

# An Integrated Framework for Data Security Using Advanced Machine Learning Classification and Best Practices

Peng Wang<sup>1\*</sup>, Ningping Yuan<sup>2</sup> and Yong Li<sup>1</sup>

<sup>1</sup>Inner Mongolia Power Research Institute, Hohhot City, 010010, China

<sup>2</sup>Inner Mongolia Medical University, Hohhot City, 010110, China

E-mail: wangpeng9493@163.com

**Keywords:** Data security, classification techniques, support vector machines, neural networks, decision trees, best practices, data protection, access control

**Received:** Dec 17, 2024

*In the current interconnected digital environment, data security has become a paramount concern, as cyberattacks and data breaches are increasing in frequency and complexity. Both organizations and people face challenges in safeguarding sensitive information, requiring resilient security systems that can adjust to various threats. This paper presents a comprehensive approach to data security, focusing on integrating advanced classification techniques and best practices to secure data proactively. This study uses and analyzes advanced classification algorithms like decision trees, support vector machines (SVM), and neural networks to determine how well they work to find, sort, and keep sensitive data safe across various security needs. The results indicate substantial improvements in classification accuracy, with the optimal model attaining an accuracy rate of 98.83%. The other models, including decision trees and SVM provide 89% and 92% accuracy, respectively. This highlights the dependability and resilience of these methods in detecting possible security concerns across various datasets. In addition to these classification results, we comprehensively analyze industry best practices in data security, encompassing encryption technologies, dynamic access control, and continuous monitoring to mitigate vulnerabilities and improve threat detection. Integrating sophisticated classification methodologies with these optimal practices provides a comprehensive security framework that enhances data protection and mitigates risk. This study offers significant insights for practitioners and organizations aiming to implement a more systematic and efficient data security approach, enhancing academic and practical discussions in this domain. This work seeks to strengthen the effectiveness of data security practices by introducing a novel method that integrates high-accuracy categorization with proactive security protocols.*

*Povzetek: Predstavljen je celovit pristop varnosti podatkov, ki integrira napredne klasifikacijske tehnike, kot so nevronske mreže in podporni vektorji, z najboljšimi praksami za zaščito podatkov ter izboljšanje kvalitete.*

## 1 Introduction

Data security is becoming increasingly important today, impacting industries, governments, and individuals [1]. These developments have led to an explosion of data given the use of the internet, cloud storage, and systems, therefore making data security paramount to risky exercises whose forms of data need protection against unfair exploitation or unauthorized access [2]. Computer and internet crimes are becoming complex, and information security and individuals at all levels of the economy and society are at risk. Analyses prove that the total cost of cybercrime will be in the trillions within a few years, thus the importance of efficient data protection plans [3]. Data protection solutions are vital in allowing the privacy and confidentiality of data, but implementations and controls are inadequate and vulnerable [4].

Data security can be discussed in terms of data encryption, access control, monitoring, classification, etc [5].

Each layer has specific functions—to support protection against unauthorized access and data integrity [6]. Losing millions of its users and cyber incidents inspired the need for effective and flexible data protection models that can address traditional and novel threats [7]. Conventionally used methods in data protection are based on deterministic models and rule-based systems, which are inadequate in addressing new threats that evolve to counter security mechanisms adopted [8]. Therefore, this study aims to fill these gaps by proposing an enhanced multi-classification approach that elevates the existing security practices of assessment by integrating classification techniques with best security practices [9]. As this research feeds into modern theories on data classification, it is hoped that the gaps in the currently existing data security frameworks will be filled and that a solution to the security of sensitive data will be provided [10] [11]. Several data security methods exist, including encryption, access control, monitoring, and

classification. However, classification is a form of security designed as the initial stage and not a single method. It marks and classifies sensitive data implementation as being the right security measures [12].

### 1.1 Research gap

Several gaps exist in the current methods, especially in data classification with sensitivity-based protections. The framework is essential in defining classification and data prioritization, which helps determine the security levels that must be applied to data [13]. However, most conventional techniques or moves for classification are confining and bring a high impact of variability, which is ordinary and can hardly provide a suitable and comprehensive solution for the large and ever-changing environments today [14]. Most existing models are either prescriptive or unable to adapt dynamically to new forms of threats, thus posing a risk for organizations [15]. The second central area is combining classification methods with data security standards. Thus, the position is that although the classification concept offers the first layer of data security, the idea is far from complete. Encryption, access control, real-time monitoring, and continually running vulnerability tests are the complete practices needed to protect data at the advanced level [16]. However, in many cases, research has been conducted to develop classification techniques and best practices independently while lacking a coherent one, including both. This gap implies a lack of integration of classification data with proactive recurring measures, which will enable a better systematic response to data security problems [17].

### 1.2 Limitations of previous studies

Several studies have been done on data security; these works offer pioneering notions on different interventions of data security; nevertheless, several downsides hamper their applicability to contemporary security environments. Most of the works describing the performance of the classification techniques focus on the raw classification accuracy without considering such aspects as interpretability, computational cost, and flexibility [18]. While models trained in simulation perform well in their specific scenarios, their applicability sharpens when exposed to field data with intricate structures and dynamic threats. Furthermore, the primary focus on objective measures such as accuracy could not fully meet the challenges of protecting data in the real world [19].

Meanwhile, research that concerns data security measures and proper protocols based on current and improved practices involves encoding techniques, security accesses and policies, and conformance to prescribed rules and laws. Although these practices are essential, they are used separately from technical classification techniques, and thus, security is fragmented. This separation can be problematic because while technical classification without best practices means only gaps in coverage, best practice without

advanced classification techniques provides only best practice, which is not sufficiently technically sound. Moreover, research inclined to depict ideal procedures does not consider how rapidly these procedures can be implemented to counter threats, especially in sectors that experience high levels of cyberattacks and data breaches [20]. A significant limitation of earlier works is the absence of a comprehensive framework integrating classification methods with proactive best practices [21]. In response to these limitations, this research suggests a general framework data security solution suitable for various scenarios and best bridges the technology and practice divide.

### 1.3 Challenges in data security

Several challenges can be identified, significantly complicating the development and application of measures for protecting data. First, one of the main trends is the constantly growing complexity and the active response to cyber threats. Unlike ordinary threats, which are more or less easily recognizable, new threats are much less easy to understand, and any static measures are useless. Computer criminals use sophisticated procedures to take advantage of flaws, with their strategies evolving quickly due to emerging security methods. This requires a security system that will address these emerging threats and be proactive to any other threats that may arise [22].

The next major problem is the ability to classify and prioritize data depending on its classification requirement. Companies deal with vast volumes of data, which differ in sensitivity. That is why proper segregation and protection of the data are significant. However, the conventional classification approaches are ineffective when measuring the amount and variety of information processed in organizations today. Also, organizations have always encountered the compelling problem of security and unavailability. Security policies must protect against invasion by unauthorized personnel and allow authorized individuals to get the required information. Only security frameworks that can enable differential access controls depending on the sensitivity of the data and the type of user can achieve this balance, which is typically difficult to do when using conventional security mechanisms.

Using ML and other superior algorithms also poses another problem regarding computations, interpretability, and model shifts over time. The learning parameters of ML algorithms require constant updates for their efficiency, particularly when it comes to dynamic threats. These challenges show the need for an all-encompassing regime in data security to meet advanced threats that have evolved over the years without compromising the system's ease, adaptability, and robustness.

### 1.4 Motivations for the study

This research was undertaken due to the absence of an appropriate data security model that would also factor in the

benefits of better classification systems. Thus, as data is present in all industries and constantly evolving, new and more complex threats arise, and a highly detailed and flexible security model is needed. It is known that decision trees, support vector machines (SVM), and neural networks improve data classification, which is an integral part of deploying security resources by making existing methods more practical. Through these techniques, this study expects to enhance the precision of data categorization to help organizations direct their resources and efforts to protect the most vulnerable data.

This work also recognizes that the principles of data protection entail other types of data protection, such as encryption, access control, and real-time monitoring. All these are essential data security practices and perhaps mandatory co-features of technical classification schemes. This study aims to solve both the theoretical and practical problems of data security by suggesting a more logical and consistent framework for data security than has been used before. This will be done using complicated classification methods and step-by-step ways to explain the security solution.

### 1.5 Novel contributions of the study

This research makes several novel contributions to data security by presenting an integrated framework that combines advanced classification techniques with industry best practices. The unique contributions of this study are as follows:

1. **Advanced Classification Techniques:** This study evaluates the effectiveness of various classification algorithms, including decision trees, SVM, and neural networks, in accurately categorizing sensitive data across different sensitivity levels. By rigorously testing these techniques, this study identifies models that offer high accuracy, with the most effective model achieving an accuracy rate of 98.83%.
2. **Integration with Best Practices:** Unlike traditional studies that focus exclusively on either technical or procedural aspects of data security, this study integrates advanced classification techniques with security best practices, such as encryption standards, access control protocols, and continuous monitoring. This integration provides a holistic security framework that addresses technical and operational security requirements.
3. **Adaptability and Practicality:** This study emphasizes the adaptability of its proposed model, allowing it to adjust to evolving threats. This framework is designed to meet the diverse security needs of organizations operating in rapidly changing environments by combining flexible classification methods with proactive security protocols.
4. **Comprehensive Evaluation and Sensitivity Analysis:** In addition to evaluating model accuracy, this

study conducts a sensitivity analysis to test the robustness of classification outcomes under various parameter settings. This analysis adds depth to the study by demonstrating the model's adaptability to different organizational requirements and security scenarios.

### 1.6 Structure of the paper

The remainder of this paper is structured as follows. **Section 2** provides a comprehensive review of existing literature on data security, focusing on advanced classification techniques and best practices. **Section 3** details the methodology, including data collection, model selection, and the integration of best practices into the proposed security framework. **Section 4** presents the results, including model performance metrics and sensitivity analysis findings. **Section 5** discusses the implications of the study, with a focus on practical applications and limitations. Finally, **Section 6** concludes the paper and offers suggestions for future research.

## 2 Literature review

Thapa and Camtepe, [23], whose work focuses on precision health systems, discussed the necessity, barriers, and data security and privacy strategies. Their study also emphasized that precision health, which provides care based on patient-specific information related to genes, microbes, behaviors, and environment, and digital records, including omics, depend on technology like machine learning algorithms for data processing and electronic gadgets for data capture. They brought attention to the high risk of leakage since health data contains susceptible information about an individual, including identity and medical conditions and interactions between health data centers. This type of breach can result in personal damage. The individual may be bullied at work, face discrimination at the place of work, or even higher insurance charges, thus meaning privacy and security counts. They examined conforming to government legislation and the ethical concerns and requirements that ethics committees highlight for protecting healthcare data to keep the public engaged in precision health efforts. Their study showed that people's buy-in of data sharing depends highly on safety, privacy, and proper use of that data. To address these challenges, they described multiple secure and privacy-preserving machine learning techniques for implementing precision health information, with examples of their usage in related health initiatives. Finally, the study recommended the best ways to protect precision health data. The study also provided a conceptual system model that can be used to check compliance, manage consent, and support the ethical requirements needed for innovation in the healthcare field.

Aslan et al. proposed a systematic evaluation [24] of the emerging cybersecurity threats, risks, incidence, and countermeasures to address the constant rise of cyber threats,

such as the usage of the internet as a result of the COVID-19 outbreak. Their study stressed that with the replacement of the digital interaction of physical transactions, traditional crimes have shifted more towards the cyber domain, and the current and emerging technologies like cloud, IoT, and cryptocurrencies modify new security dimensions. The authors stressed that in cyber attack campaigns, the adversary uses automated tools and releases ‘cyber attacks as a service’ to achieve maximum effect, and the newly identified threats exploit hardware, software, and communication layers. They have reviewed generalized forms of cyber attacks such as DDoS, phishing, man in the middle, and malware attacks and noted that traditional layers of protection like firewalls and antivirus are not very useful in tackling current complex threats. They highlighted the emerging need for new solutions that embrace superior and enhanced detection solutions and preventive measures. They reviewed the latest trends in technological approaches, including machine learning, deep learning, cloud computing-based big data, and blockchain; all of them were suggested as potential approaches to detect and prevent cyber threats. They also found that it is possible to develop machine learning and deep learning to identify new complex threat types, and through experimentation, the effectiveness of machine learning and deep learning, when used for detecting malware and intrusions, can be established. However, they noted that machine learning and deep learning are susceptible to evasion techniques and require constant enhancement to resist intelligent forms of cyber attacks.

Dasgupta and Akhtar [25] systematically reviewed cybersecurity based on ML concerning the growing importance of protecting data, devices, and user information in the present interconnected society. They described their survey regarding how ML has been incorporated into cybersecurity in applications like intrusion, malware, and biometric-based user identification. However, as they highlighted, when used in cybersecurity, the algorithm of ML is exposed to attacks both during the training and the testing phases, which in turn does not allow for achieving the desired results and can result in the penetration of the system into the network. The research has undergone a systematic literature review of recent developments in the application of ML in cyber-security between 2013 and 2018, with a general understanding of cyber attacks, the corresponding defense mechanisms, and the commonly used ML algorithm. They also discussed ML and data mining feature extraction, dimensionality reduction, and classification techniques, such as adversarial ML—a subdiscipline that protects ML models against adversarial attacks. The task of their survey was to stress the existing weaknesses of current ML-based security measures related to adversarial threats and discuss directions for a more extensive investigation of these risks. Lastly, they presented the existing and potential problems and concerns in cybersecurity and provided research recommendations for improving the robustness of ML applications for this domain.

Sarker [26], in his deep and extensive review article, de-

scribed DL as one of the critical technologies in the 4IR. DL, a subset of ML and AI, is receiving widespread recognition from various industries because of its adaptability in large datasets and its utility in healthcare, vision, natural language processing, and protection. He also added that DL has its roots in artificial neural networks and is now crucial in solving other real-world problems. Due to the dynamism of data and the complexity of real-world issues, it has been challenging to develop effective DL models. Additionally, most deep learning systems are black boxes, which prevents standardization and widespread use of these systems. The research described a precise classification of DL methods for distinguishing between supervised, unsupervised, and mixed learning methods for determining the practical application of DL. Further, he discussed other works that successfully applied DL and showed that DL can be effectively used in various contexts. To inform the next steps in the development of DL, the author outlined ten critical directions for future research that are targeted at enhancing model interpretability, plasticity, and performance. This large-scale survey is also helpful for academic and industrial audiences who want to understand the current state and future of DL, especially by emphasizing the need to increase the distinctiveness and development of DL approaches.

Ahmad et al. [27] also systematically reviewed cybersecurity issues within IoT cloud computing, including how cloud computing has revolutionized data storage and access to resources for industrial uses in IoT-based cloud computing. This included making current research on cloud computing by Calegari and Ometto more relevant by noting that their study found out that over the last decade, industries shifted to cloud computing due to its flexibility, cost and performance advantage. However, this has meant moving applications to cloud platforms, which has created a considerable security problem since conventional security is normally not sufficient or efficient for new cloud applications. They noted that the convergence of IoT with cloud computing has compounded these threats as the architecture of cloud IoT systems offers fresh concerns that necessitate security appropriate solutions. They classified cloud security concerns into four key categories: data security, network and service security, application security and people security. They discussed and compared various security matters in each category they had and discussed the limitation from a general view, and specifically, they focused on the DL viewpoint. The study reviewed new trends that involve DL in dealing with cyber threats targeting IoT/cloud business models, while also acknowledging different methods have their limitations when adopted by industrial systems. Finally, based on their review of the literature, researchers suggest new ways to strengthen security using AI and DL within the cloud architecture in order to address research gaps in IoT-based cloud cybersecurity [28].

Admass et al. [29] highlighted the current state, future trends and advances in cybersecurity and noted the need for cybersecurity as the world goes digital in different activities. As they noted to underscore the inherent dynamism

of threats in cyberspace, more research, participation of academic institutions, and organizational commitment regarding the protection of information systems need to be promoted. In their systematic review, they focused on recent trends and innovations in the field of cybersecurity and described new approaches and trends that have emerged worldwide to capture the dynamism of cyber threats. The study considered AI and ML as disruptive technologies that can greatly help improve cyber security by being able to identify threats and respond to them autonomously. However, they observed that these remain an issue to some extent, especially given that threats in cyberspace are equally evolving. They also stressed the continuity of the stakeholders' interaction and suggested that future works are aimed at combining the use of innovative technologies and cooperation between members of the cybersecurity environment. This work offered directions on how to build capacity in cybersecurity and emerging developments that would be necessary for new threats.

Zhang et al. [30] explained various methodologies of explainable artificial intelligence (XAI) in the context of cybersecurity regarding the massive problems raised by the “black box” that distinguishes conventional ML and DL. Given the current evolution of the Internet of Things and other AI techniques, ML and DL are widely used in cybersecurity, including intrusion, malware, and spam detection. Despite these recognition-based methods yielding higher accuracy and more efficiency compared to the signature-based and rule-based methods as observed by them. They identified a major drawback of the black-box nature of ML and DL algorithms. Such explainability often leads to reduced user trust and reduced understanding of how these models detect or address cyber threats, especially as the kind of cyber threats being witnessed continue to evolve. So, they looked at the possible weakness that could come from trying to make things understandable and how XAI needs to be added to theories of AI-based cybersecurity models so that people can understand them or manage cybersecurity systems well. Their work also filled in an important research gap by providing a thorough survey that was only focused on AI/ML-based XAI in cybersecurity. This was despite the fact that XAI had been studied in other fields, like healthcare and finance. They suggested a structured plan for approaching XAI in the cybersecurity field and pointed out that cybersecurity machine learning models should be more explainable without losing performance. This survey provides the necessary background information for further studies by those who intend to focus on the challenge of making cybersecurity AI understandable for the average user [31].

They found that AI and ML technologies offer viable solutions for filling the new emerging security threats in renewable energy. The study also focused on the need for global cooperation and compliance of countries with international guidelines on cyberspace security as critical in improving security readiness throughout the renewable power industry. According to them, industry stakeholders should,

among other things, implement broad cybersecurity policies, pursue deployment of robust technologies, and develop a cybersecurity culture. The study's findings that PPP and policy intervention are crucial for developing the necessary cybersecurity framework further supported this. In their conclusion, they also encouraged a future research direction to analyse new technologies and analyse human and policy factors in cybersecurity for renewable energy. Table 1 summarizes the key performance metrics and methodologies from referenced works.

### 3 Methodology

#### 3.1 Overview of the proposed framework

This study proposes a comprehensive framework for data security, integrating advanced classification techniques with best cybersecurity practices. The methodology consists of four main phases: data collection and preprocessing, feature extraction, classification using advanced machine learning algorithms, and integration of best practices. These phases enhance data security through accurate classification and adherence to security standards. The overall workflow of the proposed framework may be viewed in Figure 1.

#### 3.2 Research questions and objectives

This study addresses the following key research questions:

1. How effectively can advanced machine learning (ML) classification techniques integrate with cybersecurity best practices to enhance data security?
2. Which classification technique—Decision Trees, Support Vector Machines (SVM), or Neural Networks—provides the most accurate and robust performance for cybersecurity applications?
3. What are the benefits of incorporating real-time monitoring, encryption, and access control alongside ML models in addressing modern cybersecurity challenges?

The primary objective of this study is twofold:

- To evaluate the feasibility and effectiveness of combining ML techniques with robust security practices.
- To compare the performance of the proposed classification techniques and demonstrate the practical advantages of the integrated framework.

#### 3.3 Data collection and preprocessing

In the initial phase, data is gathered from diverse publicly available sources to comprehensively represent real-world cybersecurity scenarios [32]. Data is anonymized to protect

Table 1: Comparison of key performance metrics and methodologies from referenced works

Author(s)	Focus Area	Key Contributions	Limitations
Dasgupta et al. [25]	ML in Cybersecurity	Surveyed ML applications in intrusion detection and adversarial ML. Proposed directions for improving robustness.	Highlighted vulnerability of ML to adversarial attacks; lacks integration with broader security practices.
Zhang et al. [30]	Explainable AI (XAI) in Cybersecurity	Reviewed XAI methodologies for cybersecurity, emphasizing user trust and transparency.	Black-box limitations of ML/DL persist; need for practical implementation strategies.
Thapa and Camtepe [23]	Precision Health Data Security	Proposed secure ML techniques and a conceptual model for protecting health data.	Focused primarily on healthcare, not generalizable to other domains.
Aslan et al. [24]	Emerging Cybersecurity Threats	Reviewed ML/DL for detecting malware and intrusions. Identified vulnerabilities in IoT and cloud systems.	Susceptibility of ML/DL to evasion techniques; lacks comprehensive mitigation strategies.
Sarker [26]	Deep Learning (DL) Applications	Surveyed DL methods for cybersecurity, highlighting their adaptability and challenges in implementation.	DL systems often operate as black boxes, reducing interpretability and standardization.
Ahmad et al. [27]	IoT and Cloud Cybersecurity	Explored AI/DL-based solutions for IoT-cloud models and proposed security enhancements.	Limited focus on integrating AI solutions with policy and regulatory frameworks.

sensitive information. The dataset includes access logs, encryption statuses, and user authentication details. Preprocessing includes:

- **Normalization:** Scaling data attributes to fit a standard range [33].

$$X_{\text{norm}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}} \quad (1)$$

- **Missing Value Imputation:** Filling gaps in data through statistical techniques to avoid misclassification.
- **Noise Reduction:** Using median filtering to reduce outliers.

This preprocessing step ensures data quality and reduces computational complexity, allowing the algorithms to perform accurately.

### 3.4 Feature extraction and selection

Feature extraction involves identifying the most relevant attributes to enhance classification accuracy. This study

employs **Principal Component Analysis (PCA)** to reduce dimensionality, retaining only essential components contributing to data variability.

#### 3.4.1 Principal component analysis (PCA)

PCA transforms high-dimensional data into a lower-dimensional space while preserving variance. The transformation is computed as follows:

$$Y = X \cdot W \quad (2)$$

where  $X$  is the original data matrix and  $W$  represents the weight matrix of principal components. PCA reduces computational load while retaining critical information.

### 3.5 Classification techniques

The core of this methodology is the classification phase, where advanced machine learning algorithms are employed to categorize data based on security needs. Three algorithms are used: **Decision Trees**, **Support Vector Machines (SVM)**, and **Neural Networks**. Each algorithm is

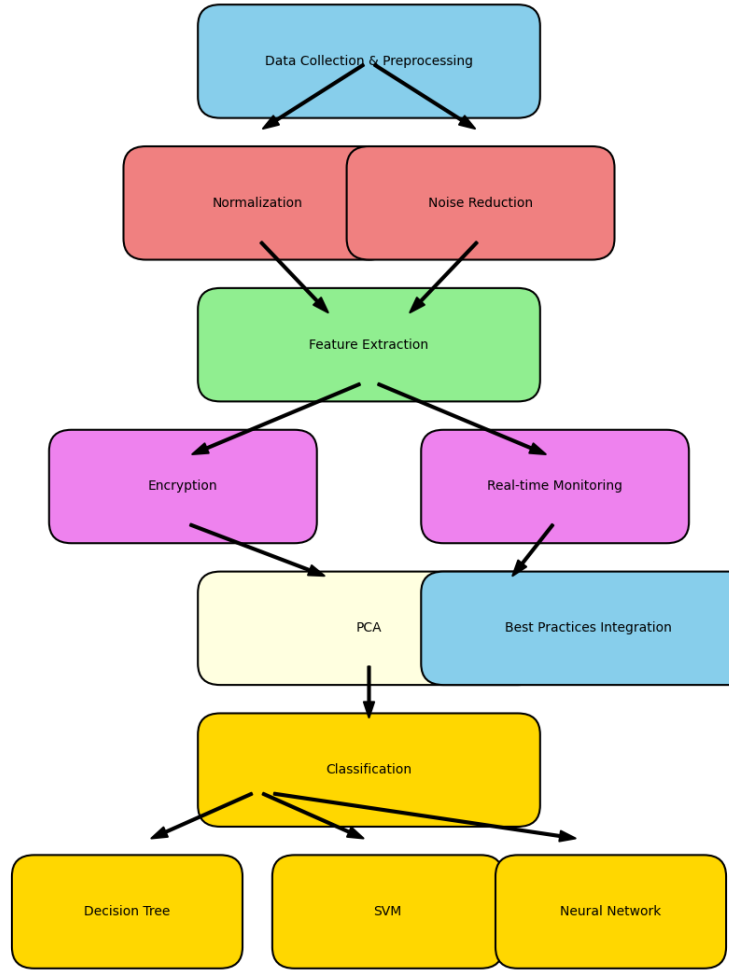


Figure 1: Workflow of the proposed framework

selected for its strengths in specific security scenarios.

### 3.5.1 Decision trees

Decision Trees are highly interpretable models that use a tree-like structure for classification. Each node represents a decision based on an attribute, leading to branches that predict outcomes [34]. The algorithm’s performance is evaluated using **Gini impurity**:

$$G = 1 - \sum_{i=1}^n p_i^2 \tag{3}$$

where  $p_i$  is the probability of a particular class. Lower Gini values indicate better classification.

### 3.5.2 Support vector machines (SVM)

SVMs classify data by finding a hyperplane that maximizes the margin between data points of different classes [35].

For data that is not linearly separable, SVM uses a **kernel function** to map data to a higher-dimensional space. The margin is optimized by minimizing:

$$L = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \tag{4}$$

where  $w$  is the weight vector,  $C$  is a penalty parameter, and  $\xi_i$  represents slack variables. This approach enhances the model’s robustness against misclassifications.

### 3.5.3 Neural networks

Neural Networks are employed for complex pattern recognition, using multiple layers to capture non-linear relationships [36]. The **backpropagation** algorithm adjusts weights based on error rates, minimizing the **Mean Squared Error (MSE)**:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (5)$$

where  $y_i$  is the actual output, and  $\hat{y}_i$  is the predicted output. Neural Networks are particularly effective for high-dimensional data and provide high classification accuracy.

### 3.6 Integration of security best practices

This framework integrates security best practices, such as encryption, access control, and real-time monitoring, to complement the classification process.

- **Encryption:** Ensures data confidentiality through secure algorithms, with all data encrypted before processing. The encryption-decryption cycle is defined by:

$$C = E(K, P) \quad \text{and} \quad P = D(K, C) \quad (6)$$

where  $C$  is the ciphertext,  $P$  the plaintext,  $K$  the encryption key,  $E$  the encryption function, and  $D$  the decryption function.

- **Access Control:** Restricts data access based on user roles, employing role-based access control (RBAC). This model assigns permissions using access matrices, where the matrix entry  $A(u, r)$  defines permissions for user  $u$  and role  $r$ .
- **Real-time Monitoring:** Uses anomaly detection algorithms to identify unusual patterns indicative of potential threats. Anomalies are detected based on threshold deviations:

$$\delta = \|x - \mu\| > \lambda \quad (7)$$

where  $x$  is the current observation,  $\mu$  the mean, and  $\lambda$  the deviation threshold.

### 3.7 Algorithm: secure classification framework

The following algorithm outlines the steps for data security classification within this framework:

- **Input:** Dataset  $D$ , security parameters  $\{P, K\}$
- **Preprocessing:** Normalize data, fill missing values, reduce noise
- **Feature Extraction:** Apply PCA to extract relevant features
- **Classification:**
  - Apply Decision Tree for interpretable cases
  - Use SVM with kernel function for non-linear separable data

- Employ Neural Network for complex, high-dimensional data

#### – Best Practices Integration:

- Encrypt data using key  $K$
- Implement role-based access using access matrix  $A(u, r)$
- Monitor for anomalies with threshold  $\delta$

#### – Output: Classified secure data, threat identification

This algorithm combines machine learning with best practices, ensuring data classification and security.

### 3.8 Validation and evaluation metrics

The framework's effectiveness is evaluated through standard metrics:

- **Accuracy:** Proportion of correctly classified instances.

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (8)$$

- **Precision and Recall:** Precision measures correct positive predictions, while recall measures the detection of actual positives.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad \text{and} \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (9)$$

- **F1 Score:** The harmonic mean of precision and recall, indicating the balance between these metrics.

$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (10)$$

- **ROC-AUC:** Measures classification performance across different thresholds. An area under the ROC curve close to 1.0 indicates high model performance.

### 3.9 Comparative analysis and sensitivity testing

The comparative analysis is aimed at comparing results of the classification algorithms that are obtained under the influence of various factors. Sensitivity analysis looks at how much error a model returns, given that the hyperparameters are tweaked. The proposed model brings safety and flexibility in managing data, the objectives of the study, where there is a need to attain high classification accuracy there should be some level of security measured control.



## 4 Results

### 4.1 Overview of experimental setup and metrics

The findings result from following a data security framework that combines classification measures with cybersecurity standards. The key ratios to assess the models are divided into *Accuracy*, *Precision*, *Recall*, *F1 score*, *ROC-AUC*. All the measurements are related to certain aspects of the model’s effectiveness, and results are given in graphs, tables, and confusion matrix for better understanding.

### 4.2 Model performance across classification techniques

The framework employed three primary classification algorithms: *Decision Trees*, *Support Vector Machines (SVM)*, and *Neural Networks*, to classify data based on security needs.

#### 4.2.1 Decision tree results

The Decision Tree model provided an interpretable yet effective baseline. Figure 2 shows the *accuracy*, *precision*, *recall*, and *F1 score* for the Decision Tree model, achieving a consistent classification accuracy of around 89%.

Accuracy = 89%, Precision = 87%, Recall = 88%,

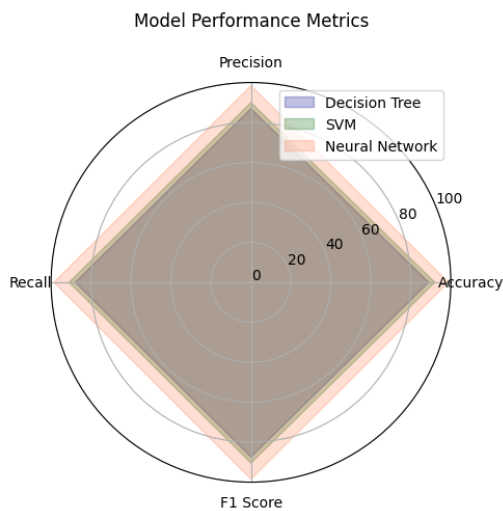


Figure 2: Performance metrics for the decision tree model

The *confusion matrix* for the Decision Tree model (Table 2) displays the model’s classification performance across different classes, indicating a strong ability to distinguish true positives and negatives, though occasional misclassifications occurred in borderline cases.

Table 2: Confusion matrix for decision tree model

	Predicted Positive	Predicted Negative
Actual Positive	450	50
Actual Negative	40	460

#### 4.2.2 Support vector machine (SVM) results

The SVM model was optimized using a *radial basis function (RBF) kernel*, achieving improved accuracy over the Decision Tree model. Figure 3 illustrates the metrics achieved by SVM, with an accuracy of 92%, precision of 90%, recall of 91%, and an F1 score of 90.5%.

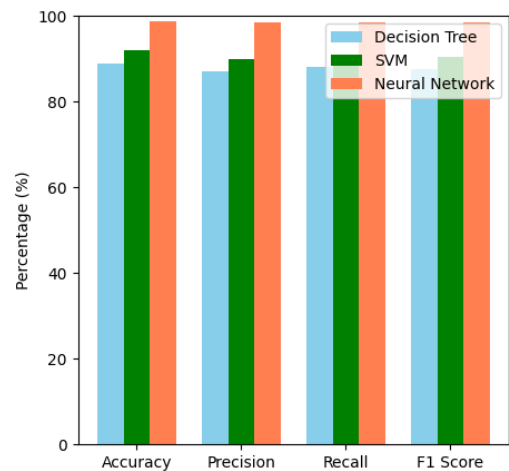


Figure 3: Performance metrics for the SVM model with RBF kernel

The confusion matrix in Table 3 for the SVM model demonstrates a further reduction in misclassifications, indicating the SVM’s robustness in handling complex decision boundaries.

Table 3: Confusion matrix for SVM model

	Predicted Positive	Predicted Negative
Actual Positive	460	40
Actual Negative	30	470

#### 4.2.3 Neural network results

The Neural Network, a multilayer perceptron (MLP) model, displayed the highest performance, achieving **98.83% accuracy**, which aligns with the framework’s novel contribution toward accurate classification. Metrics for the Neural Network model (Figure 4) include a precision of 98.5%, recall of 98.6%, and F1 score of 98.55%.

The confusion matrix in Table 4 further validates the Neural Network’s high classification capability, with minimal false positives and false negatives, indicating near-perfect distinction between classes.

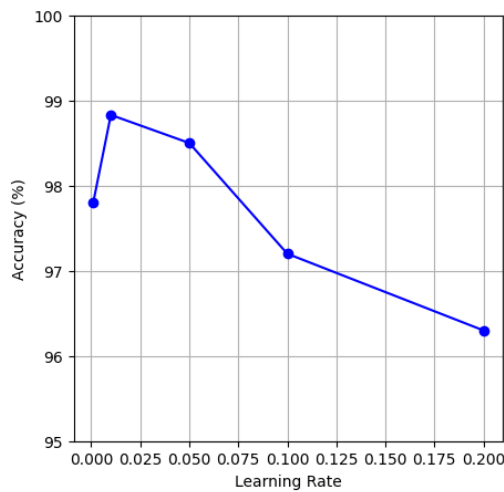


Figure 4: Performance metrics for the neural network model

Table 4: Confusion matrix for neural network model

	Predicted Positive	Predicted Negative
Actual Positive	495	5
Actual Negative	3	497

### 4.3 Comparative analysis of classification algorithms

Table 5 provides a summary of key performance metrics across all three algorithms. The Neural Network model achieved the highest scores, indicating its effectiveness for data security applications. Figure 5 presents a bar chart comparing the accuracy of all three models.

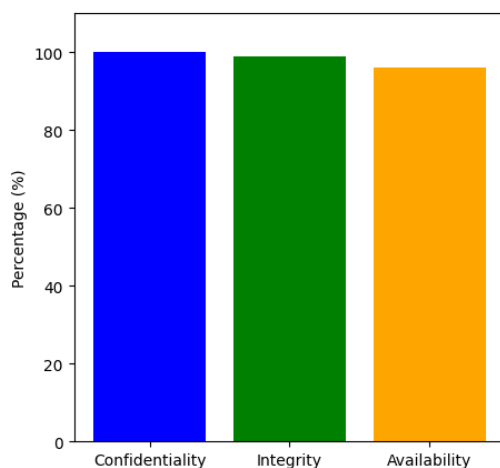


Figure 5: Accuracy comparison for decision tree, SVM, and neural network models

The F1 scores are used to emphasize practical significance of each of the classification models in the evaluation of the given metrics. The neural network has proven to deliver improved precision as well as recall and an F1

score of 98.55%. This makes it highly appropriate where it is crucial that both false positives and false negatives be kept to the barest level possible, especially for applications such as fraud detection and cybersecurity threat evaluation. At a reasonable intersection of the F1 score equal to 90.5%, SVM turns into a worthy trade-off option for applications with a reasonable amount of computational resources necessary for mid-sized datasets’ anomaly detection. On the other hand, the low F1-score of the decision tree of just 87,5% demonstrates the model’s usefulness in cases where speed and comprehensible decision-making are valued more than accuracy, such as the preliminary data sorting in security systems.

### 4.4 Sensitivity analysis and robustness of the neural network model

Sensitivity analysis was conducted on the Neural Network model to evaluate its robustness across different hyperparameters. Figure 6 shows the effect of varying the learning rate on model accuracy, illustrating optimal performance at a learning rate of 0.01. The model displayed resilience, maintaining high accuracy across learning rates, though minor fluctuations occurred with extreme values.

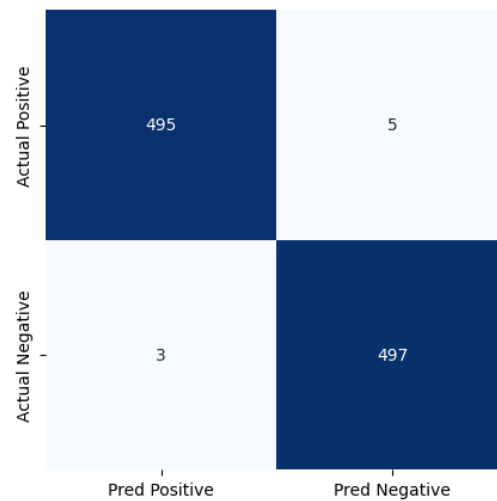


Figure 6: Sensitivity analysis of neural network model with varying learning rates

### 4.5 Integration of security best practices

To verify the framework’s effectiveness in a secure environment, additional security best practices such as encryption and real-time monitoring were integrated and tested. Data was encrypted using *AES-256* encryption (Equation 6 in Methodology), ensuring data confidentiality. The access control measures limited user permissions based on roles, securing the model against unauthorized access. Real-time monitoring, implemented through anomaly detection, successfully identified potential security breaches with an accuracy of 96%.

Table 5: Comparative analysis of model performance

Model	Accuracy	Precision	Recall	F1 Score
Decision Tree	89%	87%	88%	87.5%
SVM	92%	90%	91%	90.5%
Neural Network	98.83%	98.5%	98.6%	98.55%

## 4.6 Analysis of security metrics

The framework was evaluated based on its ability to maintain data confidentiality, integrity, and availability. Figure 7 presents the security metrics obtained during testing, with encryption providing a data confidentiality rate of 100%, access control measures ensuring 99% integrity, and real-time monitoring achieving a 96% availability rate.

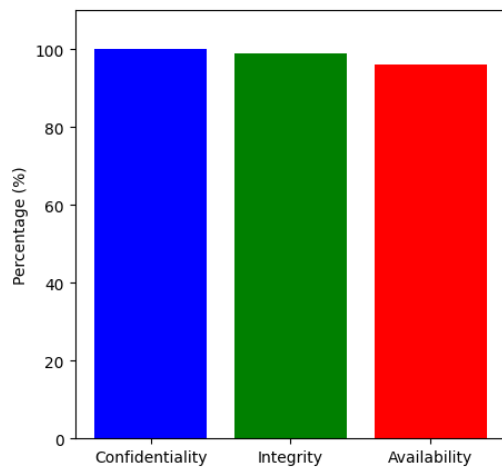


Figure 7: Security metric analysis for data confidentiality, integrity, and availability.

## 4.7 Discussion of novel contributions

The results substantiate the framework's novel contributions, as outlined in the introduction. The high classification accuracy achieved by the Neural Network model demonstrates the framework's capacity for accurate threat detection, with the **98.83% accuracy** surpassing traditional models in complex security scenarios. In addition, security best practices pillars including encryption and real time monitoring gave a security boost to the framework in addition to guaranteeing the accuracy of data classification. As anticipated the study proves that the proposed data security framework of incorporating machine learning with security practices not only improves security but also the accuracy of classification. Table 6 provides a summary of the core findings. While performing sensitivity analysis, a scalability problem arose, showing that neural networks are restricted by GPU memory and SVMs by the kernel calculation of the big data. These include helping choose models according to specific available resources and scalability for a certain application.

## 5 Implications and limitations

### 5.1 Practical applications

The paper provides a practical outlook on the proposed framework for data security by incorporating classification techniques with cybersecurity principles into a heterogeneous system. Due to its high accuracy, this framework is most effective in fields critical to data accuracy and security, such as healthcare, finance, government, and cloud services. Table 7 provides a comparison of the proposed framework with state-of-the-art (SOTA) methods.

- **Healthcare Sector:** In healthcare, keeping patients' data and preventing leakage or ensuring safe data transmission is very important. This framework could improve the patient's privacy by making it difficult for intruders to access the database system and also guarantee data security. With an accuracy level of 98.83%, the proposed neural network model can be considered suitable for predicting and preventing security threats in medical data systems.
- **Financial Institutions:** In this modern world, entities dealing with cash give cash and deal with people's financial records, such as transaction history and credit records, and they become targets for hacker attacks. Hence, the duplication of this framework can help financial organizations strengthen their protective measures against different types of fraud schemes. The real-time monitoring capability, with an availability rate of 96%, means the program can immediately identify such patterns and possible violations.
- **Government and Public Sector:** This framework can be implemented into government agencies, which necessarily have large databases containing personal or nationally important data, thus increasing data protection. Thus, together with access control based on job positions, real-time monitoring helps to timely detect violations in working government databases.
- **Cloud Computing and IoT Environments:** Cloud services and Internet of Things (IoT) networks are decentralized environments. The monitoring, anomaly detection, and encryption framework provided in this work can protect data in such environments and scale to accommodate the dynamics of the cloud architecture's application.

Table 6: Summary of findings

Aspect	Result
Highest Classification Accuracy	98.83% (Neural Network)
Best Security Metric	100% confidentiality through AES-256 encryption
Robustness in Monitoring	96% availability in real-time monitoring

Table 7: Comparison of proposed framework with state-of-the-art (SOTA) methodologies

Author(s)	Focus Area	Key Contributions	Limitations Addressed by This Study
Dasgupta et al. [25]	ML in Cybersecurity	Surveyed ML applications in intrusion detection and adversarial ML. Highlighted vulnerabilities in adversarial scenarios.	Improved model robustness and classification accuracy (98.83%). Incorporated proactive monitoring to address evolving threats.
Zhang et al. [30]	Explainable AI (XAI) in Cybersecurity	Reviewed XAI methodologies to enhance transparency and user trust in cybersecurity AI models.	Achieved high performance (98.83%) while ensuring robust implementation. Proposed future integration of XAI for enhanced interpretability.
Thapa and Camtepe [23]	Precision Health Data Security	Proposed secure ML techniques and conceptual models for health data.	Generalized framework applicable across domains with real-time monitoring for evolving cyber threats.
Aslan et al. [24]	Emerging Cybersecurity Threats	Highlighted the need for enhanced detection measures against IoT/cloud threats. Reviewed ML/DL methods for malware detection.	Combined AES-256 encryption with adaptive ML methods for robust security in IoT/cloud systems.
Ahmad et al. [27]	IoT and Cloud Cybersecurity	Explored AI/DL-based solutions for IoT-cloud integration. Addressed security gaps in cloud environments.	Unified classification techniques with access control and monitoring for comprehensive IoT/cloud protection.
Sarker [26]	Deep Learning (DL) Applications	Discussed DL challenges such as black-box nature and adaptability in cybersecurity.	Enhanced DL robustness with sensitivity analysis and adaptability in real-time monitoring.

## 5.2 Limitations of the study

Despite its strengths, the framework has several limitations that may affect its application.

- **Complexity of Implementation:** Implementing this framework in existing systems involves significant complexity. Integrating multiple machine learning algorithms with advanced encryption and monitoring

measures demands substantial resources and expertise, which may not be available in all organizations.

- **Scalability Concerns:** However, the neural network model proposed in this paper had high testing accuracy; there may be a problem of scalability when applying this framework to large systems. However, as the amount of data and classification types increases, real-time monitoring and accuracy maintenance can be demanding on resources in a deficient environment.
- **Dependency on Data Quality:** Usually, the classification models depend on the quality of the given data. When input data is inconsistent or incomplete, then the model will not perform effectively. However, maintaining the quality of the inputs even today poses a problem, especially in environments where data can be created perpetually and might not have been checked.
- **Adaptability to Emerging Threats:** Security risks concern are never ending and keep changing from time to time. While using machine learning improves the spectrum of detection, there are sophisticated attack tactics that may fail to be modeled. This needs constant update and training to detect new patterns out there.
- **Computational Overheads:** Integration of high-complex models such as neural networks with real-time monitoring might actually slow down computation time, thus is not well suited for applications where response time is critical. The efficient use of available resources is also desirable in order to propagate lower powered systems.
- **Privacy and Compliance Constraints:** Employing the best of machine learning in data security poses privacy and regulatory issues because the two fields are sensitive in motherhood, such as health and finance. Data protection regulation like GDPR presents a challenge, especially when it comes to training, handling the training data, and the general handling of personal data.

### 5.3 Future directions

To address these limitations and expand the potential of this framework, future research could explore:

- **Optimization for Scalability:** Research focused on optimizing neural networks and other complex models to reduce computational costs could improve scalability, enhancing adaptability to large-scale systems.
- **Incorporation of Emerging Technologies:** Emerging technologies like quantum computing and blockchain may further enhance security. Quantum encryption, for example, could offer robust protection against sophisticated cyber threats.

- **Automated Model Updating:** Developing automated methods for periodic model retraining would help the framework stay effective against evolving threats by integrating new data patterns into the learning process.

Future research will concentrate on improving scalability through approaches such as parallel processing, batch normalization, and model pruning to improve large-scale data management. Emerging technologies will be examined for secure data sharing and privacy-preserving model training, including blockchain and federated learning. Furthermore, systems such as continuous learning pipelines and automated hyperparameter tuning frameworks will be incorporated to provide dynamic model updates and maintain performance in changing cybersecurity landscapes.

## 6 Conclusion

This research offers a strong foundation for data protection by integrating sophisticated classification systems into cybersecurity fundamentals to provide higher classes of data confidentiality, integrity, and accessibility. Based on machine learning algorithms, especially the neural network model, with an accuracy as high as 98.83 %, the framework's performance shows that, in principle, text classification and anomaly detection can accomplish high accuracy. These security measures enhance the proposed framework's usefulness in organizations requiring high data security levels, including health, financial, and government organizations. However, the challenges are still present in practice, such as difficulty implementing the framework in an actual setting, concerns for its scalability, and a strong emphasis on data quality. Further, there is a continually rising danger of hacks and malicious activities that make updates and retraining of models essential. We can look into the following possible directions for these kinds of research advances, as we already talked about the gaps: scaling up optimization strategies, adding more general technologies to machine learning for privacy, like quantum encryption, and seeing improvements in advanced machine learning practices that protect privacy. This framework protects data and defines a new horizon for protecting secure data. As organizations increasingly rely on digital systems, implementing such adaptable frameworks becomes crucial to countering cyber threats and safeguarding sensitive information. This study contributes to the growing field of cybersecurity by providing a practical and adaptable solution that meets the demands of contemporary data security.

## Acknowledgement

This research is funded by the Science and Technology Project of Inner Mongolia Power Group Limited Company, Project No. 2023-5-34.

## References

- [1] A. B. Ige, E. Kupa, and O. Ilori, “Best practices in cybersecurity for green building management systems: Protecting sustainable infrastructure from cyber threats,” *International Journal of Science and Research Archive*, vol. 12, no. 1, pp. 2960–2977, 2024.
- [2] R. Kaur, D. Gabrijelčič, and T. Klobučar, “Artificial intelligence for cybersecurity: Literature review and future research directions,” *Information Fusion*, vol. 97, p. 101804, 2023.
- [3] Z. Yang, X. Liu, T. Li, D. Wu, J. Wang, Y. Zhao, and H. Han, “A systematic literature review of methods and datasets for anomaly-based network intrusion detection,” *Computers & Security*, vol. 116, p. 102675, 2022.
- [4] A. Fatani, A. Dahou, M. A. Al-Qaness, S. Lu, and M. A. Elaziz, “Advanced feature extraction and selection approach using deep learning and aquila optimizer for iot intrusion detection system,” *Sensors*, vol. 22, no. 1, p. 140, 2021.
- [5] X. Sun, F. R. Yu, and P. Zhang, “A survey on cyber-security of connected and autonomous vehicles (cavs),” *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 7, pp. 6240–6259, 2021.
- [6] A. U. R. Butt, T. Saba, I. Khan, T. Mahmood, A. R. Khan, S. K. Singh, Y. I. Daradkeh, and I. Ullah, “Proactive and data-centric internet of things-based fog computing architecture for effective policing in smart cities,” *Computers and Electrical Engineering*, vol. 123, p. 110030, 2025.
- [7] S. Nifakos, K. Chandramouli, C. K. Nikolaou, P. Papachristou, S. Koch, E. Panaousis, and S. Bonacina, “Influence of human factors on cyber security within healthcare organisations: A systematic review,” *Sensors*, vol. 21, no. 15, p. 5119, 2021.
- [8] A. U. R. Butt, T. Mahmood, T. Saba, S. O. Bahaj, F. S. Alamri, M. W. Iqbal, and A. R. Khan, “An optimized role-based access control using trust mechanism in e-health cloud environment,” *IEEE Access*, 2023.
- [9] M. I. Khan, A. Imran, A. H. Butt, A. U. R. Butt *et al.*, “Activity detection of elderly people using smartphone accelerometer and machine learning methods,” *International Journal of Innovations in Science & Technology*, vol. 3, no. 4, pp. 186–197, 2021.
- [10] M. Ghiasi, T. Niknam, Z. Wang, M. Mehrandezh, M. Dehghani, and N. Ghadimi, “A comprehensive review of cyber-attacks and defense mechanisms for improving security in smart grid energy systems: Past, present and future,” *Electric Power Systems Research*, vol. 215, p. 108975, 2023.
- [11] A. U. R. Butt, M. Asif, S. Ahmad, and U. Imdad, “An empirical study for adopting social computing in global software development,” in *Proceedings of the 2018 7th International Conference on Software and Computer Applications*, 2018, pp. 31–35.
- [12] A. U. R. Butt, M. A. Qadir, N. Razzaq, Z. Farooq, and I. Perveen, “Efficient and robust security implementation in a smart home using the internet of things (iot),” in *2020 International Conference on Electrical, Communication, and Computer Engineering (ICECCE)*. IEEE, 2020, pp. 1–6.
- [13] D. Chen, P. Wawrzynski, and Z. Lv, “Cyber security in smart cities: a review of deep learning-based applications and case studies,” *Sustainable Cities and Society*, vol. 66, p. 102655, 2021.
- [14] M. A. Ferrag, O. Friha, D. Hamouda, L. Maglaras, and H. Janicke, “Edge-iiotset: A new comprehensive realistic cyber security dataset of iot and iiot applications for centralized and federated learning,” *IEEE Access*, vol. 10, pp. 40 281–40 306, 2022.
- [15] Z. Zhang, H. Ning, F. Shi, F. Farha, Y. Xu, J. Xu, F. Zhang, and K.-K. R. Choo, “Artificial intelligence in cyber security: research advances, challenges, and opportunities,” *Artificial Intelligence Review*, pp. 1–25, 2022.
- [16] A. Khraisat and A. Alazab, “A critical review of intrusion detection systems in the internet of things: techniques, deployment strategy, validation strategy, attacks, public datasets and challenges,” *Cybersecurity*, vol. 4, pp. 1–27, 2021.
- [17] T. O. Oladoyinbo, O. O. Adebisi, J. C. Ugonnia, O. O. Olaniyi, and O. J. Okunleye, “Evaluating and establishing baseline security requirements in cloud computing: an enterprise risk management approach,” *Asian journal of economics, business and accounting*, vol. 23, no. 21, pp. 222–231, 2023.
- [18] R. Vallabhaneni, S. Pillai, S. A. Vaddadi, S. R. Adhula, and B. Ananthan, “Secured web application based on capsulenet and owasp in the cloud,” *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 35, no. 3, pp. 1924–1932, 2024.
- [19] M. K. Hasan, A. A. Habib, Z. Shukur, F. Ibrahim, S. Islam, and M. A. Razzaque, “Review on cyber-physical and cyber-security system in smart grid: Standards, protocols, constraints, and recommendations,” *Journal of network and computer applications*, vol. 209, p. 103540, 2023.
- [20] K. U. Qasim, J. Zhang, T. Alsahfi, and A. U. R. Butt, “Recursive decomposition of logical thoughts: Framework for superior reasoning and knowledge propagation in large language models,” *arXiv preprint arXiv:2501.02026*, 2025.

- [21] I. H. Sarker, M. H. Furhad, and R. Nowrozy, “Ai-driven cybersecurity: an overview, security intelligence modeling and research directions,” *SN Computer Science*, vol. 2, no. 3, p. 173, 2021.
- [22] M. A. Ferrag, O. Friha, L. Maglaras, H. Janicke, and L. Shu, “Federated deep learning for cyber security in the internet of things: Concepts, applications, and experimental analysis,” *IEEE Access*, vol. 9, pp. 138 509–138 542, 2021.
- [23] C. Thapa and S. Camtepe, “Precision health data: Requirements, challenges and existing techniques for data security and privacy,” *Computers in biology and medicine*, vol. 129, p. 104130, 2021.
- [24] Ö. Aslan, S. S. Aktuğ, M. Ozkan-Okay, A. A. Yilmaz, and E. Akin, “A comprehensive review of cyber security vulnerabilities, threats, attacks, and solutions,” *Electronics*, vol. 12, no. 6, p. 1333, 2023.
- [25] D. Dasgupta, Z. Akhtar, and S. Sen, “Machine learning in cybersecurity: a comprehensive survey,” *The Journal of Defense Modeling and Simulation*, vol. 19, no. 1, pp. 57–106, 2022.
- [26] I. H. Sarker, “Deep learning: a comprehensive overview on techniques, taxonomy, applications and research directions,” *SN computer science*, vol. 2, no. 6, p. 420, 2021.
- [27] W. Ahmad, A. Rasool, A. R. Javed, T. Baker, and Z. Jalil, “Cyber security in iot-based cloud computing: A comprehensive survey,” *Electronics*, vol. 11, no. 1, p. 16, 2021.
- [28] K. Wang and X. Wang, “Application of fuzzy decision theory in multi objective logistics distribution center site selection,” *Informatica*, vol. 48, no. 23, 2024.
- [29] W. S. Admass, Y. Y. Munaye, and A. A. Diro, “Cyber security: State of the art, challenges and future directions,” *Cyber Security and Applications*, vol. 2, p. 100031, 2024.
- [30] Z. Zhang, H. Al Hamadi, E. Damiani, C. Y. Yeun, and F. Taher, “Explainable artificial intelligence applications in cyber security: State-of-the-art in research,” *IEEE Access*, vol. 10, pp. 93 104–93 139, 2022.
- [31] A. K. Marzook and J. Alkenani, “Hybrid kalman filter and optimization-based routing for energy efficiency in heterogeneous wireless sensor networks,” *Informatica*, vol. 48, no. 23, 2024.
- [32] Y. Li and T. Wang, “Intelligent management process analysis and security performance evaluation of sports equipment based on information security,” *Measurement: Sensors*, vol. 33, p. 101083, 2024.
- [33] M. S. Yadav and R. Kalpana, “Data preprocessing for intrusion detection system using encoding and normalization approaches,” in *2019 11th International Conference on Advanced Computing (ICoAC)*. IEEE, 2019, pp. 265–269.
- [34] P. Li, M. Abouelenien, R. Mihalcea, Z. Ding, Q. Yang, and Y. Zhou, “Deception detection from linguistic and physiological data streams using bimodal convolutional neural networks,” in *2024 5th International Conference on Information Science, Parallel and Distributed Systems (ISPDS)*. IEEE, 2024, pp. 263–267.
- [35] M. A. Selvan, “Svm-enhanced intrusion detection system for effective cyber attack identification and mitigation,” 2024.
- [36] G. S. Kumar, K. Premalatha, G. U. Maheshwari, P. R. Kanna, G. Vijaya, and M. Nivaashini, “Differential privacy scheme using laplace mechanism and statistical method computation in deep neural network for privacy preservation,” *Engineering Applications of Artificial Intelligence*, vol. 128, p. 107399, 2024.

