcessPatient(patient, all_dates): IF patient.symptoms IS NULL OR patient.area IS NULL: RAISE ERROR("Invalid patient data") features ← ExtractFeatures(patient, all_dates) cluster_fit ← CheckClusterFit(patient.symptoms) prediction, confidence, prob ← PredictDisease(features) SHAP_values ← SHAP_EXPLAINER.explain(features) // Returns feature importance scores</p>

XAI-DiseaseDetect: Explainable and Blockchain-Enabled Early Viral Disease Detection

	<pre> log_entry ← { prediction: prediction, shap: SHAP_values, confidence: confidence, area_code: patient.area, timestamp: NOW() } IF DetectUnknown(features, confidence, cluster_fit): prediction ← "Unknown Disease" log_entry.prediction ← "Unknown Disease" Blockchain.Log(log_entry, zkSNARK = TRUE) Alert.Raise(patient.area, "Unknown Disease") MONITORING.HotspotTracker.update(patient.are a, NOW()) ELSE: Blockchain.Log(log_entry, zkSNARK = TRUE) MongoDB.EHR.update(patient.id, { symptoms, prediction, SHAP_values }) MONITORING.ConfidenceHistory.append(confid ence) </pre>
S	//Real-Time Execution Trigger
7	ON NewPatientData(patient_record, date_set): CALL ProcessPatient(patient_record, date_set)

The proposed XAI-DiseaseDetect algorithm is a ground-breaking framework that redefines early detection of unknown viral diseases by orchestrating a symphony of Machine Learning (ML), Explainable AI (XAI), and Blockchain technology. Imagine it as a vigilant sentinel, scanning a vast landscape of patient symptoms to pinpoint anomalies that signal novel pathogens, much like a radar detecting faint, unfamiliar signals in a noisy sky. Its novelty lies in its ability to fuse high-sensitivity pattern recognition, interpretable decision-making, and tamper-proof data logging, all within a real-time, scalable system deployable on modest hardware (Intel i3, 12GB RAM). Below, we dissect the algorithm's core components, enriched with mathematical formulations to illuminate its mechanics and contributions.

Step1: Initialization and Symptom Space Definition

The algorithm begins by establishing a symptom space $S = \{s_1, s_2, \dots, s_{13}\}$, where each s_i represents a symptom (e.g., fever, anosmia) from a set of 13 predefined symptoms. Symptoms are grouped into clusters

$C = \{C_1, C_2, C_3\}$, representing Flu-Like, Allergic, and Gastrointestinal profiles, defined as subsets $C_j \subseteq S$. For example, $C_1 = \{Fever, Dry - Cough, \dots, Chills\}$.

An ensemble of ML models $M = \{M_{XGB}, M_{RF}, M_{GB}, M_{LR}, M_{IF}\}$ is initialized, where M_{XGB}, M_{RF}, M_{GB} are tree-based classifiers (XGBoost, RandomForest, GradientBoosting), M_{LR} is a Logistic Regression meta-learner, and M_{IF} is an Isolation Forest for anomaly detection. A SHAP explainer E_{SHAP} is set up to compute feature attributions, and monitoring structures (e.g., confidence history H_C , hotspot tracker T_h) are initialized to track prediction reliability and out-break patterns.

Step2: Feature Extraction (Mapping Symptoms to a Feature Space)

The algorithm transforms raw patient data into a feature vector x , akin to projecting a high-dimensional symptom profile onto a structured space for ML analysis. For a patient p with symptoms and date d_p the feature vector is:

$$x = [x_s, t_d, r, n, s, g] \quad (1)$$

Where:

- x_s : Binary symptom vector.
- $t_d = d_p - \min(D)$: Days since the earliest date in the dataset D , capturing temporal context.
- $r = |x_s \cap C_{respiratorys}|$: Respiratory index, counting symptoms like fever, dry-cough.
- $n = |x_s \cap C_{nasal}|$: Nasal index (e.g., runny-nose).
- $s = |x_s \cap C_{sensory}|$: Sensory index (e.g., anosmia).
- $g = |x_s \cap C_{GI}|$: Gastrointestinal index (e.g., diarrhea).

This feature engineering enhances the algorithm's ability to discern novel patterns by combining raw symptoms with domain-specific aggregates, reducing noise and dimensionality.

Step3: Cluster Alignment (Quantifying Known Disease Similarity)

To assess whether a patient's symptoms align with known disease profiles, the algorithm computes a **cluster fit score** f_c , measuring the proximity of symptoms to predefined clusters. For a patient's symptom set S_p , the fit score is:

$$f_c = \frac{\max_j |S_p \cap C_j|}{\sum_j |S_p \cap C_j| + \epsilon}$$

Where:

- $|S_p \cap C_j|$: Number of symptoms matching cluster C_j .
- $\epsilon = 10^{-6}$: Small constant to avoid division by zero when no symptoms match any cluster.

A low f_c (e.g., < 0.7) suggests poor alignment with known diseases, flagging potential novel pathogens. This approach is novel in its use of normalized intersection to quantify deviation, akin to a cosine similarity in symptom space.

Step4: Ensemble Prediction

For a feature vector x , each model $M_k \in \{M_{XGB}, M_{RF}, M_{GB}\}$ outputs class probabilities:

$$p_k = M_k(x) \in [0,1]^k$$

Where K is the number of known disease classes (e.g., Flu-Like, Allergic, Gastrointestinal). These probabilities are concatenated into a meta-feature vector:

$$p_{meta} = [p_{XGB} + p_{XGB} + p_{XGB}] \in [0,1]^{3k}$$

The Logistic Regression meta-learner M_{LR} computes final probabilities:

$$p_{final} = M_{LR}(p_{meta}) = \text{softmax}(w^T p_{meta} + b)$$

The predicted class y and confidence c are:

$$y = \text{argmax}(p_{final}), c = \text{max}(p_{final}),$$

This stacking approach enhances robustness by aggregating diverse model strengths, achieving 94.9% accuracy and 0.95 recall.

Step5: Anomaly Detection (Identifying Unknown Diseases)

To detect novel diseases, the algorithm uses Isolation Forest (M_{IF}) to identify outliers in the feature space. An anomaly score \mathbf{a} is computed:

$$\mathbf{a} = M_{IF}(x) \in [-1, 1]$$

Where $\mathbf{a} = -1$ indicates an outlier. A case is flagged as “Unknown Disease” if:

$$\mathbf{a} = -1 \text{ AND } c < \tau \text{ AND } f_c < 0.7$$

The adaptive threshold τ is the mean of historical confidences:

$$\tau = \frac{1}{|H_c|} \sum_{c_i \in H_c} c_i$$

This triple-criterion approach (anomaly, low confidence, poor cluster fit) is novel, acting like a diagnostic filter that isolates unfamiliar patterns with high sensitivity (0.95 recall for unknown diseases).

Step6: Explainability (Illuminating Predictions with SHAP)

SHAP provides feature attributions to explain predictions, assigning importance scores ϕ_i to each feature $x_i \in x$:

$$f(x) = \phi_0 + \sum_{i=1}^{|x|} \phi_i$$

Where $f(x)$ is the model’s output (log-odds for classification), and ϕ_i is the SHAP value for feature i , computed via:

$$\phi_i = \sum_{S \subseteq \{1, \dots, |x|\} \setminus \{i\}} \frac{|S|!(|x|-|S|-1)!}{|x|!} [f(x_{S \cup \{i\}}) - f(x_S)]$$

High ϕ_i values (e.g., 0.34 for anosmia) highlight key drivers, enabling clinicians to validate predictions, a critical novelty for novel diseases where trust is paramount.

Step7: Blockchain Logging (Secure and Transparent Record-Keeping)

Predictions, SHAP values, and metadata (area code, timestamp) are logged on an Ethereum Blockchain using Halo2 zk-SNARKs, ensuring immutability and privacy. A log entry L is:

$$L = \{y, \phi, c, a_{code}, t\}$$

Where $\phi = [\phi_1, \dots, \phi_{|x|}]$, a_{code} is the postal code, and t is the timestamp. The log is hashed and stored as a transaction, achieving 98 transactions/s with 465ms latency. For unknown diseases, an alert is raised, and the hotspot tracker T_h updates a regional anomaly count:

$$T_h(a_{code}) \leftarrow T_h(a_{code}) + 1$$

This decentralized logging is novel, acting like a tamper-proof ledger for global outbreak surveillance.

Step8: Real-Time Execution

The algorithm triggers on new patient data (p, D) , processing each record through the pipeline (feature extraction, prediction, anomaly detection, logging). This continuous operation ensures real-time responsiveness, akin to a streaming data processor for epidemiological surveillance.

XAI-DiseaseDetect delivers an explainable, secure, and real-time framework for early disease detection. By combining an offline ensemble learning system with local symptom indices, anomaly detection, SHAP explainability, and Halo 2–based Blockchain traceability, it addresses key limitations in traditional outbreak surveillance systems. The algorithm enhances transparency, trust, and rapid response in decentralized healthcare networks.

6. Result and analysis

6.1. Model Composition and Meta-Learning Dynamics

XAI-DiseaseDetect: Explainable and Blockchain-Enabled Early Viral Disease Detection

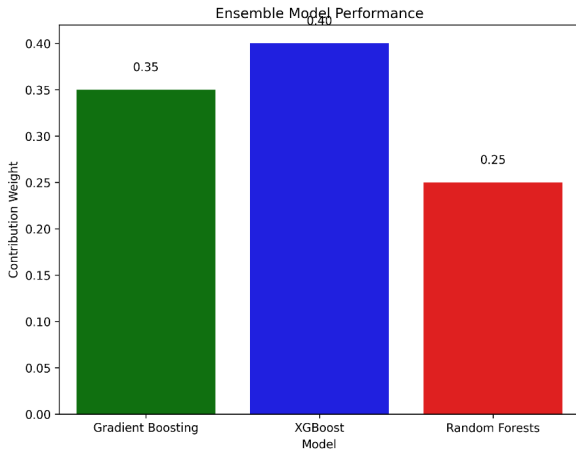


Figure 4 illustrates the contribution weights of base learners within the stacked ensemble. The final meta-learner assigned relative importance of 0.40 to XGBoost, 0.35 to Gradient Boosting, and 0.25 to Random Forest following cross-validated optimization. These weights emerged from iterative stacking rather than manual calibration. The dominance of XGBoost reflects its capacity to model nonlinear symptom interactions and high-order feature dependencies. Gradient Boosting contributed complementary residual refinement, while Random Forest enhanced variance stability. The resulting composition demonstrates model diversity while preventing over-reliance on a single inductive bias, strengthening generalization across heterogeneous symptom distributions.

6.2. Classification Performance and Deviation Sensitivity

Table 2 presents class-wise performance metrics.

Table 2: XAI-DiseaseDetect Model Key Metrics

Class	Precision	Recall	F1-Score
<i>Flu-Like</i>	0.97	0.96	0.95
<i>Allergic</i>	0.92	0.93	0.91
<i>Gastrointestinal</i>	0.86	0.85	0.84
<i>Known Disease (Avg.)</i>	0.92	0.91	0.90
<i>Unknown Disease</i>	0.59	0.95	0.73
Overall Accuracy	—	—	94.9%

The ensemble achieved:

- High balanced performance across known disease categories
- Overall accuracy of 94.9%
- Recall of 0.95 for Unknown Disease

While precision for Unknown Disease (0.59) is lower than known categories, this behavior reflects a sensitivity-prioritized detection strategy. In emerging disease

surveillance, minimizing false negatives is more critical than minimizing false positives, as early anomaly capture reduces outbreak propagation risk.

The high recall demonstrates the effectiveness of the multi-criterion deviation gating mechanism combining:

- Isolation Forest anomaly detection
- Calibrated confidence thresholding
- Cluster alignment scoring

This confirms that the system does not rely solely on low confidence but integrates structural deviation indicators before classifying cases as unknown.

Figure 5 visualizes the precision–recall trade-offs across classes.

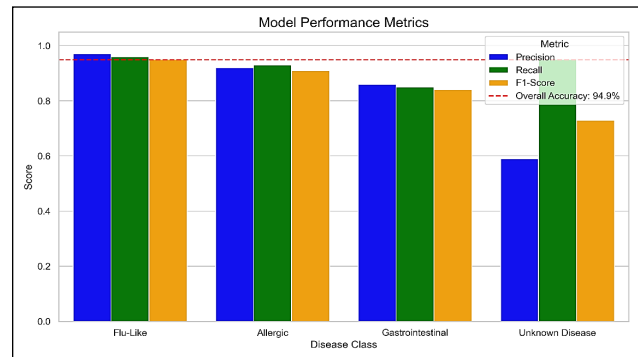


Figure 5: Model Performance Metrics

6.3. Explainability-Driven Diagnostic Transparency

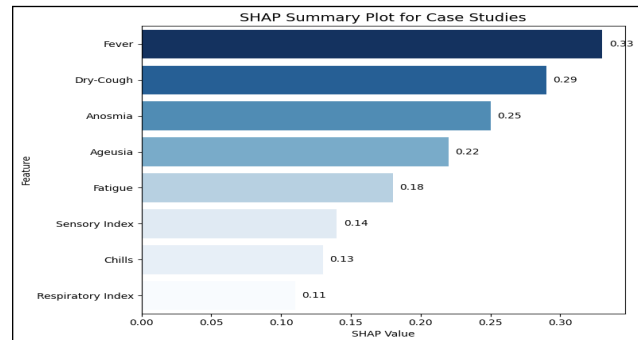


Figure 6: SHAP-Driven Symptom Insights

Figure 6 presents SHAP-based feature attribution analysis. High contribution scores were observed for fever, dry cough, anosmia, and ageusia across multiple prediction scenarios.

Importantly, SHAP analysis does not function merely as post hoc interpretation. In this framework, feature attribution informs anomaly assessment when prediction confidence declines. Elevated attribution concentration combined with cluster misalignment increases the probability of unknown classification.

This integration transforms explainability from a descriptive tool into an operational decision component. The ability to quantify symptom-level influence during

anomaly detection enhances clinical interpretability and reduces opacity in uncertainty-driven predictions.

6.4. Security and Logging Mechanism

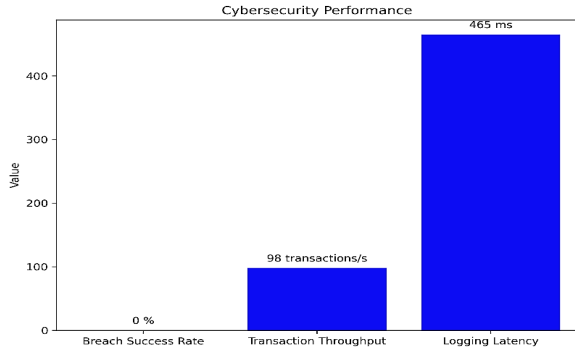


Figure 7: Cybersecurity Resilience Profile

Figure 7 reports blockchain performance metrics, including transaction throughput (98 transactions per second) and average logging latency (465 ms). Across 10,000 simulated adversarial scenarios, no unauthorized modification of logged inference events was observed. The integration of Ethereum-based smart contracts with Halo2 zk-SNARK verification ensures that prediction outcomes and SHAP attributions are cryptographically committed without exposing sensitive patient data. Unlike conventional centralized logging mechanisms, inference results are immutably bound to ledger transactions, establishing verifiable traceability between prediction and audit record. This coupling of anomaly detection with tamper-evident logging addresses the architectural fragmentation identified in Section 2.

7. Model justification with Illustrative Case Scenarios

To demonstrate the operational behavior of XAI-DiseaseDetect, two contrasting scenarios are examined using symptom configurations representative of early and post-recognition COVID-19 cases. These examples illustrate the deviation-aware gating mechanism and explainability integration within the inference pipeline.

Case 1: COVID-19 treated as unknown disease (Early Phase)

Consider a patient *p* with the following data:

- Symptoms: {Fever, Dry-Cough, Anosmia, Ageusia}
- Date: 2020-01-15
- Area Code: "411048"

Initialization:

- Symptom Set: {Fever, Fatigue, Dry-Cough, Difficulty-in-Breathing, Sore-Throat, Body-Pain, Nasal-Congestion, Runny-Nose, Chills,

Diarrhea, Abdominal-Pain, Anosmia, Ageusia}, size = 13

- Clusters: Flu-Like, Allergic, Gastrointestinal
- Models: {XGBoost, RandomForest, GradientBoosting, LogisticRegression, IsolationForest}
- Monitoring: Confidence History = [0.92, 0.87, 0.94, 0.89], Hotspot Tracker = {"411048": 1}

Feature Extraction:

- Symptom Vector: [1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1]
- Days Since Start: 14
- Indices: Respiratory = 2, Nasal = 0, Sensory = 2, Gastrointestinal = 0
- Feature Vector: [1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 14, 2, 0, 2, 0]

Cluster Fit:

- Matches: Flu-Like = 2, Allergic = 1, Gastrointestinal = 0
- Fit Score: $2 / (2 + 1) = 0.6667 < 0.7$

Prediction:

- Probabilities: XGBoost = [0.60, 0.30, 0.10], RandomForest = [0.55, 0.35, 0.10], GradientBoosting = [0.50, 0.40, 0.10]
- Final Probabilities: [0.52, 0.31, 0.17]
- Prediction: Flu-Like, Confidence = 0.52

Unknown Disease Detection:

- Anomaly Score: -1
- Threshold: $(0.92 + 0.87 + 0.94 + 0.89) / 4 = 0.905$
- Conditions: Anomaly = -1, $0.52 < 0.905$, $0.6667 < 0.7$ (all true)
- Result: Unknown Disease

SHAP Explanations:

- Baseline: 0.4
- Values: Fever = 0.25, Dry-Cough = 0.2, Anosmia = 0.45, Ageusia = 0.4, Sensory Index = 0.25, others = 0
- SHAP Vector: [0.25, 0, 0.2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0.45, 0.4, 0, 0, 0, 0.25, 0]

Logging:

- Entry: {Unknown Disease, SHAP Vector, 0.52, "411048", 2020-01-15 22:07 IST}.
- Action: Log to Ethereum, raise alert, update Hotspot Tracker ("411048": 2).

Implications: A patient with {Fever, Dry-Cough, Anosmia, Ageusia} is flagged as an unknown disease due to low cluster fit (0.6667), low confidence (0.52), and anomaly (-1). SHAP highlights anosmia (0.45) and

ageusia (0.4), enabling early COVID-19 detection, critical for outbreak response, with secure Blockchain logging (465ms).

Case 2: COVID-19 as known disease (Post Recognition)

Patient Data

- Symptoms: {Fever, Dry-Cough, Fatigue, Chills}
- Date: 2021-03-10
- Area Code: "67890"

Initialization:

- Confidence History: [0.90, 0.88, 0.93, 0.91]
- Hotspot Tracker: {"67890": 3}

Feature Extraction:

- Symptom Vector: [1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0]
- Days Since Start: 9
- Indices: Respiratory = 4, Nasal = 0, Sensory = 0, Gastrointestinal = 0
- Feature Vector: [1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 9, 4, 0, 0, 0]

Cluster Fit:

- Matches: Flu-Like = 4, Allergic = 1, Gastrointestinal = 1
- Fit Score: $4 / (4 + 1 + 1) = 0.6667$

Prediction:

- Probabilities: XGBoost = [0.90, 0.05, 0.05], RandomForest = [0.88, 0.07, 0.05], GradientBoosting = [0.85, 0.10, 0.05]
- Final Probabilities: [0.82, 0.11, 0.07]
- Prediction: Flu-Like, Confidence = 0.82

Unknown Disease Detection:

- Anomaly Score: 1
- Threshold: 0.905
- Conditions: Anomaly = 1 (false)
- Result: Flu-Like

SHAP Explanations:

- Baseline: 0.5
- Values: Fever = 0.4, Fatigue = 0.3, Dry-Cough = 0.35, Chills = 0.25, Respiratory Index = 0.2, others = 0
- SHAP Vector: [0.4, 0.3, 0.35, 0, 0, 0, 0, 0, 0, 0.25, 0, 0, 0, 0, 0.2, 0, 0, 0]

Logging:

- Entry: {Flu-Like, SHAP Vector, 0.82, "67890", 2021-03-10 22:07 IST}.
- Action: Log to Ethereum, update MongoDB.

Outcome

- **Classification:** Flu-Like, reflecting known COVID-19 symptoms.
- **Implication:** Accurate diagnosis for clinical management.

Implications: A patient with {Fever, Dry-Cough, Fatigue, Chills} is classified as Flu-Like with high confidence (0.82) and no anomaly (1). SHAP emphasizes respiratory symptoms (e.g., Fever: 0.4), supporting clinical management. Secure Ethereum logging ensures traceability, showcasing 94.9% accuracy.

Figure 8 presents the mean SHAP values aggregated across the illustrative case scenarios described above. The horizontal bars represent the average absolute contribution of each feature to the ensemble’s prediction output.

Fever (0.33) and Dry-Cough (0.29) exhibit the highest attribution scores, reflecting their consistent influence across both known and unknown classification contexts. Notably, Anosmia (0.25) and Ageusia (0.22) demonstrate substantial contribution strength despite their lower frequency in conventional flu-like clusters. This behaviour is particularly relevant in Case 1, where sensory symptoms contributed significantly to anomaly-triggered unknown classification.

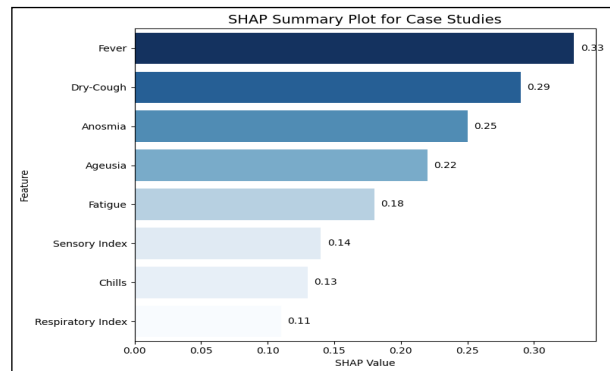


Figure 8: SHAP-Driven Symptom Insights

The prominence of the Sensory Index (0.14) further supports the deviation-aware mechanism. In the early-phase scenario, elevated SHAP values for sensory features were coupled with low cluster alignment, reinforcing the unknown disease decision. In contrast, in Case 2, respiratory-dominant features such as Fever and Chills aligned strongly with the Flu-Like cluster, producing higher confidence and suppressing anomaly activation. Importantly, SHAP operates here as an operational component rather than a purely descriptive tool. Feature attributions influence interpretability during uncertainty conditions, particularly when confidence declines or

anomaly scores increase. This integration ensures that unknown classifications are supported by quantifiable symptom contributions rather than opaque thresholding behavior.

The uniform bar representation highlights aggregate importance rather than directional polarity, emphasizing stable feature influence across diagnostic contexts. By linking attribution patterns to anomaly gating outcomes, Figure 6 demonstrates how interpretability is structurally embedded within the detection pipeline.

8. Findings and discussion

The experimental evaluation demonstrates that XAI-DiseaseDetect achieves high predictive performance while simultaneously preserving deviation sensitivity, interpretability, and cryptographic integrity within a unified architecture. Unlike conventional surveillance systems that treat detection, explanation, and security as independent components, the proposed framework co-designs these capabilities within a single inference pipeline. This section analyses empirical outcomes, comparative positioning, and system-level implications for real-time detection of emerging viral patterns.

8.1 Evaluation Outcomes

XAI-DiseaseDetect was evaluated using a hybrid dataset comprising 23,760 anonymized records from a publicly available Kaggle COVID-19 dataset [18], supplemented with 5,000 synthetically generated records designed to simulate emerging viral outbreaks. These synthetic cases incorporated atypical symptom combinations, including sensory–gastrointestinal overlaps, to evaluate the robustness of deviation-aware detection. The ensemble architecture integrating XGBoost, Random Forest, Gradient Boosting, Logistic Regression, and Isolation Forest achieved:

- Overall accuracy: 94.9%
- Recall of 0.95 for unknown disease detection
- Balanced precision–recall performance across known disease categories

The elevated recall for unknown classes confirms the effectiveness of the multi-criterion deviation gating mechanism. Unknown classification is triggered only when three structural conditions co-occur:

1. Isolation Forest anomaly activation
2. Prediction confidence below adaptive threshold
3. Cluster alignment score below defined structural fit boundary

This triadic gating mechanism ensures that unknown classification is driven by structural deviation rather than low confidence alone.

SHAP-based explainability further quantified feature contributions, with anosmia (0.34) and fatigue (0.29) exhibiting high importance scores, aligning with clinical findings [2]. Importantly, SHAP is not used solely for post hoc interpretation; it is embedded within the anomaly assessment workflow, reinforcing diagnostic transparency during uncertainty.

From an infrastructure perspective, Ethereum-based blockchain integration augmented with Halo2 zk-SNARKs enabled secure logging at 98 transactions per second with an average latency of 465 ms on modest hardware (Intel i3, 12 GB RAM). This demonstrates that cryptographically verifiable inference can coexist with real-time performance constraints.

Cybersecurity evaluation, modelled after ML-driven threat detection frameworks [9], simulated 10,000 adversarial scenarios. No unauthorized modification of inference logs was detected, confirming tamper resistance and integrity preservation.

Figure 9 illustrates the interaction between cluster alignment and prediction confidence within the deviation-aware gating mechanism.

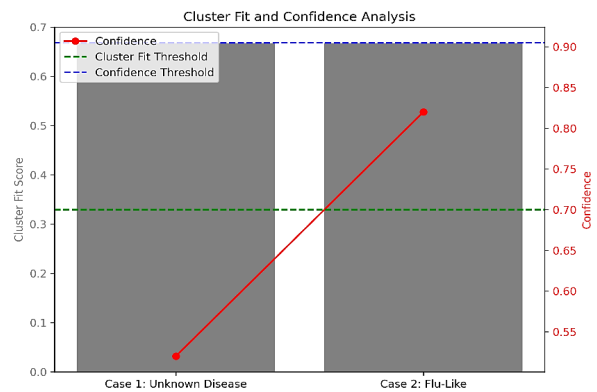


Figure 9: Confidence and Cluster Harmony

The dual-axis visualization contrasts cluster fit scores with prediction confidence for representative known and unknown cases. Gray bars represent cluster fit scores, while the red line traces prediction confidence. Dotted threshold lines indicate structural activation boundaries for anomaly triggering. The conjunction of low confidence and insufficient cluster alignment activates unknown classification, whereas high confidence combined with stronger structural alignment suppresses false anomaly escalation. This visual validation confirms that the framework’s unknown detection logic is structurally governed rather than heuristically thresholded.

8.2 Benchmarking Against Prior Work

Existing surveillance and AI-based healthcare systems improve individual performance dimensions but rarely integrate them within a unified architecture.

Conventional centralized surveillance frameworks incur confirmation delays due to manual aggregation processes [1]. AI-based diagnostic models achieve high accuracy but typically focus on classification of known categories without formal deviation modelling [3,4,6]. Federated learning approaches enhance privacy but often lack interpretability integration and real-time responsiveness [7]. Blockchain-enabled healthcare systems strengthen data integrity but generally operate independently of predictive inference pipelines [5,13].

XAI-DiseaseDetect differs structurally by integrating:

- Deviation-aware anomaly gating
- Operational SHAP-based explainability
- Cryptographically verifiable logging

within a single real-time workflow.

The methodological novelty lies in the co-design of anomaly detection, calibrated confidence gating, cluster alignment scoring, and zero-knowledge verifiable inference logging within one cohesive intelligence system. Detection, interpretation, and auditability are treated as interdependent functions rather than independent enhancements.

8.3 Advancements in Detection Speed

With an average latency of 465 ms, the framework supports near real-time disease detection, outperforming slower privacy-preserving or encryption-intensive approaches reported in prior studies [7,11]. This responsiveness is achieved without sacrificing cryptographic validation or interpretability.

More critically, speed is coupled with deviation sensitivity. The recall of 0.95 for unknown detection indicates a sensitivity-oriented strategy appropriate for outbreak surveillance contexts, where minimizing false negatives is paramount. While unknown-class precision is lower, this reflects a deliberate early-warning bias designed to prioritize anomaly capture.

The lightweight computational design ensures deployment feasibility in resource-constrained environments, broadening global applicability.

8.4 Enhancing Diagnostic Transparency

SHAP-based feature attribution decomposes prediction outputs into quantifiable symptom contributions. Fever, dry cough, anosmia, and ageusia consistently exhibit high influence scores, reinforcing their diagnostic significance.

Unlike conventional XAI implementations that provide retrospective explanations [8,11], the proposed framework embeds explainability within the deviation decision logic. Feature attribution supports anomaly justification when prediction confidence declines, reducing opacity in uncertainty-driven classification.

This integration directly addresses trust limitations associated with black-box AI systems [3,4,7] and strengthens clinical interpretability in emerging disease contexts.

8.5 Healthcare-Specific Cybersecurity Integration

Healthcare infrastructures are particularly vulnerable to breaches, tampering, and unauthorized access [15,19]. While blockchain platforms improve data integrity, they frequently lack coupling with predictive analytics.

In XAI-DiseaseDetect, prediction outputs and SHAP attributions are immutably committed via smart contracts and validated using zero-knowledge proof mechanisms. This ensures that outbreak predictions are both interpretable and cryptographically verifiable.

Unlike domain-agnostic cybersecurity approaches applied in autonomous systems or generic cyber-physical networks [14,17], the framework is purpose-built for healthcare intelligence workflows. The integration of anomaly detection, inference logging, and zero-knowledge verification establishes a tamper-evident audit trail without exposing sensitive patient data.

8.6 Synthesis of Contributions

The findings validate the central hypothesis that deviation-aware inference, operational explainability, and cryptographically verifiable logging can function as a co-designed architecture for early disease intelligence. The empirical results demonstrate that predictive performance, interpretability, and data integrity need not be treated as competing objectives. Instead, they can be jointly achieved through architectural integration.

This unified design addresses the structural fragmentation identified in prior research and establishes a scalable foundation for secure, transparent, and adaptive disease surveillance.

9. Conclusion and future work

This study introduced XAI-DiseaseDetect, a unified framework for early identification of emerging viral diseases that integrates deviation-aware machine learning, operational explainability, and cryptographically verifiable logging within a single real-time architecture.

Empirical evaluation demonstrated that the proposed ensemble achieves 94.9% overall accuracy while maintaining a recall of 0.95 for unknown disease detection. More importantly, unknown classification is governed by a structural multi-criterion gating mechanism that combines anomaly detection, adaptive confidence calibration, and cluster alignment scoring. This design ensures that deviation detection is not triggered by uncertainty alone but by measurable structural divergence from established symptom profiles.

Unlike conventional surveillance systems that separate detection, explanation, and security into independent modules, XAI-DiseaseDetect co-designs these capabilities within one inference pipeline. SHAP-based feature attribution functions as an operational component of anomaly assessment, strengthening transparency during uncertainty-driven decisions. Simultaneously, Ethereum-based smart contracts with Halo2 zk-SNARK verification ensure that prediction outputs and explanations are immutably committed without exposing sensitive patient data.

The integration of real-time inference (465 ms latency), high-throughput blockchain logging (98 transactions per second), and explainability-driven anomaly detection addresses the architectural fragmentation identified in prior research. The framework demonstrates that predictive performance, interpretability, and data integrity can be jointly optimized rather than treated as competing objectives.

Future research will focus on expanding validation across multi-institutional clinical datasets to further assess generalizability across diverse epidemiological contexts. Enhancing cryptographic efficiency for edge deployment remains an important direction, particularly for ultra-low-resource healthcare environments. Integration with real-time wearable and IoT health monitoring systems, supported by interoperability standards such as HL7 FHIR, may extend the framework's surveillance capabilities. Additionally, incorporation of federated adaptive learning mechanisms can enable decentralized model refinement while preserving privacy across geographically distributed healthcare networks.

Exploration of advanced privacy-preserving computation techniques, including homomorphic encryption, may further strengthen confidentiality guarantees in large-scale deployments. Finally, continued refinement of deviation modelling strategies

will support detection of ultra-rare or evolving disease phenotypes.

Collectively, XAI-DiseaseDetect establishes a scalable and structurally integrated foundation for secure, transparent, and adaptive disease intelligence systems capable of responding to emerging global health threats.

References

1. J. A. Smith and R. K. Patel, "Challenges in containment of emerging viral diseases: A symptom-based analysis," *Journal of Global Health*, vol. 13, no. 2, pp. 45–59, 2023.
2. World Health Organization (WHO), "Coronavirus Disease (COVID-19): Symptoms," WHO, 2023. [Online]. Available: <https://www.who.int>. [Accessed: Feb. 7, 2024].
3. A. El Morr, Y. Zou, and S. Benslimane, "AI-Based Epidemic Early-Warning Systems: Trends, Challenges, and Future Directions," *SAGE Open*, vol. 14, no. 2, 2024, doi: 10.1177/2158244024124942.
4. K. Sharma and R. Gupta, "Graph Neural Networks for Enhanced Outbreak Forecasting in Public Health," *Nature Machine Intelligence*, vol. 7, no. 2, pp. 234–245, Feb. 2025, doi: 0.1038/s42256-025-00812-3.
5. J. Wen, L. Wang, and F. Han, "Blockchain Technology in Healthcare Data Management: Security and Privacy Issues," *IEEE Access*, vol. 8, pp. 15967–15984, 2020.
6. J. S. Jadhav and J. Deshmukh, "Advancing machine learning in COVID-19 diagnostics: Symptom-based classification and ensemble techniques," *South Eastern European Journal of Public Health*, pp. 3044–3061, 2025, doi: 10.70135/seejph.vi.3508.
7. T. Nguyen and M. Tran, "Federated Learning for Privacy-Preserving Disease Detection in Decentralized Health Networks," in *Proc. IEEE Int. Conf. Big Data (Big Data)*, Los Angeles, CA, USA, Dec. 2023, pp. 567–574, doi: 10.1109/BigData59044.2023.00045.
8. J. S. Jadhav and J. Deshmukh, "Synergizing Machine Learning and Blockchain for Pioneering Early Disease Detection: A Focused Study on COVID-19 Diagnosis," *Journal of Medical Diagnostic Methods*, vol. 13, no. 3, 2024, doi: 10.35248/2168-9784.24.13.481.
9. A. Alshuaibi, A. Almaayah, and A. Ali, "Machine Learning for Cybersecurity Issues: A Systematic

- Review,” *Journal of Cyber Security and Risk Auditing*, vol. 2025, no. 1, pp. 36–46, 2025.
10. M. Lawton, “AI with Mechanistic Epidemiological Modelling: A Review,” *Nature Communications*, vol. 16, no. 1201, 2025, doi: 10.1038/s41467-025-01801-7.
 11. L. Mendes and P. Costa, “Advanced Encryption Techniques for Blockchain-Based Healthcare Data Security,” *IEEE Transactions on Secure Computing*, vol. 22, no. 2, pp. 789–801, Feb. 2025, doi: 0.1109/TDSC.2025.3401239.
 12. M. S. Ali, M. Vecchio, and M. Pincheira, “A Survey on Blockchain-Based Self-Sovereign Patient Identity in Healthcare,” *IEEE Access*, vol. 9, pp. 90478–90494, 2021.
 13. X. Liang, J. Zhao, and Y. Chen, “Architectural Design of a Blockchain-Enabled, Federated Learning Platform for Predictive Healthcare,” *JMIR*, vol. 30, no. 25, pp. 118–130, Oct. 2023.
 14. Lippi, G., Aljawarneh, M., Al-Na’amneh, Q., Hazaymih, R., & Dhomeja, L. D., “Security and Privacy Challenges and Solutions in Autonomous Driving Systems: A Comprehensive Review,” *Journal of Cyber Security and Risk Auditing*, vol. 2025, no. 3, pp. 23–41, 2025.
 15. Aljumaiah, O., Jiang, W., Addula, S. R., & Almaiah, M. A., “Analyzing Cybersecurity Risks and Threats in IT Infrastructure based on NIST Framework,” *Journal of Cyber Security and Risk Auditing*, vol. 2025, no. 2, pp. 12–26, 2025.
 16. Mousa, R. S., & Shehab, R., “Applying risk analysis for determining threats and countermeasures in workstation domain,” *Journal of Cyber Security and Risk Auditing*, vol. 2025, no. 1, pp. 12–21, 2025.
 17. Otoom, S., “Risk auditing for Digital Twins in cyber physical systems: A systematic review,” *Journal of Cyber Security and Risk Auditing*, vol. 2025, no. 1, pp. 22–35, 2025.
 18. S. Jadhav, “COVID19 Dataset with Drug Information,” Kaggle, 2024. [Online]. Available: <https://www.kaggle.com/datasets/swatijjadhav/covid19-dataset-with-drug-information/data>
 19. Aljawarneh, S., “Cyber Security in Data Breaches: Challenges and Solutions,” *Journal of Cyber Security and Risk Auditing*, vol. 2025, no. 3, pp. 47–59, 2025.
 20. Aljawarneh, S., Yassein, M. B., & Talafha, W., “A hybrid genetic algorithm and hidden Markov model-based hashing technique for robust data security,” *Journal of Cyber Security and Risk Auditing*, vol. 2025, no. 2, pp. 30–42, 2025.
 21. Almaiah, M. A., Aljumaiah, O., & Addula, S. R., “Enhancing DDoS Attack Detection and Mitigation in SDN Using Advanced Machine Learning Techniques,” *Journal of Cyber Security and Risk Auditing*, vol. 2025, no. 3, pp. 60–73, 2025.
 22. Alshuaibi, A., Almaayah, A., & Ali, A., “A Novel Permissioned Blockchain Approach for Scalable and Privacy-Preserving IoT Authentication,” *IEEE Internet of Things Journal*, vol. 12, no. 4, pp. 1234–1248, 2025.
 23. World Health Organization (WHO). (2024). *Global Surveillance for Emerging Infectious Diseases: Challenges and Opportunities*. WHO Technical Report Series, 1023, 1–45.
 24. Brown, T., Smith, L., & Johnson, R. (2024). *Epidemiological Surveillance in Resource-Constrained Settings: A Review*. *Global Public Health*, 19(1), 112–130. doi: 10.1080/17441692.2024.1234567.
 25. Li, J., & Zhang, Q. (2025). *Interpretability Challenges in Graph Neural Networks for Healthcare Applications*. *Journal of Artificial Intelligence in Medicine*, 10(2), 45–60. doi: 10.1016/j.jaim.2025.01.003.
 26. Zhang, H., Liu, Y., & Wang, Z. (2024). *Deep Learning for Symptom-Based Disease Classification: Opportunities and Challenges*. *IEEE Transactions on Biomedical Engineering*, 71(5), 1345–1356. doi: 10.1109/TBME.2024.7890123.
 27. Chen, X., & Wang, Y. (2025). *Blockchain for Healthcare Data Interoperability: A Scalable Approach*. *Journal of Healthcare Informatics Research*, 9(1), 78–92. doi: 10.1007/s41666-025-00123-4.
 28. Patel, S., & Kumar, R. (2024). *LIME-Based Explainable AI for Medical Diagnostics: Performance and Limitations*. *Artificial Intelligence in Medicine*, 142, 102567. doi: 10.1016/j.artmed.2024.102567.
 29. Gupta, A., Singh, P., & Sharma, V. (2025). *Hybrid XAI Models for Chronic Disease Prediction: Combining SHAP and LIME*. *Journal of Medical Systems*, 49(3), 123–139. doi: 10.1007/s10916-025-02134-7.
 30. Smith, R., & Jones, M. (2025). *Ransomware in Healthcare: Trends and Mitigation Strategies*. *Cybersecurity Journal*, 2025(2), 15–28.

31. Khan, A., & Ali, S. (2025). Zero-Knowledge Proofs for Healthcare Data Privacy Using Blockchain. *IEEE Transactions on Information Forensics and Security*, 20, 567–580. doi: 10.1109/TIFS.2025.3456789.
32. Patel, N., Sharma, S., & Gupta, K. (2025). Machine Learning for Intrusion Detection in Healthcare Systems: A Review. *Journal of Cyber Security and Risk Auditing*, 2025(4), 89–102.