**Module Code:** CSMRS16

**Type of Assignment: Research report -** Research Project Proposal

**Individual or Group Assignment:** Individual Report

**Student Number:** 28812368

**Date (when the work completed):** 05-Dec-2021

**Actual hrs spent for this assignment:** 12 hrs (excluding the weekly learning and practical hrs)

## I. RESEARCH BACKGROUND

In today's world, DATA has become an organization's most valuable and critical asset to the organization's business. Everything an organization does involves using DATA in some way or another. As more and more enterprise data are growing exponentially, and the DATA is being stored and processed on network-based computers, the majority of the organizations move away from traditional methods of computing and storing data to Cloud Computing. The Forbes council reported that Cloud adoption is already becoming mainstream. A vast majority of the enterprise workloads and sensitive customer data is already on the Cloud for various advantages such as ease of usage, reliability, cost efficiency, economies of scale, availability / rapid information processing. As a result, the usage of cloud services across the globe is increasing, and the Cyber Defense Report, 2021 summarizes that 38% of UK organizations deliver their services via cloud computing services, including servers, storage, databases, networking, software, analytics and intelligence - over the Internet ("the cloud"). The Cloud services are provided by notable large-scale enterprise pioneers such as Microsoft Azure, IBM, Google, and Amazon across the globe.

However, Cybersecurity is considered one of the most significant threats in a global network of large-scale enterprise organizations; a record 86% of organizations suffered from a successful cyberattack in 2021, according to the Cyberthreat Defense Report, 2021. The report demonstrates that security is one of the main concerns and issues organizations face when they look forward to embracing Cloud's merits. There are several challenges in the Cloud, such as Data Breaches, Data Loss, Insecure Access Control Points, Denial of Service (DOS) and Distributed Denial of Service (DOS). A DOS attacks are the most common type of cyberattack in Cloud Computing. Attackers create colossal traffic to force the target routers and network to consume its bandwidth and resources with the virtually created network, thus causing overloading to prevent the Cloud servers from serving or even shutting them down. Thus, the target system cannot provide services to its legitimate users, resulting in an authorized denial of service.

A Firewall only may not defend Cloud Data Centre against Cybersecurity threats adequately. For example, it cannot detect insider attacks, either over a physical or virtual network and outside attacks such as IP spoofing, DNS poisoning, Denial of Service (DoS) / Distributed Denial of service (DDoS) attacks, phishing attack, user to root attack, port scanning, attack on the virtual machine (VM). Therefore, the second line of defence control after the firewall to detect cyber-attacks is mandatory to identify various network attacks.
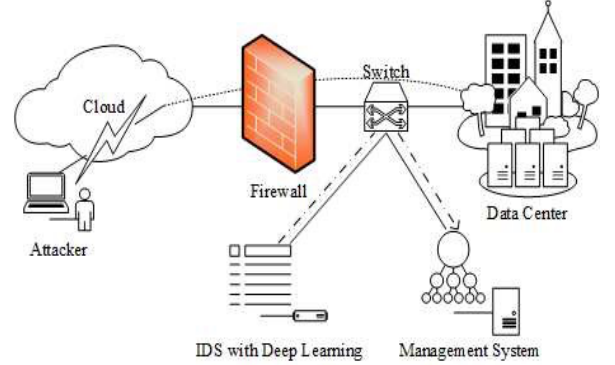


Fig. 1. Intrusion Detection System for Cloud [4] Data Center

According to the CDR Report 2021 survey, four out of five organizations globally prefer security products that feature machine learning (ML) and artificial intelligence (AI) technology. The machine learning methodologies have been widely used to identify various types of attacks, respond to them automatically, and help the network administrator take the corresponding measures and prevent intrusions. On the other hand, deep learning technologies can extract better representations from an extensive dataset and create a much better model. The Convolutional Neural Network (CNN) and Deep Reinforcement Learning (DRL) is the most representative neural networks in the field of deep learning technology research survey [1].

## II. LITERATURE REVIEW

In recent years, numerous studies using machine learning and AI methods for intrusion detection have taken primary focus and surpassed traditional intrusion detection methods.

Yang H Et al. [2] proposed a combined wireless network intrusion detection model based on the deep learning method. The combined wireless network research was studied using various methods like multi-restricted Boltzmann Machine (RBM) and the Back Propagation (BP) network, and finally, a support vendor machine is used to train the detection method. The Deep Belief Network-Support Vector machine (DBN-SVM ) IDS model is measured using benchmark datasets NSL-KDD dataset, and the accuracy is compared with other four detection

methods such as SV DBN, PCA-SVM and DBN-SVM. The paper denotes the DBN-SVM method outperformance in terms of accuracy, recall, precision, and F1 compared to other methods.

Xu C Et al. [3] propose a Deep Neural Network (DNN) model to identify the network intrusion detection and uses Gated Recurrent Units (GRU) as the central memory unit along with Multi-Layer Perceptron (MLP) and SoftMax modules to increase the performance accuracy of intrusion detection. They use the well-known KDD and NSL-KDD datasets. The experiments are compared on LSTM and GRU, and according to the experimental results, the overall detection rate was 99.42% on KDD 99 dataset and 99.31% on the NSL-KDD dataset, and the detection rates for DoS attacks were 99.98% on the KDD99 dataset, and 99.55% on NSL-KDD dataset. According to the test results, the bi-directional GRU system has better results than LSTM for intrusion detection systems.

Hizal Et al. [4] proposed a combined approach using a lightweight model based on convolutional and recurrent layers. The performance of the CNN-RNN based IDS model is tested on the NSL-KDD dataset, and the performance is evaluated using both binary and multiclass label classification and depicted an accuracy of 99.86%.

Chen L Et al. [5] propose a novel NIDS system based on a convolutional neural network. The performance of this system is measured using the CICIDS2017 raw dataset. The performance is also evaluated in terms of throughput. Furthermore, the model was compared with Support Vendor Machines (SVM) and Deep Belief Networks (DBN). The CNN based model results in 99.56% accuracy with the raw dataset compared to SVM and DBN models.

Tuan A Tang Et al. [6] proposed Gated Recurrent Unit Recurrent Neural Network (GRU-RNN) anomaly-based intrusion detection systems for Software Defined Networking (SDNs.) The performance of the GRU-RNN based IDS model is tested on the NSL-KDD dataset and achieve an accuracy of 89% with only six limited raw features.

## III. RESEARCH SCOPE

One of the most common security challenges in Cloud computing is the DoS / DDoS attack, a debasement of network protocol, which involves pumping a massive number of packets rather than the contents in the packets. The DoS attack makes computing and memory resources too busy and denies the legitimate user access to a machine, server, or host. Overthrowing network protocols such as Transmission Control Protocol (TCP) and

User Datagram Protocol (UDP) enable the attacker to thwart the services by generating vast amounts of traffic on the Cloud network, enabling them to crash instantly. The packets involved in these attacks are not easily traceable because of the large volume of the packets. For example, the TCP SYN flooding is used to subvert the TCP protocol exploiting the three-way handshake principle. Also, there are various challenges and issues in the deep learning algorithms based on intrusion or misuse detection; for example, most studies' intrusion detection approaches are analyzed with older data sets such as KDDCup99 that does not cover current and emerging data sets incorporating new DOS attack types. The most critical issue in the deep learning-based intrusion detection approach is to conduct the training process on massive datasets, and the extremely long training time of the deep learning models remains a significant problem and alarming concern. Thus, a constant need to improve the accuracy of the deep intrusion detection scheme is evident.

The key objective of the research proposal is to focus on creating the six-sigma standard intrusion detection system based on the CNN and DRL dual learning model and ensure that not approved network traffic is prevented from reaching the Cloud Servers. In addition, the dual learning model is to improve the accuracy of intrusion detection and less false positive rate, thus providing a new research method for intrusion detection.

The main objective of this research is summarized as follows

- Design and Implement the dual model intrusion detection system based on CNN and DRL deep learning networks for Cloud.

- Study the performance of the dual model CNN-DLR intrusion detection system in binary classification and multiclass classification, and different learning rate impacts on the accuracy and false positive rate on the benchmark datasets NLS-KDD.

- Compare the performance of dual model CNN-DRL IDS with other machine learning methods in binary and multiclass classification to improve intrusion detection accuracy and less false positive rate, thus providing a new research method for intrusion detection.

## IV. INTELLECTUAL CHALLENGES

The new industrial paradigm, foresees the use of Information and Communications Technology (ICTs) for remote monitoring and control critical

infrastructures, unifying the operational technologies (OTs) with the new information technologies (ITs) such as edge computing infrastructures composed of cloud servers.

Although the technical evolution is relevant for the modernization of control processes, new security challenges and risks arise within the new digital transformation era. Many of these challenges generally come from typical vulnerabilities of the cyber domains (e.g., accessible critical ports and services, irregular access control, lack of isolation measures or uncontrolled network sections, lack of auditing and accountability, irregularity in the governance processes, and incompatibilities), and they need to be properly managed 24/7 through security controls as specified by standards or specific cybersecurity frameworks (e.g., NIST). The Intellectual Challenges in the research proposal is mitigated by applying below principles;

**Choosing the Right Topic:** Developing a doable topic, reading everything on the areas related to the topic in IEEE and Elsevier etc, huge survey of the literature and developing an overarching theoretical context of the research topic.

**Choosing the Right Methodology:** The design is detailed out of the research study, and theoretical context with clear quantitative or qualitative directions. Consideration made to create a Proof of Concept in the planning process.

## V. METHODOLOGIES AND RESEARCH DESIGN

Figure 2 illustrates the high-level view of the system model, and Figure 3 details the proposed dual model system based on CNN and DRL deep learning methodologies.

The first half involves Data Pre-processing, Training, and Testing. The Data pre-processing stage comprises three layers: Data labelling, One Hot Encoding, and Normalization. The second half / last stage involves detection of intrusion and evaluation of results.

At the prediction stage, reinforcement learning is applied to anomaly-based events based on policy definition, whereas the convolutional neural network detects signature-based intrusion detection.
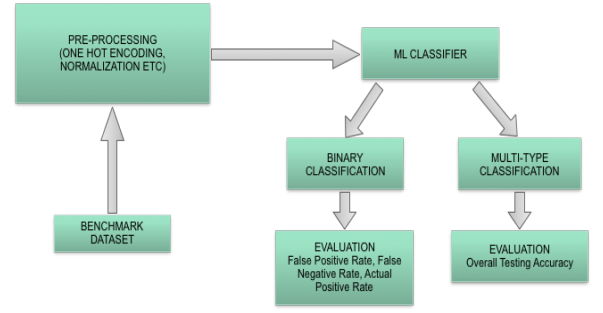


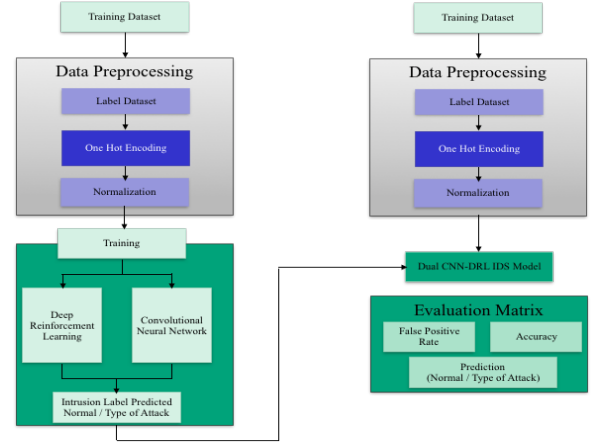Fig. 2. High-Level View of Intrusion Detection System Model



Fig. 3. Block Diagram of the Proposed Dual Model System based on CNN and DRL Methodologies

### A. Dateset Description

The NSL-KDD dataset was created in 2009 by the University of New Brunswick, Canada, and the Institute for Cybersecurity, widely used in intrusion detection research experiments.

The NSL-KDD dataset is created to overcome the inefficiencies [7], such as a large number of redundant records in the KDDCup99 dataset, but also makes the number of records reasonable in the training set and testing set so that the classifier does not favour more frequent records.

In the NSL-KDD dataset, there are 41 features and 1 class label for every traffic record, and the features include basic features (No.1- No.10), content features (No.11 - No.22), and traffic features (No.23 - No.41) as shown in Table I - Features in NSL-KDD Dataset

| Type | Features |
|------|----------|
| Nominal | Protocol_type(2),Service(3), Flag(4) |
| Binary | Land(7),logged_in(12), root_shell(14),su_attempted(15), is_host_login(21), is_guest_login(22) |
| Numeric | Duration(1),src_bytes(5), dst_bytes(6),wrong_fragment(8),urgent(9),hot(10), |

| | num_failed_logins(11), num_compromised(13), num_root(16), num_file_creations(17), num_shells(18), num_access_files(19), num_outbound_cmds(20), count(23),srv_count(24), serror_rate(25),srv_serror_rate(26), rerror_rate(27), srv_rerror_rate(28), same_srv_rate(29) diff_srv_rate(30), srv_diff_host_rate(31), dst_host_count(32), dst_host_srv_count(33), dst_host_same_srv_rate(34), dst_host_diff_srv_rate(35), dst_host_same_src_port_rate(36), dst_host_srv_diff_host_rate(37), dst_host_serror_rate(38), dst_host_srv_serror_rate(39), dst_host_rerror_rate(40), dst_host_srv_rerror_rate(41) |
|---|---|

TABLE I.      FEATURES IN NSL-KDD DATASET

According to their characteristics, attacks in the benchmark dataset are categorized into four attack types: DoS (Denial of Service attacks), R2L (Root to Local attacks), U2R (User to Root attack), and Probe (Probing attacks). In addition, the testing set has some specific attack types that disappear in the training set, which allows it to provide a more realistic theoretical basis for intrusion detection.

| DoS | R2L | U2R | Probe |
|---|---|---|---|
| back,land, neptune, pod, smurf, teardrop, mailbomb, apache2, processtable, udpstorm | ftp_write, guess_passwd, imap, multihop, phf, spy, warezclient, warezmaster, sendmail, named, snmpgetattack, snmpguess, xlock, xsnoop, worm | buffer_overf low, loadm odule, perl, rootkit ,http, tunnel, ps, sqlatta ck, xterm | Ipswee p, nmap, port, sweep, satan, mscan, saint |

TABLE II.      ATTACK TYPES IN NSL-KDD DATASET

### B. Data Preprocessing

The Data pre-processing stage comprises three layers: Data labelling, One Hot Encoding, and Normalization.

**Data Labelling:** The data is classified into standard data and four attack types (Normal, DoS,

Probe, R2L and U2R). All the data types are mapped to the data classifications that are cleansed and ready to train.

- **Denial of Service attack (DoS):** In this type of attack, the attacker successfully makes computing and memory resources too busy or denies the legitimate user access to a server is called a DoS attack

- **Remote to Local attack (R2L):** In this type of attack, the attacker gains local access as a user of the target server/machine without any account by sending vulnerable packets to a remote machine over a network.

- **Root to Local (R2L):** In this type of attack, the attacker accesses the regular user account on the target system and gain root access to exploit some vulnerability such as ransomware.

- **Probe:** This technique involves scanning a network of computers (host); an attacker gathers all necessary information about the target system or finds existing threats or vulnerabilities.

**One Hot Encoding:** To comply with the benchmark datasets' classifiers, the features are preferred to be either Boolean or Numerical but not strings of characters. Thus, the proposed model is said to use a one-hot encoder and encode the resounding features of the benchmark dataset.

**Normalization:** The next step in the data processing stage is the normalization of the features. During the normalization process, all features in the dataset are proposed to scale (Mean = 0 and Standard Deviation =1) to improve the evaluation results.

### C. Methodology

Deep Reinforcement Learning (DRL) is employed to predict the anomalies at the prediction stage, whereas Convolutional Neural Network (CNN) compares against known threats. DRL is a universal framework for learning sequential decision-making tasks. RL is defined using a Markov decision process consisting of five entities, i.e., state, action, reward, policy, and value. Here, the agent learns its behaviour based on the environment's feedback and subsequently improves its action. The basis of solving the reinforcement learning problem is to find a policy, i.e. mapping from state to action. The CNN is a feed-forward neural network with profound structure and convolution calculation.

## D. Deep Reinforcement Learning

Reinforcement Learning (RL) is a general and valuable technique for decision making and solving problems under uncertainty. DRL is an area of machine learning that applies neuron-like structure for learning tasks is the closest form of human learning because it can learn by its own experience through exploring and exploiting the unknown environment. Difference from the other ML techniques such, i.e., supervised methods learning, RL represents an agent by creating their own learning experience through interacting directly with the environment. RL uses the concept of state, reward and action as shown in Fig - 4, Overview of Deep Reinforcement Learning.
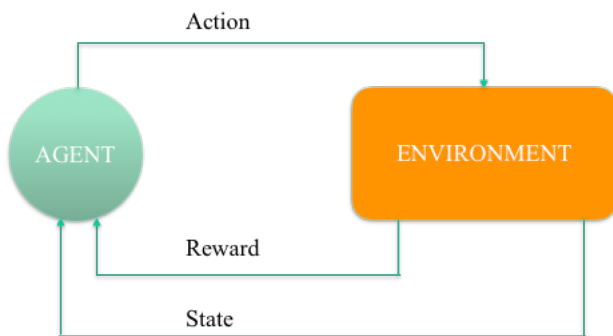


Fig. 4. Overview of Deep Reinforcement Learning

## E. Convolutional Neural Network

A CNN is another deep learning algorithm that can take in an input image, assign importance (learnable weights and biases) to various aspects/objects in the image, and differentiate one from the other. The pre-processing required in a ConvNet is much lower than other classification algorithms. While in primitive methods, filters are hand-engineered, with enough training, ConvNets can learn these filters/characteristics. The architecture of a ConvNet is analogous to that of the connectivity pattern of Neurons in the Human Brain. Individual neurons respond to stimulation only in a restricted visual field region known as the Receptive Field. A collection of such fields overlaps to cover the entire visual area. The ConvNet has three main layers: input layer, an output layer, and multiple hidden layers, consisting of a series of convolutional layers like pooling layers, fully connected layers, and normalization layers. These layers are called hidden layers since the activation function and inputs veil the final convolution. The most widely used activation function is the ReLU function.
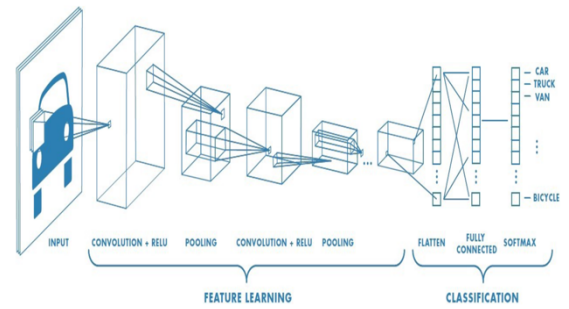


Fig. 5. A Convolutional Neural Network

## F. The Architecture Design

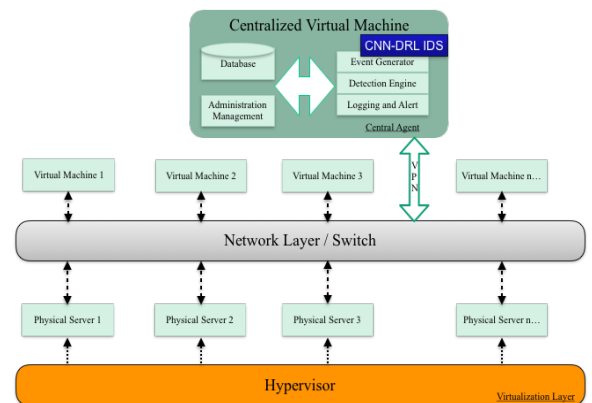Figure 6 is the proposed dual model IDS architecture based on CNN and DRL for the Cloud data centres.



Fig. 6. A CNN DRL dual model IDS Architectural Design for Cloud Data Centers

In the proposed architecture, both CNN and DLR deep learning is employed. The CNN-DRL IDS system comprises three components: Event Generator, Detection Engine, Logging and Alert. The core component is the detection engine that detects events by applying a dual learning technique. The dual model is designed to process large datasets in the distributed Cloud network simultaneously. The Detection Engine compares the network packets with the trained dataset and produces the detecting results.

The CNN-DRL specific architectural design is made up of, namely, host network (list of virtual machines hosted above network layer) and centralized agent network (combined agent and administrator management system) connected through a virtual public network (VPN). There can be two attack scenarios, i.e., (1) a single or a group of VM is compromised, (2) underlying host physical server/machine is compromised. In case 2, any virtual machines can be compromised as the attacker will have access to the hypervisor. The proposed model separates the centralized agent network from the host network by enabling communication from the host network to the agent network through a virtual public network (VPN). The VPN will prevent

the centralized agent network with IDS from being compromised by the attacker. Although, it is assumed that the attacker cannot corrupt or compromise the system logs generated within the Cloud system, which is generally true for most Cloud IDS attacks. The VMs communicate to the hypervisor through system calls, and these system calls may carry necessary log information. Log information may include when the instance was created, the VM's capabilities, the memory usage or CPU utilization of a machine, and how a VM is communicating in the network. The Logging and Alert component connects to administration management in the event of intrusion detection, and an appropriate comms mechanism is adapted (E-Mail or SMS) to the Cloud administrator.

## G. Evaluation Metric

The most critical performance indicator (Accuracy, AC) of intrusion detection is used to measure the performance of the CNN-DRL dual model IDS. In addition to the accuracy, the detection rate and false-positive rate is measured. The True Positive (TP) is the same as those correctly rejected network packets, and it denotes the number of anomaly records identified as an anomaly. The False Positive (FP) is the same as the correctly rejected network packets, and it denotes the number of standard records identified as anomalies. The True Negative (TN) is the same as the correctly admitted network packets, and it denotes the number of standard records identified as usual. The False Negative (FN) is the same as the incorrectly admitted packets, and it denotes the number of anomaly records identified as usual.

**Accuracy:** The number of records classified correctly versus the total number of records (1)

$$AC = TP+TN \; / \; TP+TN+FP+FN \qquad (1)$$

**Actual Positive Rate (TPR):** This is also called Detection Rate (DR). The DR is calculated based on the percentage of the number of records identified correctly over the total number of anomaly records, as shown in (2)

$$TPR = TP \; / \; TP+FN \qquad (2)$$

**False Positive Rate (FPR):** The percentage of the number of records rejected incorrectly is divided by the total number of standard records (3)

$$TPR = FP \; / \; +TN \qquad (3)$$

**False Negative Rate (FNR):** The percentage of the number of records admitted incorrectly is divided by the total number of standard records (4)

$$TPR = FN \; / \; +TN \qquad (4)$$

## H. Next Steps - Tasks and Deliverables

Table III captures the list of activities, deliverables, and the effort represented in the number of weeks.

| Tasks (tasks to be conducted to achieve the dissertation objectives) | Deliverables (e.g. data analyzed, design, algorithm, codes, test plan…) | Effort (person-weeks) (e.g. 1p-2w 1 person takes 2 weeks) |
|---|---|---|
| Task 1 | Create POC (Proof of Concept) for dual model intrusion detection system based on CNN and DRL deep learning Algorithms | Proof of Concept | 1p- 2w |
| Task 2 | Development / Coding and Testing CNN-DRL dual model Intrusion detection system | Source Code / Executables | 1-8w |
| Task 3 | Study the Performance of the CNN-DRL dual model IDS, and evaluate different learning rate impacts on the accuracy and false positive rate using benchmark datasets NLS-KDD | Data Analyze Outcome | 1-2w |
| Task 4 | Compare the performance of dual model CNN-DRL IDS with other machine learning methods in binary and multi-class classification. | Data Analyze Outcome | 1-2w |
| Resources required by the project: | UoR Lab | |
| Costing for this project (if any): | None (Assumption no software license costs incurred as the required software is downloaded or available from University of Reading Students login) | |

TABLE III.       PLAN - TASKS AND DELIVERABLES

## VI. ETHICAL AND RISK CONSIDERATIONS

The research proposal is written considering not to duplicate or reflect existing research topic. To the best of the knowledge, the dual model based on CNN and DRL for Cloud is a unique topic. A thorough research analysis is made, and literature review is conducted.

## VII. REFERENCES

[1]    Jauro F, Chiroma H, Gital AY, Almutairi M, Shafi'i M A, Abawajy J H. Deep learning architectures in emerging Cloud computing architectures: Recent development, challenges and next research trend. Applied Soft Computing. 2020 Nov 1; 96:106582.

[2]    Yang H, Qin G, Ye L. Combined wireless network intrusion detection model based on deep learning method. IEEE Access. 2019 Jun 19;7:82624-32.

[3]    Xu C, Shen J, Du X, Zhang F. An intrusion detection system using a deep neural network with gated recurrent units. IEEE SAccess. 2018 Aug 28;6:48697-707

[4]    Hizal S, Cavusoglu U, Akgun D. A new Deep Learning Based Intrusion Detection System for Cloud Security. In 2021 3rd International Congress on Human-Computer Interaction, Optimization and Robotic Applications (HORA) 2021 Jun 11 (pp. 1-4). IEEE

[5]    Chen L, Kuang X, Xu A, Suo S, Yang Y. A Novel Network Intrusion Detection System Based on CNN. In 2020 Eighth International Conference on Advanced Cloud and Big Data (CBD), 2020 Dec 5 (pp. 243-247). IEEE.

[6]    Tang TA, Mhamdi L, McLernon D, Zaidi SA, Ghogho M. Deep recurrent neural network for intrusion detection in SDN-based networks. In 2018 4th IEEE Conference on Network Softwarization and Workshops (NetSoft), 2018 Jun 25 (pp. 202-206), IEEE.

[7]    Revathi S, Malathi A. A detailed analysis on NSL-KDD dataset using various machine learning techniques for intrusion detection. International Journal of Engineering Research & Technology (IJERT). 2013 Dec;2(12):1848-53.