# Enhancing Cloud Security: Artificial Intelligence-based Data Classification Model for Cloud Computing

**[1]Bhuvana Jayabalan, [2]Vaibhav Srivastav, [3]Dr. Poonam Singh, [4]Prof. Awakash Mishra, [5]Savinder Kaur**

**Abstract**: Cloud computing (CC) is the Internet-based delivery of computer services, including data retention, processing and programmers. This permits users to develop and improve their electronic devices by providing them with instant utilization of communal assets. Categorizing data using CC is important due to organizes and protects information according to its level of sensitivity. By developing Intelligent Rat Swarm Optimized Adaptive Boosting (IRO-Adaboost), an innovative AI-based data classification approach, we hope to enhance CC environments' security. In order to train our proposed data categorization method, that we initially gathered a collection of data involving various types of data from numerous organizations. The Box-Cox Transformation (BCT) procedures are used for processing the raw data that has been obtained. We employed the Term Frequency-Inverse Document Frequency (TF-IDF) method for extracting useful features for the data that is analyzed. Our suggested approach uses swarm intelligence based on rat behaviour to enhance the performance of the Adaboost algorithm. To evaluate the proposed IRO-Adaboost technique to different standard methods, a number of metrics are employed in the outcome assessment stage, including sensitivity (92%), accuracy (96%), False Negative Rate (FNR-0.2064), False Positive Rate (FPR-0.05), and specificity (95%). The experiment's findings indicate that the recommended IRO-Adaboost strategy worked better than other conventional approaches to increase security in a cloud computing environment.

*Keywords:* Artificial Intelligence (AI), Cloud Computing, Cloud security, Data classification, Intelligent Rat swarm Optimized Adaptive Boosting (IRO-Adaboost).

## 1. Introduction

A service distribution model known as cloud computing gives customers access to resources including computers, networks, storage, software and programming frameworks as needed and charges them based on cloud computing utilized. Data classification is the procedure of choosing a threshold degree and generating many data categories [1]. Whether information is being made, altered, stored or transferred, it is a necessary activity at every stage. The data categories specify security is required and how much it is worth in terms of business assets. The different characteristics are used to classify the data. Certain classifications of data, based on the degree of risk involved in their disclosure [2]. Data

classified as restricted, secure in the organization, highly classified, or controlled. Several categories of the data are based on the generation process, user specifics and patterns of usage. Concerns about the accessibility of data are essential and crucial for any company in cloud computing. Connecting to the cloud raises safety concerns over data [3]. Data in the cloud should be protected from hostility and catastrophes. Cloud providers need to be informed and take the appropriate action to guarantee that information is constantly accessible. Data classification is employed to predict the unclassified data class. Data mining employs specialized approaches to identify previously undiscovered valid correlations and patterns in the data collection [4]. The techniques include factual models, numerical computations, data assessment, prediction and factual modeling. In cloud computing environments, data assets can have a significant impact based on the service delivery and business models [5]. In this instance, the study identified a set of attributes predicated on cloud data in security specifications to provide security based on use and the use of controlling access in the cloud. One of the dangers is the possibility of a hacker or cracker defending a network against attackers [6]. The objective of the study extends to the AI-driven data categorization approach that works to strengthen the preservation, accessibility, and integrity of essential data assets by adjusting in cloud computing changes the security concerns and dangers.

[1]*Associate Professor, Department of Computer Science and Information Technology, Jain (Deemed to be University), Bangalore, Karnataka, India, Email Id- j.bhuvana@jainuniversity.ac.in, Orcid Id- 0000-0002-8372-6311*

[2]*Assistant Professor, Department of Computer Science & Engineering, Vivekananda Global University, Jaipur vaibhav.srivastav@vgu.ac.in, Orcid Id- 0009-0000-4136-6314*

[3]*Associate Professor, Department of ISME, ATLAS SkillTech University, Mumbai, Maharashtra, India, Email Id- poonam.singh@atlasuniversity.edu.in, Orcid Id- 0009-0006-6491-4962*

[4]*Professor, Maharishi School of Engineering and Technolgy, Maharishi University of Information Technology, Lucknow, India -226036, Email Id- awakash.mishra@muit.in, Orcid Id- 0009-0002-5447-7499*

[5]*Centre of Research Impact and Outcome, Chitkara University, Rajpura-140417, Punjab, India, Email ID- savinder.kaur.orp@chitkara.edu.in, Orcid Ids- https://orcid.org/0009-0004-6109-0772*

**Key contribution**

- The report highlights the importance of cloud computing as a platform for dispersing computer services and emphasizes, it enables the provision of shared resources with streaming services.

- The BCT technique normalizes data distribution, while the TF-IDF method extracts features from processed data. AI-based data categorization model designed to improve security in cloud computing settings.

- The experiment reveals that the IRO-Adaboost strategy outperforms than other traditional methods in enhancing security in cloud computing settings.

## 2. Related works

The paper [7] presented a novel strategy aimed at enhancing cloud service providers ' capacity to simulate customer behavior. The detection and recognition procedure was carried out using a particle swarm optimization-based probabilistic neural network (PSO-PNN). The study [8] investigated a Deep Learning (DL) that can perceive and counteract known as new threats. Thus, the massive volume of data known as Big Data (BD) produced by numerous devices may be too much for conventional analytical programmers to handle. The paper [9] proposed the Fuzzy Min-Max Neural Networks (FMMNN) of the Intrusion Detection System (IDS) based on the DL-based intrusion detection technique. An essential component of data protection is intrusion detection and a crucial element of this procedure is the capacity to accurately identify various kinds of network assaults. The paper [10] proposed an approach that completes the forecast a predetermined length of time before the expected time point, allowing enough time for task scheduling based on the expected workload. The performance of exemplary cutting-edge workload prediction techniques was first compared in the research. The article [11] used edge nodes to minimize transmission of data to integrate edges and cloud-based computing for IoT analytics of data. The study [12] explored the use of edge content delivery (ECD) technology and Machine Learning (ML) methodologies to create an integrated new predictive model-based quality inspection system for industrial production. The paper [13] presented an integrated architecture that uses Google Collaborator for computations and the Convolutional Neural Networks (CNNs) method to determine the polarity of words on the cloud. The paper [14] demonstrated the viability of an additional practical, precise, and cost-effective remote sensing image classification method by focusing on cloud computing. The study [15] provided a successful DoS assault detection technique that integrates the Opposition based

Learning (OBL) methodology to produce the Oppositional Crow Search Algorithm (OCSA), which tackles such kinds of issues. The study [16] looked at the diversity of edges and centralized cloud computing infrastructure that influences the choice of transferring locations. To generate decisions regarding offloading that are nearly optimal between the MDs, edge cloud server, and centralized cloud server, the decentralized DL-driven task offloading (DDTO) technique was recommended.

## 3. Methodology

In this study, we introduced the IRO-Adaboost for enhancing security in a cloud computing environment. The following phase represents the collection of dataset; the study used BCT for pre-processing procedures. The TF-IDF is employed for the feature extraction process. The flow of this investigation of the suggested approach is shown in Figure 1.
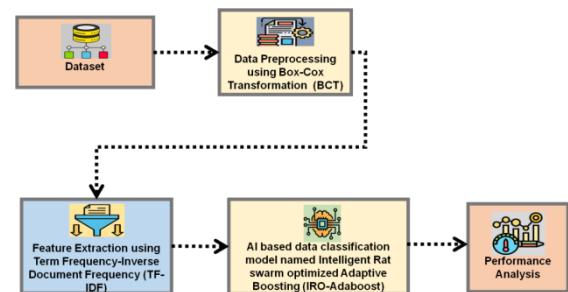


**Fig 1.** Block diagram of the suggested approach

### 3.1 Dataset

The UNSW dataset, which could be utilized in modern anomalous detection techniques which has been made available. It contains ten distinct kinds of traffic packets. Combined with regular packets, they include 9 different attack such as shell code, worms, analysis, denial-of-service, backdoor, exploits, fuzzing devices and surveillance [17]. There are two components to the UNSW-NB-15 training-set.csv training collection, which was utilized to create the approach and the UNSW-NB-15 test-set.csv testing collection, which was utilized for modeling the real-time packet that received and utilized for testing. (175,341) and (82,332) documents, respectively, make up the training and testing sets.

### 3.2 Preprocessing using Box-Cox Transformation (BCT)

In our study, the Box-Cox method is employed to change the data to make normally distributed. The variance equation also minimizes heteroscedasticity and corrects normalcy and linearity. Equation (1) illustrates the BCT data categorization model for cloud computing is calculated.

$$z_s^* = \begin{cases} \frac{z_s^\lambda - 1}{\lambda} \, if \, \lambda \neq 0 \\ In(z_s) \, if \, \lambda = 0 \end{cases} \qquad (1)$$

Where $z_s$ are real data in the present period $s$ $z_s$ determines the altered data in real time $s$. The least mean square error of residuals is represented by $\lambda$. Only time series with positive values ($z_s > 0$) can use the transformation of the preceding equation. The transformation takes on the following form if the time series data additionally include negative values for equation (2):

$$z_s^*(\lambda) = \begin{cases} \frac{(z_s + \lambda_2)^{\lambda_1} - 1}{\lambda_1} \, if \, \lambda_1 \neq 0 \\ In(z_s + \lambda_2) \, if \, \lambda_1 \neq 0 \end{cases} \qquad (2)$$

Where the transformation parameter is denoted by $\lambda_1$ and $z_s > -\lambda$ is the reason for selecting $\lambda_2$. The calculation or estimation of the $\lambda$ parameter is the process used in the Box-Cox approach for cloud computing of data classification.

## 3.3 Feature extraction using Term Frequency-Inverse Document Frequency (TF-IDF)

This approach uses the TF-IDF to demonstrate how feature extraction functions in cloud computing for data categorization. TF-IDF is employed to determine a word's importance inside a text. To perform the procedure, the average usage of a particular word in the study and that word's inverted frequency in the document over a group of texts are compounded. It is used in applications such as automated text evaluation as well as extremely beneficial in figuring out the significance of phrases in ML methods for NLP. The TF-IDF values are evaluated based on the preprocessing findings. Transforming a set of raw documents into a matrix of TF-IDF characteristics is the same goal pursued by TF-IDF applications. The TF-IDF methodology involves a systematic computation of word counts, generation of IDF values, and subsequent calculation of a TF-IDF score or group of scores. This is the sole distinction. TF-IDF is more effective if we comprehend the fundamental ideas behind TF-IDF operation in equations (3) to (4) and (5):

$$TFIDF(s,c) = TF(s,c) * IDF(s) \qquad (3)$$

$$TF(s,c) = \left(\frac{freq(s,c)}{\sum_j^m freq(s_j,c)}\right) \qquad (4)$$

$$IDF(s) = \log\left(\frac{M}{count(s)}\right) \qquad (5)$$

$freq(s,c)$ indicates the number of times the word $s$ appears in document $c$. $Count\ (s)$ The number of documents inside the compilation that have the phrase $s$. Employing TF-IDF for data categorization in cloud computing allows enterprises to set up access restrictions, impose security measures that are specific to

the value and sensitivity of the data in the cloud and efficiently organize and manage their data assets.

## 3.4 Intelligent Rat swarm Optimized Adaptive Boosting (IRO-Adaboost)

IRO-Adaboost offers an effective combination of intelligence Rat swarms and adaptive learning capabilities, making an effective solution to data categorization in cloud computing environments. IRO-Adaboost improves the quality of classification and scalability while showcasing the capacity that adjusts to the changing needs of contemporary cloud-based data mining by using the advantages between the two models.

### 3.4.1 Adaboost

The standard boosting algorithm is called an Adaboost. In the addition model, many reduced in a straight line are combined to create one effective classifier. Equation (6) provides the following expression for the addiction model:

$$G(w) = \sum_{s=1}^{S} \alpha_s g_s(s) \qquad (6)$$

Here, $g_s(s)$ represents a weak classifier, $G(w)$ is an inverse combination of weak classifiers, and $\alpha_s$ is the powerful classification's rating of the harmful classifier. The classifier created by the subsequent iteration in the progressive step-by-step method receives training using the classification algorithm from the preceding iteration. It can be stated in equation (7):

$$G(w)_n = G(w)_{n-1} + \alpha_s g_s(w) \qquad (7)$$

The least powerful classifier in the $w$-th iteration is denoted by $g_s(w)$, its weight is $\alpha_s$, and the combined result of all weak classifiers from the preceding iteration is represented by $G(w)_{n-1}$. There is an exponential loss function used in the Adaboost algorithm. Thus, the weight calculation process for each data classifier in the equation below (8) as follows:

$$\alpha_s = \frac{1}{2} in\left(\frac{1-\epsilon_s}{\epsilon_s}\right) \qquad (8)$$

Here, $s$ denotes the $s$-th iteration of the error rate. The training distribution of samples is modified according to $\alpha_s$.

$$x_{s+1} = \frac{x_s}{y_j} exp\left(-\alpha_s z g_s(s)\right) \qquad (9)$$

Where $z$ is the categorization label, $y_j$ is a normalization factor, and $x_s$ is the average weight from the preceding round of training data. The last powerful classifier is shown in equation (10):

$$e(w) = sign\left(G(w)\right) \qquad (10)$$

Where the $sign$ is the symbolic value used to translate the robust classifier's output into classification results, by computing the weighted average of the data classification, predictions are generated.

### 3.4.2 Intelligent Rat swarm Optimization (IRO)

Cloud computing platforms frequently handle remote computer resources and large-scale datasets. To optimize classification models in cloud environments, scalability and efficiency are critical IRO's capacity to parallelize computation. The newest metaheuristic algorithm, IRO was motivated by the following and attacks activities of rats. Living in swarms of both men and females, rats are a native species. In several situations, rats exhibit extremely aggressive behavior, which can cause some animals to die. The beginning placements of eligible solutions, or the positions of the rats are chosen at random in the search space in the original IRO approach as follows in equation (11).

$$w_j = w_{jmin} + rand \times (w_{jmax} - w_{jmin}), j = 1,2, \dots M \quad (11)$$

Where $w_{jmin}$ and $w_{jmax}$ are the $jth$ variable's minimum and maximum limits, correspondingly. $M$ is the total number of data. It has been proposed that this equation illustrates which rats use bait to attack and estimates the rats in updated positioning following the cloud attack as given in equation (12).

$$\vec{O}_j(w+1) = |\vec{O}_q(w) - \vec{O}| \quad (12)$$

Where $\vec{O}_j(w+1)$ describes the revised positions of $j$th rats and $\vec{O}_j(w)$ is the most ideal option discovered. $\vec{O}$ may be acquired by applying the below equation (13):

$$\vec{O} = B \times \vec{O}_j(w) + D \times (\vec{O}_q(w) - \vec{O}_j(w)) \quad (13)$$

Where $\vec{O}_j(w)$ explains the positions of $j$thrats, $B$ and $D$ are computed using the formula (14) and (15) that follows:

$$B = Q - w \times \left(\frac{Q}{Iter_{max}}\right), \ w = 1,2,3,..,Iter_{max} \quad (14)$$

$$D = 2 \times rand \quad (15)$$

In this case, the random numbers [1, 5] make up the parameter $Q$, whereas the random numbers [0, 2] make up the parameter $D$. $Iter_{max}$ is the maximum number of iterations possible in the optimization process and $w$ is the iteration that is in progress at this time. The approved standard provides high levels of assurance and the required degree of connection for the data, the outcome on various data classifiers and contrasts with proposed method.

## 4. Result and discussion

A range of experiments are presented in the following section to evaluate the suggested data categorization

initiatives. In comparison to other existing approaches like K-Nearest Neighbor (KNN) [18], Support Vector Machine (SVM) [18], and Artificial Neural Network (ANN) [18], the study's results show the effectiveness of our proposed strategy, IRO-Adaboost of cloud security architecture. To evaluate how our suggested model differs from the results of various existing approaches in terms of sensitivity, accuracy, FPR, FNR, and specificity.

### 4.1 Accuracy

The accuracy of a classifier is determined by multiplying its total number of correct predictions over the entire sample. False positive (FP), true positive (TP), false negative (FN) and true negative (TN). In mathematical terms, it is stated in equation (16):

$$Accuracy = (TN + TP)/(TN + TP + FN + FP) \quad (16)$$

To preserve the data's integrity and use, as well as the performance of data categorization and cloud security, accuracy and reliability are essential. A comparison is made between the accuracy levels of the suggested and other existing methods. The proposed method (96%) is superior to other existing methods (SVM 74%, KNN 71.82%, and ANN 75.6%). Figure 2 shows the outcome of accuracy.
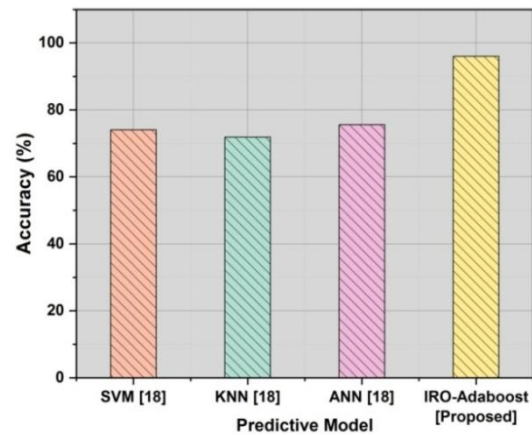


**Fig 2.** Outcome for accuracy

### 4.2 Sensitivity

The Sensitivity is the total number of data that were successfully recovered and correctly classified by several relevant documents. The higher Sensitivity indicates that the model performs better at identifying and resolving security concerns. The IRO-Adaboost value of the suggested approach attains 92%, which is higher than other methods of SVM, KNN and ANN, which attains 48%, 43.64%, and 51.19%, respectively for the sensitivity. Our proposed method has a higher sensitivity than other current approaches. Figure 3 shows the results

of Sensitivity. Sensitive data is protected in the premises in the cloud by an organization using data security technologies and procedures.

$$Sensitivity = \frac{TP}{TP+FN} \qquad (17)$$



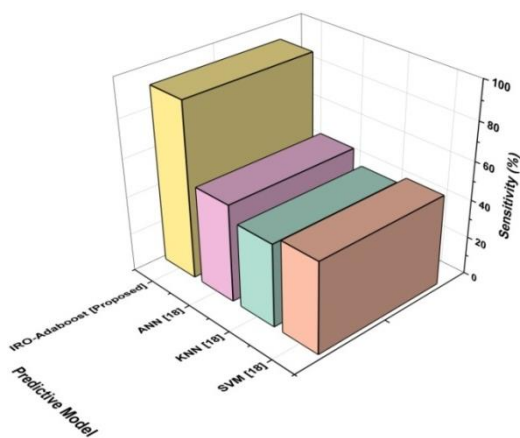**Fig 3**. Outcome for Sensitivity

## 4.3 Specificity

An important performance statistic for assessing binary classification models is specificity known as the true negative rate (TNR). It gauges the percentage of actual negatives (negative class instances) that the model accurately predicted to be negatives. The model's ability to eliminate false positives and ensure that lawful network activity is not mistakenly identified as security concerns is demonstrated by the enhanced Specificity. By increasing specificity, cloud security teams can prevent the use of effort on false alarms and instead concentrate their efforts and resources on actual security concerns. The calculation of precision is expressed using equation (18):

$$Specificity = \frac{TN}{TN+FP} \qquad (18)$$

The IRO-Adaboost value of the suggested approach, which is 95%, higher than other methods of SVM, KNN, and ANN which attains 82.66%, 81.21%, and 83.73%, respectively, for the specificity. Figure 4 shows the results of specificity. Table 1 displays the outcomes of the existing approach in addition to the suggested strategy's results.
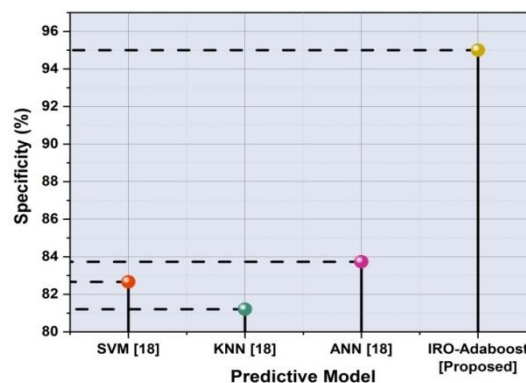


**Fig 4**. Outcome for specificity

**Table 1.** Suggested technique in comparison to the current method

| Predictive Model | Accuracy (%) | Sensitivity (%) | Specificity (%) |
|---|---|---|---|
| SVM [18] | 74 | 48 | 82.66 |
| KNN [18] | 71.82 | 43.64 | 81.21 |
| ANN [18] | 75.6 | 51.19 | 83.73 |
| IRO-Adaboost [Proposed] | 96 | 92 | 95 |

## 4.4 False Positive Rate (FPR) and False Negative Rate (FNR)

The obtained FPR values in this case such as 0.17332, 0.1878, 0.1626, and 0.05 for SVM, KNN, ANN, and IRO-Adaboost respectively. Similarly, in terms of FNR, the current model SVM achieves (0.5199), KNN (0.563), ANN (0.48805) and the suggested IRO-Adaboost achieves 0.2064. In contrast, the FPR and FNR attain the minimal value, demonstrating an effective result for the suggested approach. Table 2 and Figure 5 show the results of FPR and FNR. In contrast, the FPR and FNR attain the minimal value, demonstrating an effective result for the suggested approach. This occurs when an action attack gets detected by the IDS as acceptable.
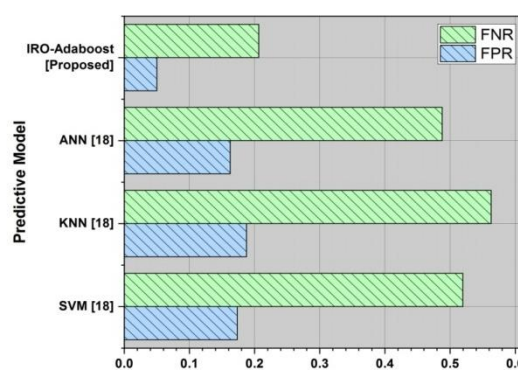


**Fig 5**. Outcome for FPR and FNR

**Table 2.** Comparison of FPR and FNR results using the proposed method

| Predictive Model | FPR | FNR |
|---|---|---|
| SVM [18] | 0.17332 | 0.5199 |
| KNN [18] | 0.1878 | 0.563 |
| ANN [18] | 0.1626 | 0.48805 |
| IRO-Adaboost [Proposed] | 0.05 | 0.2064 |

## 4.5 Discussion

In environments with complex and high-dimensional feature spaces, which are common in cloud security, KNN tends to be robust against noisy data that might provide unexpected results. Real-time threat identification and response is a critical aspect of cloud security operations, but it is hampered by KNN's inefficient processing [18]. Furthermore, by making it difficult to understand model assessments and identify any security flaws or attackers, the black-box nature of ANNs undermines the openness and dependability of cloud security systems [18]. SVMs typically need to be updated if the threat environment or data distribution changes, they are unable to handle dynamic settings and evolving threats [18]. IRO-Adaboost uses intelligent optimization algorithms and ensemble learning approaches to fortify cloud security systems against adversarial attacks and common evasion techniques used by highly competent cyber attackers.

## 5. Conclusion

The accessibility of technology has changed significantly as a result of cloud computing, which offers exceptional simulated adaptability and flexibility. The limitation could need increasing infrastructure costs and resources, optimizing model inference at certain times and reducing network latency. This study suggests cloud-based effortless text document classification and data security solutions. The IRO-Adaboost technique is to build the application's components in a technique that enables validation. Accurate data categorization is essential to its effectiveness since it guarantees appropriate organization, security and adherence to cloud security laws. Additionally, the IRO-Adaboost approach performs better than the ANN, SVM, and KNN classification algorithms in terms of specificity (95%), accuracy (96%), FPR (0.05), FNR (0.2064), and sensitivity (92%). The findings show that storing data without originally determining its security requirements is less compatible than the suggested approach.

## References

[1] Sunyaev, A., &Sunyaev, A. (2020). Cloud computing. *Internet Computing: Principles of Distributed Systems and Emerging Internet-Based Technologies*, 195-236.

[2] Ahmad, F. B., Nawaz, A., Ali, T., Kiani, A. A., & Mustafa, G. (2022). Securing cloud data: a machine learning based data categorization approach for cloud computing.

[3] Shakya, S. (2019). An efficient security framework for data migration in a cloud computing environment. *Journal of Artificial Intelligence*, *1*(01), 45-53.

[4] Afzal, S., &Kavitha, G. (2019). Load balancing in cloud computing–A hierarchical taxonomical classification. *Journal of Cloud Computing*, *8*(1), 22.

[5] Parast, F. K., Sindhav, C., Nikam, S., Yekta, H. I., Kent, K. B., &Hakak, S. (2022). Cloud computing security: A survey of service-based models. *Computers & Security*, *114*, 102580.

[6] Chadwick, D. W., Fan, W., Costantino, G., De Lemos, R., Di Cerbo, F., Herwono, I., ...& Wang, X. S. (2020). A cloud-edge-based data security architecture for sharing and analyzing cyber threat information. Future generation computer systems, 102, 710-722.

[7] Rabbani, M., Wang, Y. L., Khoshkangini, R., Jelodar, H., Zhao, R., & Hu, P. (2020). A hybrid machine learning approach for malicious behavior detection and recognition in cloud computing. *Journal of Network and Computer Applications*, *151*, 102507.

[8] Gupta, R., Tanwar, S., Tyagi, S., & Kumar, N. (2020). Machine learning models for secure data analytics: A taxonomy and threat model. *Computer Communications*, *153*, 406-440.

[9] Kumar, A., Umurzoqovich, R. S., Duong, N. D., Kanani, P., Kuppusamy, A., Praneesh, M., &Hieu, M. N. (2022). An intrusion identification and prevention for cloud computing: From the perspective of deep learning. *Optik*, *270*, 170044.

[10] Gao, J., Wang, H., & Shen, H. (2020, August). Machine learning-based workload prediction in cloud computing. In *2020 29th International Conference on computer communications and Networks (ICCCN)* (pp. 1-9). IEEE.

[11] Ghosh, A. M., &Grolinger, K. (2020). Edge-cloud computing for Internet of Things data analytics: Embedding intelligence in the edge with deep learning. *IEEE Transactions on Industrial Informatics*, *17*(3), 2191-2200.

[12] Schmitt, J., Bönig, J., Borggräfe, T., Beitinger, G., &Deuse, J. (2020). Predictive model-based quality

inspection using Machine Learning and Edge Cloud Computing. *Advanced engineering informatics*, *45*, 101101.

[13] Ghorbani, M., Bahaghighat, M., Xin, Q., &Özen, F. (2020). ConvLSTMConv network: a deep learning approach for sentiment analysis in cloud computing. *Journal of Cloud Computing*, *9*(1), 1-12.

[14] Wu, X. (2022). Big data classification of remote sensing image based on cloud computing and convolutional neural network. *Soft Computing*, 1-13.

[15] Al-Asaly, M. S., Bencherif, M. A., Alsanad, A., & Hassan, M. M. (2022). A deep learning-based resource usage prediction model for resource provisioning in an autonomic cloud computing environment. *Neural Computing and Applications*, *34*(13), 10211-10228.

[16] Wu, H., Zhang, Z., Guan, C., Wolter, K., & Xu, M. (2020). Collaborate edge and cloud computing with distributed deep learning for smart city Internet of things. *IEEE Internet of Things Journal*, *7*(9), 8099-8110.

[17] Chkirbene, Z., Erbad, A., & Hamila, R. (2019, April). A combined decision for secure cloud computing based on machine learning and past information. In 2019 IEEE Wireless Communications and Networking Conference (WCNC) (pp. 1-6). IEEE.

[18] Agarwal, A., Khari, M., & Singh, R. (2021). Detection of DDOS attack using deep learning model in cloud storage application. *Wireless Personal Communications*, 1-21.