# Machine Learning-Driven Assessment and Security Enhancement for Electronic Health Record Systems

Check for updates

## Birendra Kumar Saraswat[1*], Neeraj Varshney[2] and Prem Chand Vashist[3]

[1]Depertment of Computer Science & Engineering, GLA University, Mathura, 281406, India;
[2]Depertment of Computer Engineering & Applications, GLA University, Mathura, 281406, India; [3]Depertment of Information Technology, GL Bajaj Institute of Technology and Management, Greater Noida- 201306, India

**E-mail/Orcid Id:**

*BKS,* ✉ birendrasaraswat@gmail.com, 🆔 https://orcid.org/0000-0003-0628-6823;
*NV,* ✉ neeraj.varshney@gla.ac.in, 🆔 https://orcid.org/0000-0002-6537-7891;
*PCV,* ✉ pcvashist@gmail.com, 🆔 https://orcid.org/0000-0002-5350-4064

**Abstract:** The digitalized patient-centric system, the Electronic Health Record (EHR), is a platform where comprehensive health information is stored, managed, and accessed electronically. The primary findings of this study aim to secure sensitive patient data and increase overall system resilience by demonstrating that machine learning can evaluate vulnerabilities and improve the security of Electronic Health Record (EHR) systems. This research examines the prospects of incorporating machine learning-driven assessment tools and safety improvements in EHRs to enhance data protection in the healthcare industry. The proposed method utilizes the implementation of machine learning classifiers, specifically the XGBoost and LightGBM models. These classifiers are employed to enhance various aspects of the system, such as data protection and security, within the framework of EHRs. The study emphasizes the efficiency of these machine learning classifiers in ensuring that EHR systems are secure enough to deal with any problem that may occur due to threats posed by external factors or hackers. The findings reveal that the XGBoost model always has outstanding performance, with a near-perfect Receiver Operating Characteristic Curve (ROC) having an AUC equal to 1.00, indicating close to perfect accuracy in distinguishing positive from negative cases. Similarly, LightGBM has a perfect ROC curve as well. Therefore, its performance would be considered flawless. Consequently, future developments could lead to sophisticated machine learning models besides those that have already been developed. Improving data storage through encryption and building safer communication protocols should also be considered to make these systems withstand new security problems. Thus, this study contributes to the existing literature on applying technology to safeguard vulnerable medical records while fostering a safe and efficient healthcare ecosystem.

## Introduction

There is growing interest in using data from electronic health records (EHRs) for patient registries. This study aimed to examine how EHR interoperability impacts patient safety and other dimensions of care quality in high-income healthcare settings. EHRs are electronic systems used and maintained by healthcare organizations to collect and store patients' medical information. The study concluded that patient registries are patient-centered, purpose-driven, and designed to derive information on specific exposures and health outcomes.

These databases store a patient's medical history, diagnosis, prescription, immunization dates, allergies, radiographs, and test results. EHRs improve patient treatment coordination and medical professional communication. Healthcare administrators can make evidence-based choices, reduce medical errors, and speed up administrative operations with electronic health information. By increasing interoperability and data exchange, EHRs improve patient outcomes and healthcare efficiency, which promotes medical research, public health, and healthcare policy. EHRs, also referred

160

to as Electronic Medical Records (EMRs), are electronic counterparts of traditional paper health records created, managed, and preserved by care providers (Mayer et al., 2020). These records are exclusively accessible to patient caregivers. Personal health records (PHRs) enable patients to manage and update their medical histories. EHRs are protected by the Health Insurance Portability and Accountability Act (HIPAA), not personal records (Himabindu et al., 2024). Patients' electronic health records document several healthcare practitioner visits in diverse places. Figure 1 depicts the EHR System.

have gained prominence for overcoming difficulties encountered by conventional Natural Language Processing (NLP) strategies for extracting data (Locke et al., 2021; Osmani et al., 2018).

The ability to autonomously learn representations from data is a significant factor contributing to the increasing popularity of deep learning (DL) models (Xiao et al., 2018). Their resilience to high-complexity functions grows alongside the scale of the dataset (Esteva et al., 2019). Clinical applications of ML/DL models hold considerable promise for revolutionizing the healthcare
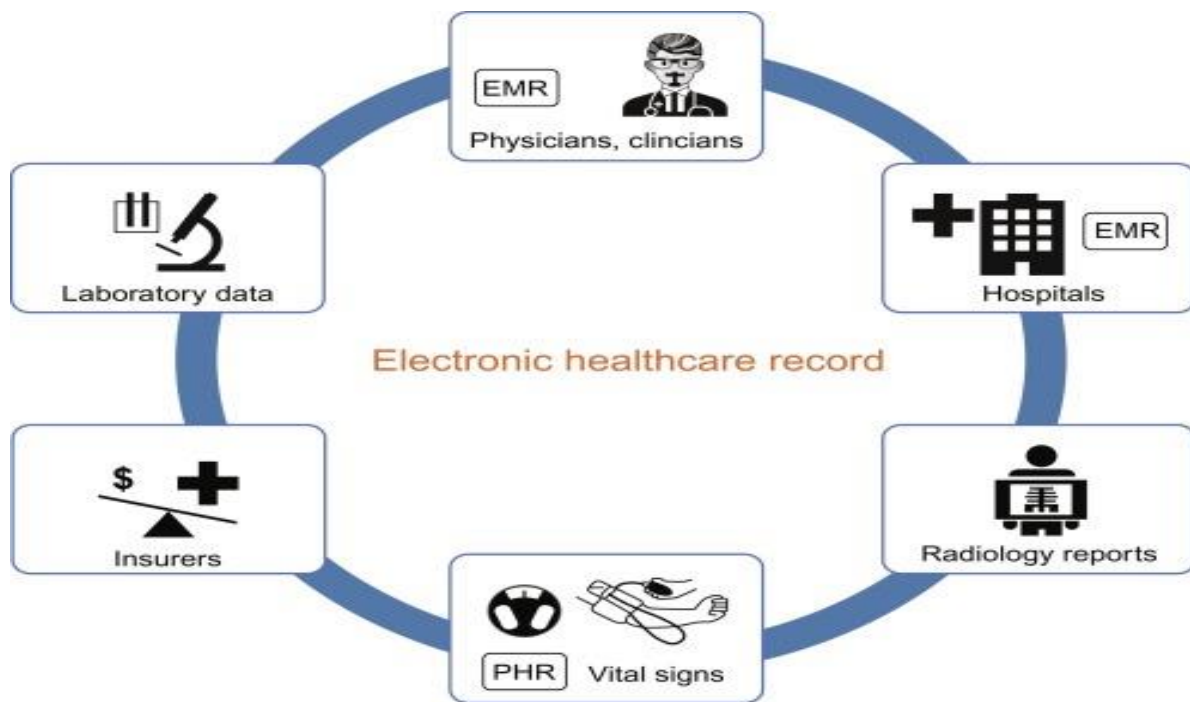


**Figure 1. Introduction to the EHR System (Kumar et al., 2029).**

EHR includes demographics, concerns, prescriptions, vital signs, medical history, immunizations, laboratory data, and imaging results. EHRs streamline clinical processes. Quality management, outcome reporting, and evidence-based decision support, which can capture a patient's clinical experience, are among the capabilities of an EHR. Medical professionals, hospitals, and other healthcare facilities utilize EHRs as the current standard. Healthcare facilities preserve patients' medical history, but they might be hard to access. Some hospitals provide patients with physical or electronic copies of their information, whereas others with more advanced systems provide secure online access (Lee et al., 2021). The healthcare sector is rapidly adopting digital technologies such as artificial intelligence (AI), machine learning, big data analytics, smart sensors, the Internet of Things (IoT), and robotics to improve the efficiency and quality of care provided. Advanced countries are leading this trend (Lee et al., 2019). Recently, Deep Learning (DL) techniques

industry. However, various privacy and security issues must be resolved before these techniques can be reliably applied in healthcare systems (Qayyum et al., 2020). ML can potentially exacerbate pre-existing health inequities, posing several ethical problems (Chen et al., 2021). Bioethics concepts can be used to create morally sound ML models in the healthcare industry (Keerthana et al., 2024; Vayena et al., 2018). Security threats like as data confidentiality, privacy, and integrity assaults are a major concern because of its ability to generate huge amounts of data at predictable intervals (Kumari et al., 2018). Strict access limits and other security measures protect patient data. Covered entities need internal operations controls. Protected Health Information (PHI) is used worldwide and supplied digitally. Healthcare privacy entails the preservation of confidential documents. Regulations and policies play a pivotal role in accomplishing this goal. Patients have the right to know who opens and uses their medical records. Patients' health

information is protected by HIPAA (Hathaliya et al., 2020).

This research signifies noteworthy progress in data security by incorporating multiple state-of-the-art characteristics. Combining the secure asymmetric encryption of Elliptical Curve Cryptography (ECC) with the powerful symmetric encryption of Advanced Encryption Standard (AES) creates a triple-layered defense mechanism with the Hyperledger blockchain. Besides, this research employs enhanced machine learning classifiers which are XG Boost and Light GBM, two widely accepted models with outstanding performances. Additionally, it enhances their predictability by employing an ensemble of two classifiers, XGBoost and LightGBM, rather than relying on a single classifier, thereby yielding a more robust and accurate outcome. Furthermore, both depend on blockchain for the same reason that they have developed encryption processes as well as ML. The encrypted EHR data is moved to a blockchain where the integrity and immutability of the data are guaranteed. This demonstrates that the increased security measures do not violate blockchain technology's decentralized or secure nature.

Moreover, a strong validity condition implies that sharing should occur between authorized entities only. Additionally, enhanced security is achieved through advanced encryption techniques, ensuring the secure transmission of data and safeguarding information effectively. Ultimately, this investigation conforms to contemporary security protocols and hence adheres to internationally acceptable encryption standards. The utilization of AES as a symmetrical encryption standard and ECC as a prevailing asymmetrical encryption method supports current security norms. In comparison to simple methodologies employed earlier, this creates a more robust secure framework for handling patient data.

## Related Work

A review of the literature analyzing the relevant work by different authors.

Huang et al. (2023) employed transparent ML methods to form a top-down arrangement of major predictors using model importance statistics like gain, cover, and frequency. The results revealed an average age of 74.05 with a standard deviation of 12.85. The AUROC (Area under the Receiver Operating Characteristic Curve) for the XGBoost model was 0.662. The SHAP explanations whose total values were greatest included urine output, leukocytes, bicarbonate, and platelets.

Yang et al. (2023) developed an application of ML for predicting acute respiratory distress syndrome (ARDS) in Intensive Care Unit (ICU) patients by creating a new model and validating it. The AUC values of the respective models were as follows: Logistic Regression (LR) was 0.664, K-Nearest Neighbour (KNN) was 0.692, support vector machine (SVM) was 0.567, Decision Trees Classifier (DTC) was 0.709, Random Forest (RF) was 0.732, XGBoost was 0.793, LightGB was 0.793, and CatBoost was 0.817.

Shah et al. (2023) presented a new approach to improving network security and analyzing data derived from Personal Health Records (PHR). When they analyzed data from individual health records, they employed neural networks with variational Boltzmann spatial encoder capabilities. They achieved a more secure network by using the decentralized blockchain architecture. The experimental investigation was conducted using data and network security. It measured random accuracy at 81%, specificity at 55%, latency at 62%, quality of service at 52%, and computational cost at 41%.

Alam et al. (2023) provided an application called FedSepsis for early sepsis detection leveraging EHRs. Several (Deep Learning) DL methods were utilized for the prediction and NLP jobs. Performance was satisfactory, and when devices were moderately numerous, the outcomes in the federated learning configuration were comparable to those in the single server-centric configuration. The most optimal approach was to use multimodality in conjunction with generative adversarial neural networks. The outcomes were a near-perfect accuracy rate of 96.55%, a receiver operating characteristic area of 99.35%, and a latency of 4.56 hours.

Corbin et al. (2022) explored the potential benefits of clinical decision support based on machine learning in the context of antibiotic prescribing management. A retrospective multi-site study was conducted, which trained ML models to anticipate antibiotic susceptibility patterns, also referred to as personalized antibiograms, using EHR data about 8342 infections at Stanford's emergency departments and 15,806 cases of uncomplicated UTIs at Boston's Massachusetts General Hospital and Brigham & Women's Hospital. Based on data from Stanford, clinicians were able to reallocate antibiotic selections with the help of tailored antibiograms, resulting in a coverage rate of 85.9%. This rate was comparable to clinician performance, which had been determined to be 84.3% (p = 0.11). The tailored antibiogram coverage percentage in the Boston dataset

was 90.4%, which was much better than the doctors' rate of 88.1% (p < 0.0001).

Tsiklidis et al. (2022) developed a model that could predict continuously the likelihood of patient death or a risk metric. The AUROC measured by the model was a measure of its accuracy. The author obtained an accuracy level for this model of 92.9%.

Pang et al. (2021) proposed seven machine learning models that utilized the EHR data from up to 2 years ago to predict the chances of children aged between 2 and 7 being obese. There were seven models, and their comparison was done using post-hoc pairwise testing as well as Cochran's Q test, while performance was evaluated using different standard classifier metrics. XGBoost outperformed all other models with an AUC of 0.81 (0.001). Besides, it performed better than the other models on traditional classifier metrics: accuracy 66.14% (0.41%), specificity 63.27% (0.41%), precision 30.90% (0.22%), and F1-score 44.60% (0.26%).

Hou et al. (2020) utilized XGboost to construct an ML model for predicting 30-day mortality in sepsis-3 patients admitted to the MIMIC-III database and determined if it outperformed traditional prediction models. According to the AUCs' results (0.819 [95% CI 0.800–0.838], 0.797 [95% CI 0.781–0.813] and 0.857 [95% CI 0.839–0.876]) and decision curve analysis of the three models, the XGboost model exhibited the best overall performance among the others. This was validated by the risk nomogram and clinical impact curve, where the XGboost model demonstrated good predictive value.

Souri et al. (2020) recommended an IoT-supported student health monitoring system where smart medical gadgets were used to trace the vital signs of students discreetly and any changes in their biology or behavior. The concept was to identify probable dangers connected with shifts in the way students behaved and what they did to their bodies by gathering important information from IoT gadgets and processing it with the help of machine learning mechanisms. The results obtained during the experiment confirmed that there was effective functioning and precision of this model concerning student health evaluations. After testing the proposed model, the SVM achieved the highest accuracy of 99.1%, which was encouraging for the aim. The outcomes were superior to those of algorithms based on decision trees, random forests, and multilayer perceptron in neural networks.

Vos et al. (2020) examined that EHRs could enhance collaboration among healthcare professionals, but their impact on teamwork remained clueless. When five outpatient clinics in a Dutch hospital with a comprehensive EHR system were examined, the research found mixed results. Although the system facilitated real-time coordination across specialties, it hindered interdisciplinary collaboration due to asynchronous access to patient records. While it streamlined certain tasks and facilitated data-based decision-making, specialized interfaces impeded data comprehension. Additionally, while it improved documentation efficiency, it also imposed rigid authorization requirements and increased administrative burdens on physicians, limiting flexibility.

Hirano et al. (2020) proposed an open-source, publicly available, and CNN-based COVID-Net model's vulnerability was examined. This model was among the first deep learning models to detect COVID-19 using chest X-ray (CXR) images. Two kinds of attacks—targeted and nontargeted—were investigated using perturbation created by the fast gradient sign technique (FGSM). The authors evaluated both the COVID-Net CXR small and CXR big models. Their results showed that both models had been able to attain success rates of >85% for non-targeted attacks and >90% for targeted attacks after adding 2% universal adversarial perturbations.

Mandair et al. (2020) investigated the development of a machine-learning model aimed at predicting the incidence of myocardial infarction (MI) within six months, utilizing harmonized electronic health record (EHR) data. The findings demonstrated that, compared to alternative models, a combination of random under-sampling with deep neural network (DNN) classification proved more effective. There were 2,531 patients with MI diagnosed in this study, while there were 2.25 million without MI diagnosis. The classification accuracy of a deep neural network trained with random under-sampling was much higher compared to other approaches. The moderate benefits of the deep neural network became apparent when compared to logistic regression using only known risk factors, namely, F1 Score is 0.092, and AUC is 0.835.

Newaz et al. (2019) proposed a new security framework called HealthGuard based on machine learning to identify malicious activities in Smart Healthcare System (SHS). The results showed that HealthGuard was an effective security framework for SHS, with an accuracy of 91% and an F1 score of 90%.

Bhattacharya et al. (2019) presented a framework called Blockchain-Based Deep Learning as a Service (BinDaaS). The integrated blockchain and DL methods for multiple sharing EHR records among several healthcare users were carried out in two phases. Different

parameters such as accuracy, end-to-end latency, mining time, computation, and communication costs were used to compare the obtained results with those of existing state-of-the-art proposals. Based on the results obtained, BinDaaS surpassed all other systems.

## Problem Statement

The major problem is providing strong security and authenticity to electronic health record (EHR) systems, considering the changing nature of cyber threats. The main concerns are preventing unauthorized access, breaches, or tampering with data that may lead to disclosure of patients' privacy and medical confidentiality. Furthermore, the complicated healthcare environments consisting of several stakeholders and interrelated systems make it difficult to ensure secure data exchange and compatibility. This investigation involves assessing factors such as key verification, clarity in representation, and ensuring the system meets modern encryption standards in healthcare data security.

## Dataset Description

A widely used medical dataset, Medical Information Mart for Intensive Care III (MIMIC-III), is available on Kaggle. The MIMIC-III dataset is massive, anonymous, and publicly available. Each entry in the dataset is accompanied by an ICD-9 code, documenting the diagnoses and procedures performed. These codes are further subdivided into sub-codes, in most cases indicating specific circumstances surrounding them. The data set is comprised of 112,000 clinical reports with an average length of 709.3 tokens and 1,159 top-level ICD-9 codes. On average, each report has been assigned to 7.6 codes. These data contain vital signs, prescriptions, laboratory measures, observations and notes recorded by healthcare professionals; fluid balance, procedure codes, diagnostic codes, imaging reports; hospital length of stay survival data; and additional patient information. This database supports applications like academic research and development or monitoring healthcare services,

## Research Methodology

An analysis of the designed architecture is conducted within the framework of the research technique.

## Technique Used

Various techniques used in the proposed method are the Hash-based Message Authentication Code (HMAC) Algorithm, AES algorithm with ECC for encryption, Blockchain, Cloud Computing, and ML Classifiers.

**HMAC Algorithm:** HMAC is a popular technique used in many different types of EHR systems and other areas of cybersecurity (Vignesh et al., 2017). When sending messages or data across different parts of an EHR system, HMAC is utilized to ensure that nothing has been tampered with along the way. EHR HMAC techniques improve patient data security. HMAC is used to fingerprint sensitive patient data in EHR systems to ensure data integrity. Intentional or not, changes create a new hash value that flags the file as compromised. Finally, HMAC can authenticate data in transit. An HMAC produced with the sender's private key can verify EHR data transfers. Recalculating and comparing the HMAC at the receiving end prevent unauthorized access to data in transit (Gabriel et al., 2021). Timestamps confirm data currency and replay attacks are prevented with HMAC. HMAC secures accounts and transmits data. When HMACs are produced using credentials, the system securely authenticates users.

**AES Algorithm:** EHR systems must implement the AES algorithm to protect patients' personal information. Data in healthcare applications like EHR systems can benefit from AES's high degree of security and efficiency because it is a frequently used symmetric encryption technique. The audit trails are unchangeable, and all parties are accountable when these logs are encrypted (McGhin et al., 2019). HIPAA and other healthcare standards require EHR systems to encrypt patient data and corporate procedures to protect liability.

**ECC:** ECC is a robust encryption method that can be used to set up safe lines of communication within EHR systems. It is computationally efficient and well-suited for resource-constrained contexts like those found in healthcare devices because of its robust security and reduced key lengths.

**Blockchain:** Blockchain technology could revolutionize healthcare EHR systems. Blockchain technology in EHRs has many benefits. Its distributed, unalterable ledger improves data security and integrity. Data breaches, hacking, and tampering with patient records are greatly reduced by blockchain technology. Each patient data transaction is saved as a block and linked to the one preceding it. Blockchain technology also solves healthcare data exchange and interoperability issues (Huang et al., 2019). The network maintains data accuracy and consistency by making patient data transfer between healthcare providers secure and fast. Smart contracts improve interoperability and permit automatic data transfer under certain conditions. Blockchain also lets patients manage their health records. Patients can grant and revoke cryptographic keys to restrict data access to those who need it. This open strategy protects

patient privacy and consent (Mishra et al., 2023). Figure 2 depicts the working of blockchain technology.

**Cloud Computing:** Cloud computing in EHR systems is altering healthcare by making patient health information management, storage, and access flexible and efficient. The cloud-based EHR technology benefits healthcare workers and changes patient data storage and access. They allow hospitals to scale their data storage and processing capacities without investing in additional facilities, which is a huge benefit. The ability to view patient records from anywhere with an internet connection improves medical staff mobility and accessibility, leading to faster and better medical decisions. Healthcare providers can improve patient care and quality by eliminating the need for pricey on-premises gear and software. Moreover, cloud-based EHR solutions must prioritize data security and compliance (Chenthara et al., 2019).

**i) XG Boost:** XG Boost (Extreme Gradient Boosting) is a machine learning classifier that is known for its efficiency and effectiveness. EHR systems have seen XGBoost shine in predictive modeling tasks such as disease diagnosis and risk prediction (Romeo et al., 2020). It can also fill gaps, deal with convoluted interactions between data points, and measure how important each feature is in healthcare analytics. The accuracy of illness forecasting could increase, high-risk patients could be identified, and healthcare practitioners could maximize limited resources using XGBoost.

**ii) Light GBM:** Light GBM is another gradient-boosting technique that can handle many features and a huge volume of datasets. Light GBM has a range of uses or applications in EHR systems, including prognosis modeling, drug response modeling, and outlier detection. It is well suited for real-time health data applications as it trains very quickly and can be applied to big datasets.
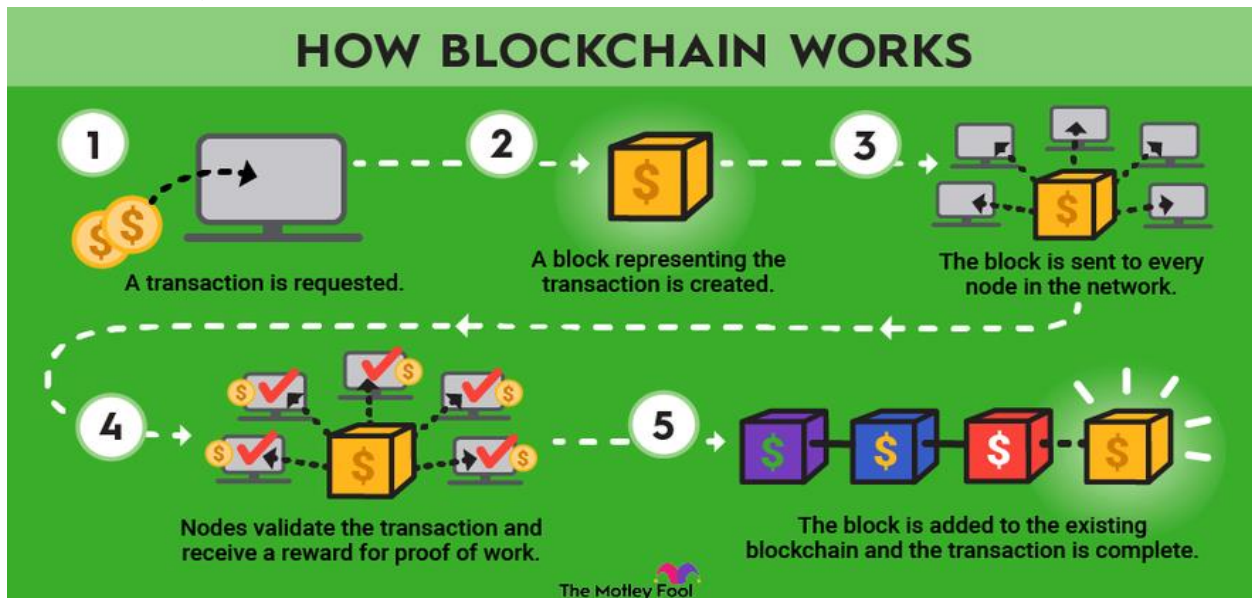


**Figure 2. Blockchain Technology (Ahmadi and Aslani, 2018).**

**ML Classifier:** Electronic Health Record machine learning classifiers play a crucial role in illness prediction, patient risk stratification, and treatment recommendation. EHR utilizes machine learning algorithms, such as neural networks, decision trees, and support vector machines (SVM), to manage extensive patient data. Those who are engaged in the health sector would finally be able to tell what is likely to happen. ML classifiers might result in a high-quality healthcare service that everyone can afford because they detect diseases before they become severe, predict hospital re-admissions, and optimize treatment options (Hasan et al., 2029). The machine learning classifiers mentioned below are XGBoost and LightGBM.

LightGBM assists healthcare organizations in improving their forecasting capabilities, leading to faster preventive actions and personalized treatment. Likewise, it finds deviations or anomalies in the patient information that could help detect health conditions early, thereby enhancing patient safety (Chami et al., 2019).

**Proposed Methodology**

Figure 3 demonstrates the proposed method in a diagrammatic form and outlines a system for storing and retrieving electronic health records (EHRs) using blockchain, HMAC authentication, encryption, and machine learning (ML) classification.
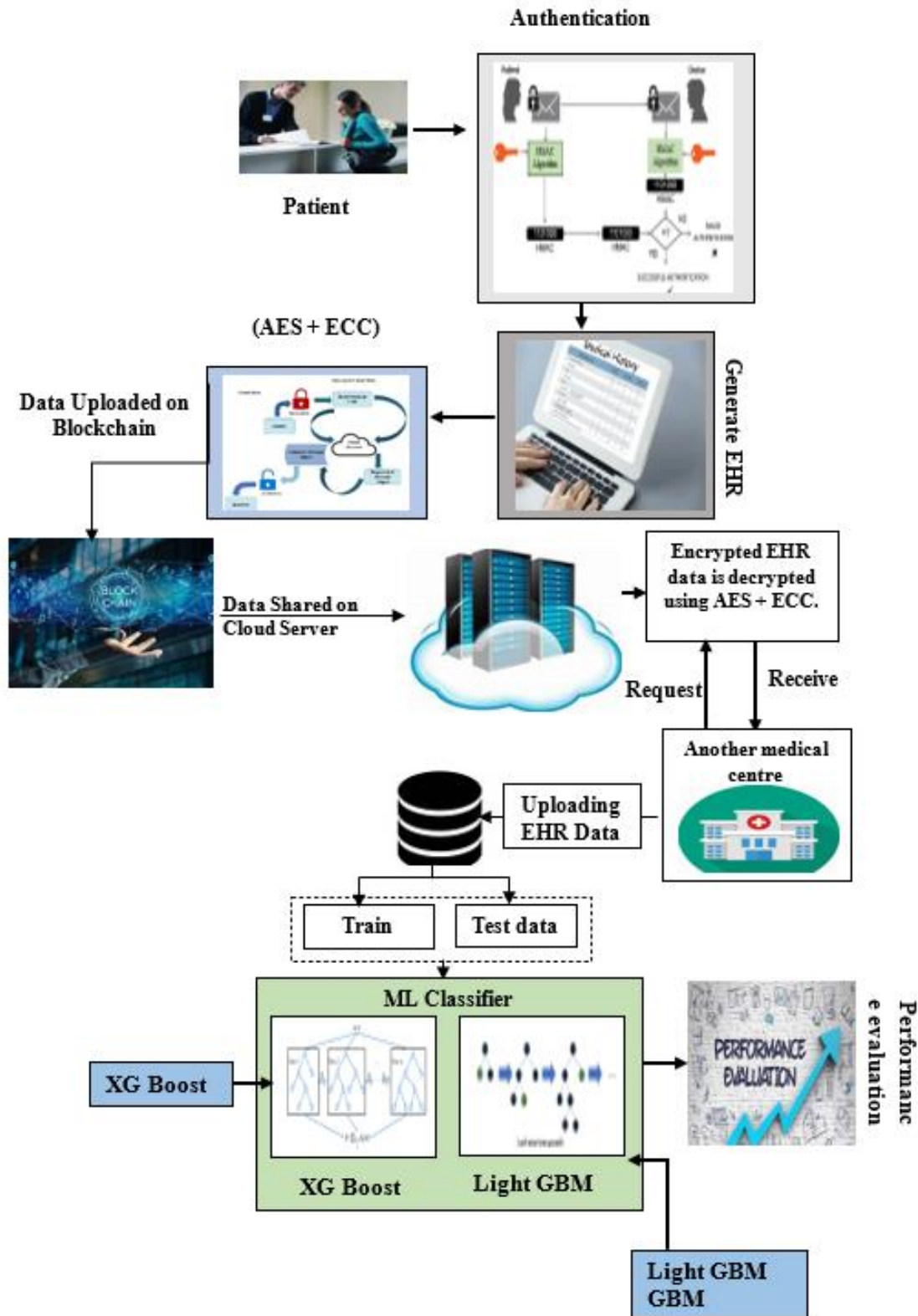
**Figure 3. Proposed Methodology.**

## Proposed Algorithm

The author uses mathematical notations to symbolize some individual steps:

**Step 1:** Registration of Patient:

Assume P represents the patient set, where P consists of elements p1 to pn.

On registering a new patient, add to the set P:

P' = P U {p$_n$+1}

**Step 2:** HMAC authentication:

Generating the authentication code involves utilizing the HMAC function in the following way:

$$HMAC(K, M) = H((K \oplus o_{pad}) \parallel H((K \oplus i_{pad}) \parallel M))$$

Where, $\oplus$ signifies the bitwise XOR operation, K denotes the secret key, M represents the message for authentication, $i_{pad}$ signifies the inner padding (repeated byte, typically 0x36), $o_{pad}$ represents the outer padding (repeated byte, typically 0x5C), $\parallel$ denotes concatenation, and H stands for a cryptographic hash function (e.g., SHA-256, SHA-512).

The authentication code can be calculated as follows:

authentication_code = H (key, patient_info)

**Step 3:** Encryption using AES + ECC:

Encrypting the message using AES: $C_{AES} = AES\_Encrypt(M, K_{AES})$

Encrypting the AES key: $C_{ECC} = ECC\_Encrypt(K_{AES}, K_{ECC}^{pub})$

Combining the ciphertexts: $Final\_Ciphertext = (C_{AES}, C_{ECC})$

Decrypting the AES key: $K_{AES} = ECC\_Decrypt(C_{ECC}, K_{ECC}^{priv})$

Decrypting the message: $M = AES\_Decrypt(C_{AES}, K_{AES})$

Where, $M$ represents the plaintext message and $K_{AES}$ is the AES symmetric key.

Let $K_{ECC}^{pub}$ and $K_{ECC}^{priv}$ denote the ECC public and private keys respectively, while $C_{AES}$ and $C_{ECC}$ denote the AES and ECC ciphertexts respectively.

Assume **E (m, k)** be the **AES + ECC encryption function**, where **'m'** as the message and **'k'** as the public key.

Encrypted EHR data can be shown as:

encrypted_data = E (EHR_data, public_key)

**Step 4:** Blockchain upload:

Let B denote the Blockchain, and T represent the transaction containing the encrypted EHR data along with the hashes of the previous and present blocks. Let's denote the input data as D, the current state of the blockchain as S, and the resulting updated state as S'.

The formula for updating the blockchain state is given by:

$$S' = Hash(Encrypt(D) + Hash(S))$$

The process of transferring encrypted EHR data into a Blockchain platform can be illustrated as follows:

T = {encrypted_data, previous_block_hash, current_block_hash}

B = B U {T}

**Step 5:** Verification of the key condition:

Let K denote the pre-shared secret key and received_key is the key received from the other medical center requesting for data.

Key verification condition can be expressed as:

IF (received_key == K)

THEN receive_data ()

ELSE end_process ()

**Step 6:** ML classification using XG Boost and Light GBM:
Let denote the output labels $X$ represent the input data, $Y$ represent the input data, $f_{XG}$ denote the XG Boost classifier function, and $f_{LGBM}$ denote the Light GBM classifier function.
The classification process using ML classifiers (XG Boost and Light GBM) is demonstrated as follows:

For XG Boost:

$Y_{predicted, \ XG} = f_{XG} \ (X_{train})$

$$(XGBoost) = \sum_{i=1}^{n} L \ (y_i, y_i^{\wedge}) + \sum_{k=1}^{K} \Omega \ (f_k)$$

For Light GBM:

$Y_{predicted, \ LGBM} = f_{LGBM}(X_{train})$

$$(LightGBM) = \sum_{i=1}^{n} L \ (y_i, y_i^{\wedge}) + \sum_{k=1}^{K} \Omega \ (f_k)$$

Where, represents the number of training samples.
$L \ (y_i, y_i^{\wedge})$ denotes the loss function measuring the difference between the true label $y_i$ and the predicted label $y_i^{\wedge}$. K is the number of trees in the ensemble. $\Omega(f_k)$ signifies the regularization term penalizing the complexity of each tree.
Subsequently, accuracy is calculated using the predicted values as follows:

$Accuracy_{XG \ Boost} = \frac{TP_{XGBoost} + TN_{XGBoost}}{TP_{XGBoost} + FP_{XGBoost} + TN_{XGBoost} + FN_{XGBoost}}$

$TP_{XGBoost} = \sum_{i=1}^{N} \| \ (y_{XG \ Boost}[i] = 1 \ and \ y_{true}[i] = 1)$

$TN_{XGBoost} = \sum_{i=1}^{N} \| \ (y_{XG \ Boost}[i] = 0 \ and \ y_{true}[i] = 0)$

$FP_{XGBoost} = \sum_{i=1}^{N} \| \ (y_{XG \ Boost}[i] = 1 \ and \ y_{true}[i] = 0)$

$FN_{XGBoost} = \sum_{i=1}^{N} \| \ (y_{XG \ Boost}[i] = 0 \ and \ y_{true}[i] = 1)$

Similarly,

$Accuracy_{Light \ GBM} = \frac{TP_{LightGBM} + TN_{Light \ GBM}}{TP_{Light \ GBM} + FP_{Light \ GBM} + TN_{Light \ GBM} + FN_{Light \ GBM}}$

$TP_{Light \ GBM} = \sum_{i=1}^{N} \| \ (y_{Light \ GBM}[i] = 1 \ and \ y_{true}[i] = 1)$

$TN_{Light \ GBM} = \sum_{i=1}^{N} \| \ (y_{Light \ GBM}[i] = 0 \ and \ y_{true}[i] = 0)$

$FP_{Light \ GBM} = \sum_{i=1}^{N} \| \ (y_{Light \ GBM}[i] = 1 \ and \ y_{true}[i] = 0)$

$FN_{Light \ GBM} = \sum_{i=1}^{N} \| \ (y_{Light \ GBM}[i] = 0 \ and \ y_{true}[i] = 1)$

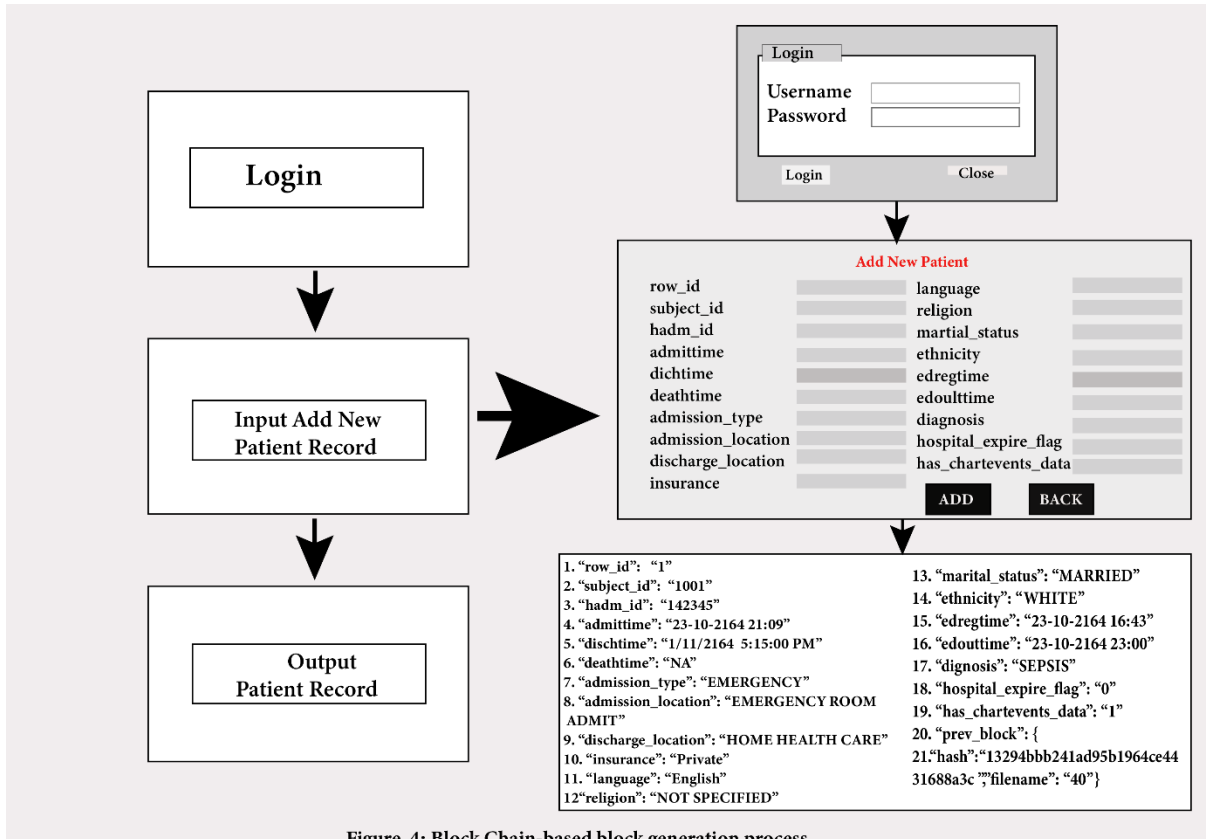$accuracy_{XG} = $ accuracy_score $(Y_{test}, Y_{predicted, \ XG})$

$$accuracy_{LGBM} = \text{accuracy\_score} \ (Y_{test}, Y_{predicted, \ XG})$$

## Result and Discussion

The efficacy of a proposed ML-driven security architecture for EHR applications was assessed in the study. This is achieved through a better encryption technique by combining multiple encryption methods and advanced machine learning classifiers, which identify and prevent possible security breaches within such vital healthcare databases. These models showed strong capabilities in detecting patterns as well as risks that could be present, making EHR systems more protected from cyber threats. This section discusses how the ML-driven security framework detects and prevents potential scams on electronic health record systems.

Figure 4 illustrates the blockchain-based method of generating blocks. It starts with a blank new patient registration form on a website that has fields for username, password, and confirm password, followed by a signup button. Then, the EHR system's Add New Patient screen is used to capture important demographic and clinical information about patients in a set P, referred to as $p_1$ to $p_n$. These are such details that are recorded on such screens as row_id, subject_id, hadm_id, admitting, dischtime, admission_type, admission_location, discharge_location, insurance details, language preferred,



Figure 4: Block Chain-based block generation process



**Figure 4. Data encryption and decryption process.**

religion, marital_status, ethnicity, edregtime, downtime, diagnosis, and hospital_expire_flag. The last figure shows a code from a medical database. The above-given code contains patient identifiers like row ID, subject ID, and admit and discharge time. This also contains other details about the patient, including his or her being admitted type, location, insurance, language, religion, marital status, and ethnicity. Further, this code has dates reflecting the time of registration of the patient at the hospital, which later led to their cancelation from hospitalization and occurrence of death since the patient was not discharged alive. At last, the code has a hash and filename.

environment. The results are shown in the figure below, such as data encryption and decryption process, encryption time vs. decryption time, and plain text vs. encrypted text.

The encryption and decryption process of data is illustrated in Figure 5. It represents a piece of text as encrypted and decrypted. An observable public key, private key, plaintext size, encryption time, encrypted size, and encrypted text are present. Thereafter, the user is prompted to provide the private key required to decrypt the text.

Figure 6 shows the number of times a piece of data is encrypted or decrypted. The x-axis displays the time
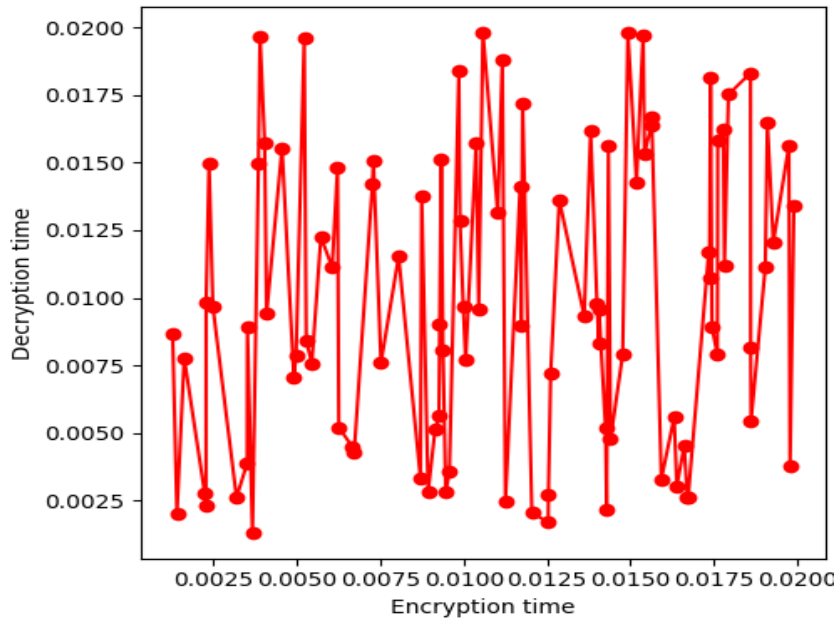


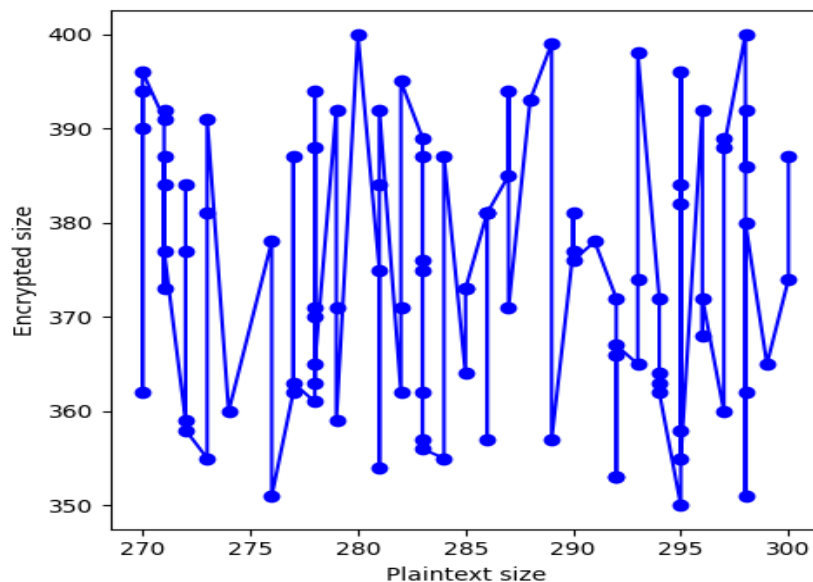**Figure 5.  Encryption time vs Decryption time.**



**Figure 6.  Plaintext size vs Encrypted size.**

These findings highlight the capability of Machine Learning classifiers to enhance the security and privacy of EHRs, hence aiding in safeguarding sensitive patient data within the ever-evolving healthcare technology

taken in seconds to encrypt, while the y-axis indicates the number of times the data is decrypted in bytes. The graph shows that the data is encrypted more times than it is decrypted. This is because encryption is a one-way

process, while decryption is a two-way process. The graph also shows that the number of times the data is encrypted or decrypted increases as the encryption time increases. This is because more complex encryption algorithms take longer to run than less complex encryption algorithms. Figure 7 (as shown in the graph) shows the size of plaintext and encrypted size. The y-axis is labeled encrypted and ranges from 350 to 400, and the other is labeled plaintext size, with a range starting at 270 and increasing to 300. These findings highlight the capacity of machine learning classifiers to enhance the security and privacy of EHRs, thereby supporting the continuous endeavors to safeguard sensitive patient data in the ever-changing field of healthcare technology. The results are shown in the confusion matrix below for the XGBoost and LightGBM models.

Figure 8 displays a confusion matrix for an XGBoost model, illustrating the model's classification performance. Rows indicate predicted labels, columns represent true labels, and each cell shows the instances where predictions differed from actual labels. Figure 9, portraying a LightGBM model's confusion matrix, reveals strong performance with most labels aligning on the diagonal, affirming the model's suitability for the task. Table 1 presents a comparative analysis of previous methodologies alongside the proposed approach, utilizing the MIMIC III dataset. The results indicate that Huang et al. (2023) achieved a 66.2% accuracy by employing the XG Boost technique. Yang et al. (2023) showcased varying accuracies, with DTC at 70.9%, RF at 73.2%, XG Boost at 79.35%, Light GBM at 79.3%, and Cat Boost at 81.7%. Tsiklidis et al. (2022) demonstrated
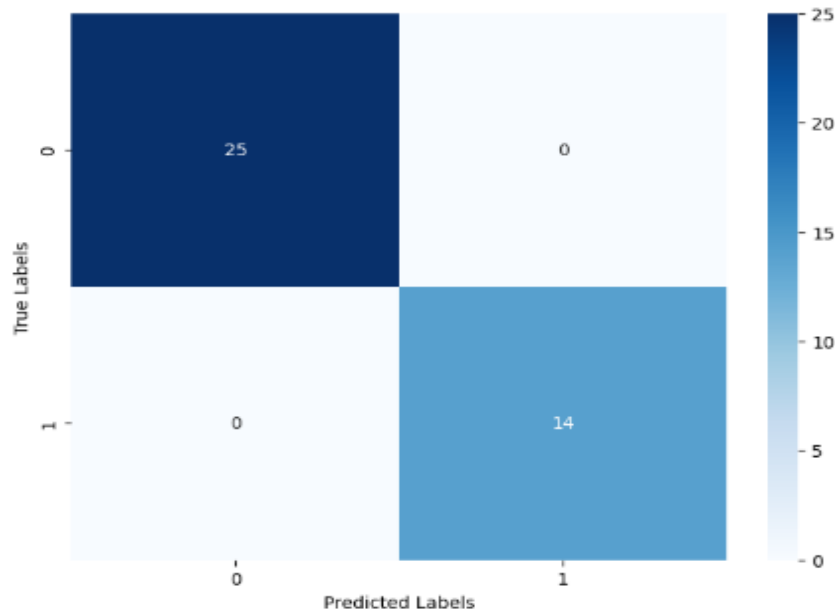


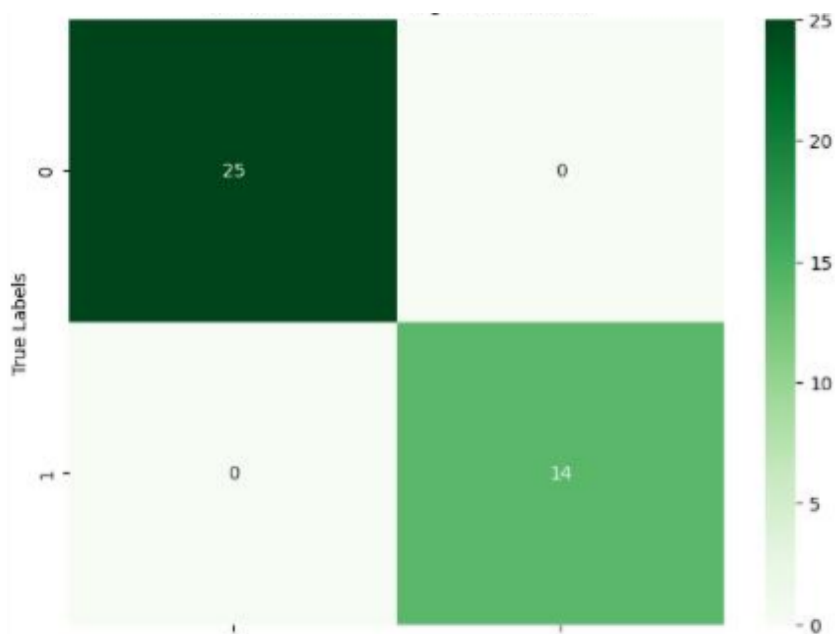**Figure 7. Confusion matrix of XG Boost model.**



**Figure 8. Light GBM model's confusion matrix.**

accuracies of 74% for SVM, 78.7% for LR, 87.2% for Gaussian Naïve Bayes, and 92.9% for GB Classifier. Hou et al. (2020) attained an 81.9% accuracy with the XG Boost technique. Notably, the proposed method achieved a remarkable 100% accuracy using both XGBoost and LightGBM models.

**Table 1. Comparative Analysis of Related Techniques.**

| Authors | Techniques | Values |
|---|---|---|
| **Huang et al. (2023)** | XGBoost | 66.2% |
| **Yang et al. (2023)** | DTC, Random Forest, XGBoost, LightGBM, CatBoost | 70.9%, 73.2%, 79.35, 79.3%, 81.7% |
| **Tsiklidis et al. (2022)** | SVM, LR, Gaussian Naïve Bayes, GB Classifier | 74%, 78.7%, 87.2%, 92.9% |
| **Hou et al. (2020)** | XGBoost | 81.9% |
| **Proposed Method** | XG Boost, Light GBM | 100%, 100% |

Figure 10 and Figure 11 shows the ROC curves of XGBoost and LightGBM. The performance of both the XGBoost and LightGBM models is exceptional, as indicated by their ROC curves.

The XGBoost model has a nearly perfect ROC curve with an AUC (Area Under the Curve) of 1.00 indicating its excellent ability to correctly classify positive and negative cases. Likewise, the ROC curve for the LightGBM model is flawless with a perfect AUC of 1.00 showcasing its flawless power to separate between two classes. Both machine learning models show outstanding binary classification performance in these situations, making them very effective for the MIMIC III dataset.

## Conclusion

Electronic health records (EHRs) incorporate a vast amount of patient information and diagnostic data, most of which are considered important health information for a person. With the advancement of technology, the emergence of advanced cyber threats has escalated, hindering health information systems' privacy and security. Due to this, privacy and security concerns
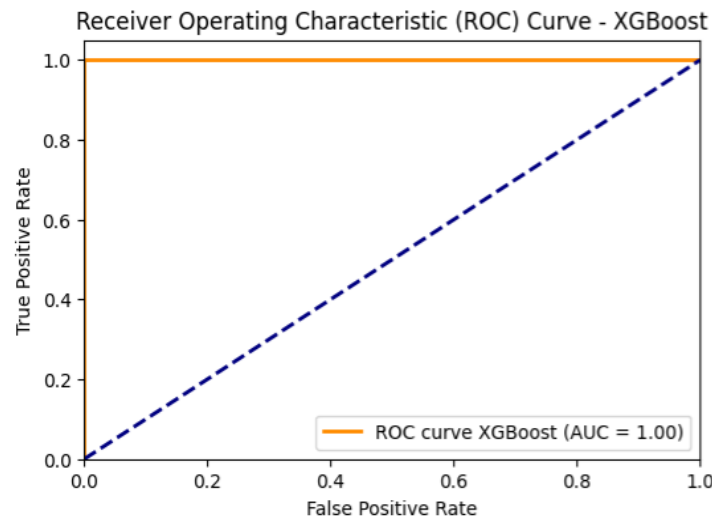


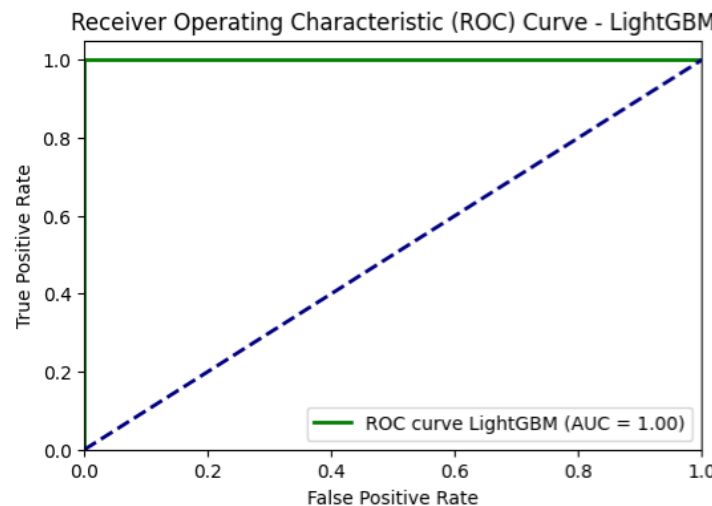**Figure 9. ROC Curve of XGBoost.**



**Figure 10. ROC Curve of LightGBM.**

present the largest and most important barrier to adopting EHRs. In conclusion, incorporating ML-driven assessments and secured EHRs would be a transformative solution for data protection in the healthcare sector. Using machine learning, blockchain and encryption algorithms to test and improve the security of EHR systems has been shown to work very well, especially when the proposed method includes XGBoost and LightGBM models. The results obtained showed that the XGBoost model had exceptional performance, with a nearly perfect ROC curve and an AUC of 1.00, thus indicating its high accuracy in classifying positive versus negative cases. As well as that, the LightGBM model had a flawless performance with a perfect ROC curve.

Furthermore, in the future, more sophisticated ML models, advanced data encryption techniques, and secure communication protocols can make this proposed model strong enough to withstand emerging threats and increase its diagnostic capabilities.

## Conflict of Interest

The authors declare no conflict of interest.

## References

Ahmadi, M., & Aslani, N. (2018). Capabilities and advantages of cloud computing in the implementation of electronic health Record. *Acta Informatica Medica, 26*(1), 24. https://doi.org/10.5455/aim.2018.26.24-28

Alam, M. U., & Rahmani, R. (2023). FedSepsis: A Federated Multi-Modal Deep Learning-Based Internet of Medical Things Application for Early Detection of Sepsis from Electronic Health Records Using Raspberry Pi and Jetson Nano Devices. *Sensors, 23*(2), 970. https://doi.org/10.3390/s23020970

Bhattacharya, P., Tanwar, S., Bodkhe, U., Tyagi, S., & Kumar, N. (2021). BINDAAS: Blockchain-Based Deep-Learning as-a-Service in Healthcare 4.0 Applications. *IEEE Transactions on Network Science and Engineering, 8*(2), 1242–1255. https://doi.org/10.1109/tnse.2019.2961932

Chami, S., & Tavakolian, K. (2019). Comparative study of Light-GBM and LSTM for early prediction of sepsis from clinical data. *Computing in Cardiology, 46*, 1-4. https://doi.org/10.22489/cinc.2019.367

Chen, I. Y., Pierson, E., Rose, S., Joshi, S., Ferryman, K., & Ghassemi, M. (2021). Ethical machine learning in healthcare. *Annual Review of Biomedical Data Science, 4*(1), 123–144. https://doi.org/10.1146/annurev-biodatasci-092820-114757

Chenthara, S., Ahmed, K., Wang, H., & Whittaker, F. (2019). Security and Privacy-Preserving challenges of e-Health solutions in cloud Computing. *IEEE Access*, 7, 74361–74382. https://doi.org/10.1109/access.2019.2919982

Corbin, C.K., Sung, L., Chattopadhyay, A., Noshad, M., Chang, A., Deresinksi, S., Baiocchi, M., & Chen, J. H. (2022). Personalized antibiograms for machine learning driven antibiotic selection. *Commun. Med., 2*, 38. https://doi.org/10.1038/s43856-022-00094-8

Esteva, A., Robicquet, A., Ramsundar, B., Kuleshov, V., DePristo, M., Chou, K., Cui, C., Corrado, G., Thrun, S., & Dean, J. (2018). A guide to deep learning in healthcare. *Nature Medicine, 25*(1), 24–29. https://doi.org/10.1038/s41591-018-0316-z

Gabriel, S. J., & Sengottuvelan, P. (2021). An Enhanced Blockchain Technology with AES Encryption Security System for Healthcare System. *2nd International Conference on Smart Electronics and Communication* (ICOSEC), *2021*, 400-405. https://doi.org/10.1109/icosec51865.2021.9591956

Hasan, M. K., Alam, M. A., Das, D., Hossain, E., & Hasan, M. (2020). Diabetes prediction using ensembling of different machine learning classifiers. *IEEE Access*, 8, 76516–76531. https://doi.org/10.1109/access.2020.2989857

Hathaliya, J. J., & Tanwar, S. (2020). An exhaustive survey on security and privacy issues in Healthcare 4.0. *Computer Communications*, 153, 311–335. https://doi.org/10.1016/j.comcom.2020.02.018

Himabindu, D. D., Pranalini, B., Kumar, M., Neethika, A., Sree N, B., C, M., B, H., & S, K. (2024). Deep CNN-based Classification of Brain MRI Images for Alzheimer's Disease Diagnosis. *International Journal of Experimental Research and Review, 41*(Spl Vol), 43-54. https://doi.org/10.52756/ijerr.2024.v41spl.004

Hirano, H., Koga, K., & Takemoto, K. (2020). Vulnerability of deep neural networks for detecting COVID-19 cases from chest X-ray images to universal adversarial attacks. *PLoS ONE, 15*(12), e0243963. https://doi.org/10.1371/journal.pone.0243963

Hou, N., Li, M., He, L., Xi, B., Wang, L., Zhang, R., Yu, Y., Sun, X., Pan, Z., & Wang, K. (2020). Predicting 30-days mortality for MIMIC-III patients with sepsis-3: a machine learning approach using XGboost. *J. Transl. Med., 18*, 462. https://doi.org/10.1186/s12967-020-02620-5

Huang, A. A., & Huang, S. Y. (2023). Dendrogram of transparent feature importance machine learning statistics to classify associations for heart failure: A reanalysis of a retrospective cohort study of the Medical Information Mart for Intensive Care III (MIMIC-III) database. *PLoS ONE*, *18*(7), e0288819.
https://doi.org/10.1371/journal.pone.0288819

Huang, X. (2019). Blockchain in Healthcare: a Patient-Centered model. *Biomedical Journal of Scientific & Technical Research*, *20*(3).
https://doi.org/10.26717/bjstr.2019.20.003448

Keerthana, B., Vamsinath, J., Kumari, C., Appaji, S. V., Rani, P. P., & Chilukuri, S. (2024). Machine Learning Techniques for Medicinal Leaf Prediction and Disease Identification. *International Journal of Experimental Research and Review, 42*, 320-327. https://doi.org/10.52756/ijerr.2024.v42.028.

Kumar, S. R., Gayathri, N., Muthuramalingam, S., Balamurugan, B., Ramesh, C., & Nallakaruppan, M. (2019). Medical big data mining and processing in e-Healthcare. In Elsevier eBooks. pp. 323–339.
https://doi.org/10.1016/b978-0-12-817356-5.00016-4

Kumari, A., Tanwar, S., Tyagi, S., & Kumar, N. (2018). Fog computing for Healthcare 4.0 environment: Opportunities and challenges. *Computers & Electrical Engineering*, *72*, 1–13.
https://doi.org/10.1016/j.compeleceng.2018.08.015

Lee, D. (2019). Effects of key value co-creation elements in the healthcare system: focusing on technology applications. *Serv Bus., 13*, 389–417.
https://doi.org/10.1007/s11628-018-00388-9

Lee, D. H., &. S. N. (2021). Application of Artificial Intelligence-Based Technologies in the Healthcare Industry: Opportunities and challenges. *ideas.repec.org*, *18*(1), 271.
https://ideas.repec.org/a/gam/jijerp/v18y2021i1p271-d473475.html

Locke, S., Bashall, A., Al-Adely, S., Moore, J., Wilson, A., & Kitchen, G. B. (2021). Natural language processing in medicine: A review. *Trends in Anaesthesia and Critical Care*, *38*, 4–9.
https://doi.org/10.1016/j.tacc.2021.02.007

Mandair, D., Tiwari, P., Simon, S., Colborn, K., & Rosenberg, M. (2020). Prediction of incident myocardial infarction using machine learning applied to harmonized electronic health record data. *BMC Med. Inform. Decis. Mak., 20*, 252 (2020). https://doi.org/10.1186/s12911-020-01268-x

Mayer, A.H., da Costa, C.A., & Righi, R. da R. (2020). Electronic health records in a Blockchain: A systematic review. *Health Informatics Journal, 26*(2), 1273-1288.
https://doi.org/10.1177/1460458219866350

McGhin, T., Choo, K. R., Liu, C. Z., & He, D. (2019). Blockchain in healthcare applications: Research challenges and opportunities. *Journal of Network and Computer Applications, 135*, 62–75.
https://doi.org/10.1016/j.jnca.2019.02.027

Mishra, V., Mishra, M., Tamrakar, A., Srikanth, T., Kumar, T., & Kumar, A. (2024). Pneumonia Detection through Deep Learning: A Comparative Exploration of Classification and Segmentation Strategies. *International Journal of Experimental Research and Review, 40*(Spl Volume), 41-55.
https://doi.org/10.52756/ijerr.2024.v40spl.004

Newaz, A. I., Sikder, A. K., Rahman, M. A., & Uluagac, A. S. (2019). HealthGuard: A Machine Learning-Based Security Framework for Smart Healthcare Systems. 019 Sixth International Conference on Social Networks Analysis, Management and Security (SNAMS), Granada, Spain. pp. 389-396.
https://doi.org/10.1109/snams.2019.8931716

Osmani, V., Li, L., Danieletto, M., Glicksberg, B., Dudley, J., & Mayora, O. (2018). Processing of Electronic Health Records using Deep Learning: A review. *arXiv (Cornell University)*.
https://doi.org/10.48550/arxiv.1804.01758

Pang, X., Forrest, C. B., Lê-Scherban, F., & Masino, A. J. (2021). Prediction of early childhood obesity with machine learning and electronic health record data. *International Journal of Medical Informatics, 150*, 104454.
https://doi.org/10.1016/j.ijmedinf.2021.104454

Qayyum, A., Qadir, J., Bilal, M., & Al-Fuqaha, A. (2020). Secure and Robust machine learning for healthcare: a survey. *IEEE Reviews in Biomedical Engineering, 14*, 156–180.
https://doi.org/10.1109/rbme.2020.3013489

Rao, K., Devi, J., Anuradha, Y., G, K., Kumar, M., & Rao, M. S. (2024). Enhancing Liver Disease Detection and Management with Advanced Machine Learning Models. *International Journal of Experimental Research and Review, 42*, 100-110. https://doi.org/10.52756/ijerr.2024.v42.009

Romeo, L., & Frontoni, E. (2021). A Unified Hierarchical XGBoost model for classifying priorities for COVID-19 vaccination campaign. *Pattern Recognition, 121*, 108197.
https://doi.org/10.1016/j.patcog.2021.108197

Shah, H., Rai, K., Singh, D., Gupta, S., BR, S. R., & Tripathi, R. C. (2023). Blockchain machine learning based personal health record data analysis with smart decentralization and security enhancement. *Research Square (Research Square)*. https://doi.org/10.21203/rs.3.rs-2653352/v1

Souri, A., Ghafour, M. Y., Ahmed, A. M., Safara, F., Yamini, A., & Hoseyninezhad, M. (2020). A new machine learning-based healthcare monitoring model for student's condition diagnosis in Internet of Things environment. *Soft Computing*, *24*(22), 17111–17121. https://doi.org/10.1007/s00500-020-05003-6

Tsiklidis, E. J., Sinno, T., & Diamond, S. L. (2022). Predicting risk for trauma patients using static and dynamic information from the MIMIC III database. *PLoS ONE*, *17*(1), e0262523. https://doi.org/10.1371/journal.pone.0262523

Vayena, E., Blasimme, A., & Cohen, I. G. (2018). Machine learning in medicine: Addressing ethical challenges. *PLoS Medicine*, *15*(11), e1002689. https://doi.org/10.1371/journal.pmed.1002689

Vignesh, R., Suja Cherukullapurath Mana, B. Keerthi Samhitha, and Jithina Jose. Integrating The Hospital Management System Based on Cloud Environment Using Hmac Algorithm. *International Journal of Computer Engineering and Technology (IJCET), 8*(4), 2017. https://iaeme.com/Home/article_id/IJCET_08_04_009

Vos, J.F.J., Boonstra, A., Kooistra, A., Seelen, M., & Offenbeek, M.V. (2020). The influence of electronic health record use on collaboration among medical specialties. *BMC Health Serv Res., 20*, 676. https://doi.org/10.1186/s12913-020-05542-6

Xiao, C., Choi, E., & Sun, J. (2018). Opportunities and challenges in developing deep learning models using electronic health records data: a systematic review. *Journal of the American Medical Informatics Association*, *25*(10), 1419–1428. https://doi.org/10.1093/jamia/ocy068

Yang, M., Hu, W., & Yan, J. (2023). Development of machine learning models for predicting acute respiratory distress syndrome : evidence from the MIMIC-III and MIMIC-IV. *Research Square (Research Square)*. https://doi.org/10.21203/rs.3.rs-3221576/v1