

# **SEMI-SUPERVISED MACHINE LEARNING APPROACH FOR DDOS DETECTION**

*Y.NARASIMHA REDDY P.VEERESH.*

*Associate Professor, Associate Professor & Head of CSE Dept,  
Department of CSE*

*St.Johns College of Engineering and Technology, Yemmiganur, Kurnool (Dist).*

**Abstract** - Distributed denial of service (DDoS) attacks are a major threat to any network-based service provider. The ability of an attacker to harness the power of a lot of compromised devices to launch an attack makes it even more complex to handle. This complexity can increase even more when several attackers coordinate to launch an attack on one victim. Moreover, attackers these days do not need to be highly skilled to perpetrate an attack. Tools for orchestrating an attack can easily be found online and require little to no knowledge about attack scripts to initiate an attack. The purpose of this paper is to detect and mitigate known and unknown DDoS attacks in real time environments. Identify high volume of genuine traffic as genuine without being dropped. Prevent DDoS attacking (forged) packets from reaching the target while allowing genuine packets to get through. A DDoS attack slows or halts communications between devices as well as the victim machine itself. It introduces loss of Internet services like email, online applications or programme performance. We apply an automatic characteristic selection algorithm primarily based on N-gram sequence to obtain meaningful capabilities from the semantics of site visitors flows. DDoS attacks are the perfect planned attacks with the aim to stop the legitimate users from accessing the system or the service by consuming the bandwidth or by making the system or service unavailable. The attackers do not attack to steal or access any information but they decline the performance of the network and the system.

**Keywords** - Distributed Denial of Service (DDoS), Malware Detection, Machine learning, NLP Method, Text semantics.

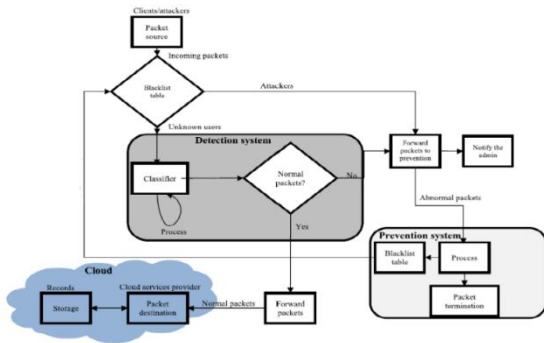
## **INTRODUCTION**

Data mining techniques have been used to develop sophisticated intrusion detection systems for the last two decades. Artificial Intelligence, Machine Learning (ML), Pattern Recognition, Statistics, Information Theory are the most used data mining techniques for intrusion detection. With the increase in dependability of the internet comes with it an important challenge: data availability. Data availability is a key requirement for a network system to be considered secure. Distributed denial of service attacks are intentional attempts by malicious users to disrupt or degrade the quality of a network or service. These attacks involve a number of compromised connected online devices, The use of botnets makes it easier for attackers to launch massive attacks due to the fact that they harness the power of a lot of devices for an attack. Attacks involving botnets also make it difficult to determine the exact source of the attack. Differentiating between flash crowds also poses a major challenge.

There are two main methods to launch DDoS attacks in the Internet. The first method is for the attacker to send some malformed packets to the victim to confuse a protocol or an application running on it (i.e., vulnerability attack). The other method, which is the most common one, involves an attacker trying to do one or both of the following:

(i) Disrupt a legitimate user's connectivity by exhausting bandwidth, router processing capacity or network resources; or (ii) Disrupt a legitimate user's services by exhausting the server resources (e.g., sockets, CPU, memory, disk/database bandwidth, and I/O bandwidth). Employing the resources of recruited computers to perform DDoS attacks allows attackers to launch a much larger and more disruptive attack. Furthermore, it becomes more complicated for the defense mechanisms to recognize the original attacker because of the use of counterfeit (i.e., spoofed) IP addresses by zombies under the control of the attacker. Most of the DDoS attacks launched to date have tried to make the victims' services unavailable, leading to revenue losses and increased costs of mitigating the attacks and restoring the services. In today's DDoS attacks, attackers use more complicated methods to launch an attack. Despite all of the efforts towards decreasing the number of DDoS attack incidents, they have expanded rapidly in the frequency and the size of the targeted networks and computers.

**SYSTEM ARCHITECTURE:**



The scheme assumes that source IP addresses are not spoofed. During detection, the detection Unit captures incoming packets within a time window and checks them against a blacklist to determine if their source addresses have been blacklisted. If the source is blacklisted, the packet is sent to the prevention unit without any further processing. If the source is not blacklisted, the packet is directed through a classifier to determine whether the packet is legitimate or malicious. A packet is considered malicious if its source requests to connect to the same destination more frequently than a set threshold. Legitimate packets are forwarded to their destination and malicious ones are forwarded to the prevention unit.

The prevention unit has three main functions; it sends a signal to the system administrator of the attack, then adds the source address to the blacklist if it does not already exist in the blacklist, and finally drops the packet. Evaluation of the system was performed under single and multiple source attack scenarios and was found to be consistent, with 97 percent and 94 percent accuracy, respectively.

**LITERATURE SURVEY:**

Various methodologies and techniques for reducing the effects of DDoS attacks in different network environments have been proposed and evaluated.

The authors identified users' requests or demands to a specific resource and their communicative data. Then samples of such requests are sent to the detection systems to be judged for abnormalities. Also, Liu and Gu have used Learning Vector Quantisation (LVQ) neural networks to detect attacks. This is a supervised version of quantisation, which can be used for pattern recognition, multi- class classification and data compression tasks.

Akilandeswari and Shalinie have introduced a Probabilistic Neural Network Based Attack Traffic Classification to detect different DDoS attacks. However, the authors focus on separating Flash Crowd Event from Denial of Service Attacks Gupta, Joshi and Misra have used a neural network to detect the number of zombies that have been involved in DDoS attacks. The process work- load is based on prediction using a feed-forward neural network.

They listed the most popular DDoS attacks on cloud systems, classified and discussed intrusion detection systems along with their challenges. Classification of the intrusion detection systems was knowledge-based and anomaly-based. Sub categorizations were done based on the scalability of the management system, user authentication and the response mechanism. Rashmi V. Deshmukh and Kailas K. Devadkar discussed DDoS attacks and provided a taxonomy of attacks in the cloud environment They classified defense mechanisms based on prevention, detection, and response to detection techniques Somani et al. presented insights into the characterization, prevention, detection, and mitigation mechanisms of DDoS attacks in the cloud environment A taxonomy of DDoS solutions was also presented and the solutions were categorized under prevention, detection, and mitigation. They concluded by discussing considerations to be made in selecting a defense solution.

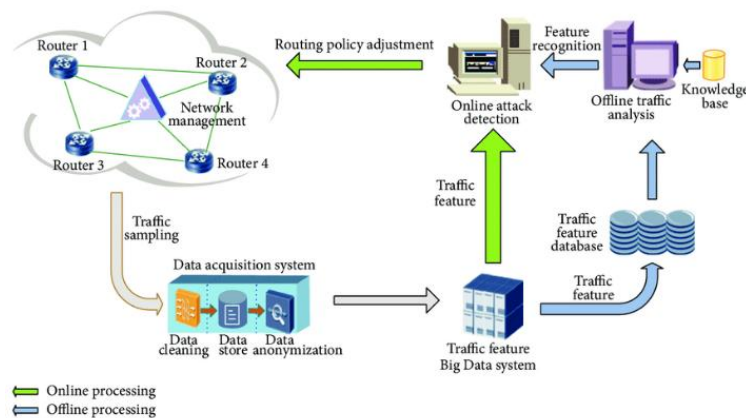
Mahjabin et al. presented a review on different DDoS attacks. They discussed attack phases in a DDoS attack, variations and evolutions of attacks as well as attackers' targets and motivations. They classified and analyzed prevention and mitigation techniques based on their underlying principle of operation. Underlying principles for the prevention methods reviewed were filters, secure overlay service, load balancing, honeypots, and awareness-based prevention systems. Mitigation techniques were also categorized broadly based on detection, response, and tolerance-based systems. They concluded by listing the key features, advantages, and limitations of the prevention and detection mechanisms reviewed Kalkan et al. presented DDoS attack scenarios in software defined networks (SDNs). Solutions to attacks were also broadly classified as intrinsic (having inherent

properties) and extrinsic (depending on external factors). Solutions were also classified according to their defense function (detection, mitigation, and both detection and mitigation) and SDN switch intelligence (capable switch vs. dumb switch).

Zare et al. came up with a paper reviewing papers on DDoS attacks and countermeasures between the years 2000 and 2016. They discussed intrusion detection systems and analyzed countermeasures against DDoS attacks based on the location of the defense mechanism; source-end, core-end, victim-end, and distributed defense.

There has been a great number of surveys on DDoS defenses in previous years. DDoS attacks are on the rise and there is a constant need to present the state-of-the-art defense mechanisms to aid researches in their attempt to combat such attacks. In previous surveys however, most defense categorizations were done based on their underlying principle of operation and not the function they perform in defending against attacks. Defense mechanisms were also not discussed into detail to give a greater understanding of their operation. Also, the individual defense mechanisms were not compared with each other, but rather, comparisons were mainly done based on the underlying principle of operation or location of deployment of the defense mechanism.

**Overview of DDOS Attack:**



Overview of real-time DDoS attack detection framework.

**TECHNOLOGIES USED**

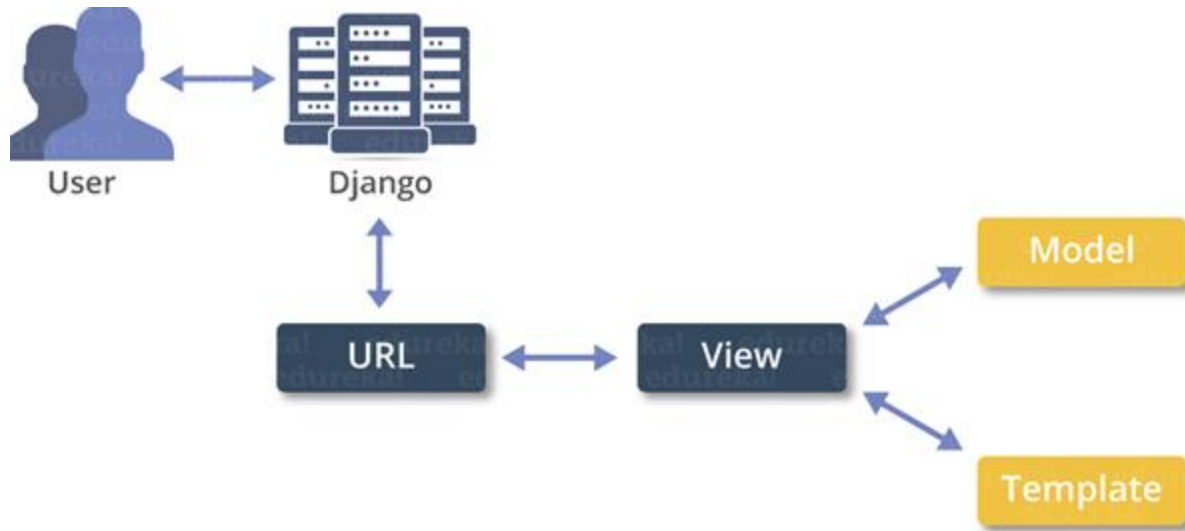
**Python**

Python is high-level, interpreted, and general purpose programming language. Released first in 1991 by Guido van Rossum. Its object-oriented approach aims the programmers to write clear, logical code for large and small scale development. It is a garbage collected and dynamically typed. This language contains a substantial body of documentation, abundant of it contributed by various authors. The markup used for the Python documentation is restructured text, developed by the docutils project, amended by custom directives and using a toolset named sphinx to post-process the Hypertext Mark-up Language (HTML) output.

**DJANGO**

Django is a high-level Python Web framework that encourages rapid development and clean, pragmatic design. Built by experienced developers, it takes care of much of the hassle of Web development, so you can focus on writing your app without needing to reinvent the wheel. It's free and open source.

Django's primary goal is to ease the creation of complex, database-driven websites. Django emphasizes reusability and "pluggability" of components, rapid development, and the principle of don't repeat yourself. Python is used throughout, even for settings files and data models.



Django also provides an optional administrative create, read, update and delete interface that is generated dynamically through introspection and configured via admin models.

**EXISTING SYSTEM:**

Malware detection methods can be static or dynamic. In dynamic malware detection approaches, the program is executed in a controlled environment (e.g., a virtual machine or a sandbox) to collect its behavioral attributes such as required resources, execution path, and requested privilege, in order to classify a program as malware or benign. Static approaches (e.g., signature-based detection, byte-sequence n-gram analysis, opcode sequence identification and control flow graph traversal) statically inspect a program code to detect suspicious applications. David et al proposed DeepSign to automatically detect malware using a signature generation method. The latter creates a dataset based on behaviour logs of API calls, registry entries, web searches, port accesses, etc, in a sandbox and then converts logs to a binary vector. They used deep belief network for classification and reportedly achieved 98.6% accuracy. In another study, Pascanu et al. Proposed a method to model malware execution using natural language modeling. They extracted relevant features using recurrent neural network to predict the next API calls. Then, both logistic regression and multi-layer perceptrons were applied as the classification module on next API call prediction and using history of past events as features. It was reported that 98.3% true positive rate and 0.1% false positive rate were achieved. Demme et al. examined the feasibility of building a malware detector in IoT nodes' hardware using performance counters as a learning feature and K-Nearest Neighbor, Decision Tree and Random Forest as classifiers. The reported accuracy rate for different malware family ranges from 25% to 100%. Alam et al. applied Random Forest on a dataset of Internet-connected smartphone devices to recognize malicious codes. They executed APKs in an Android emulator and recorded different features such as memory information, permission and network for classification, and evaluated their approach using different tree sizes. Their findings showed that the optimal classifier contains 40 trees, and 0.0171 of mean square root was achieved.

**PROPOSED SYSTEM:**

Junk code injection attack is a malware anti-forensic technique against OpCode inspection. As the name suggests, junk code insertion may include addition of benign OpCode sequences, which do not run in a malware or inclusion of instructions (e.g. NOP) that do not actually make any difference in malware activities. Junk code insertion technique is generally designed to obfuscate malicious OpCode sequences and reduce the 'proportion' of malicious OpCodes in a malware. In our proposed approach, we use an affinity based criteria to mitigate junk OpCode injection anti-forensics technique. Specifically, our feature selection method eliminates less instructive OpCodes to mitigate the effects of injecting junk OpCodes. To demonstrate the effectiveness of our proposed approach against code insertion attack, in an iterative manner, a specified proportion (5%, 10%, 15%, 20%, 25%, 30%) of all elements in each sample's generated graph were selected randomly and their value incremented by one. For example, in the 4th iteration of the evaluations, 20% of the indices in each sample's graph were chosen to increment their value by one.

In addition, in our evaluations the possibility of a repetitive element selection was included to simulate injecting an OpCode more than once. Incrementing  $E_{i,j}$  in the sample's generated graph is equivalent to injecting  $OpCode_j$  next to the  $OpCode_i$  in a sample's instruction sequence to mislead the detection algorithm. Algorithm 2 describes an iteration of junk code insertion during experiments, and this procedure should repeat for each iteration of k-fold validation. To show the robustness of our proposed approach and benchmark it against existing proposals, two congruent algorithms described in Section 1 are applied on our generated dataset using Adaboost as the classification algorithm. All evaluations were conducted using MATLAB R2015a running on a Microsoft Windows 10 Pro personal computer powered by Intel Core i7 2.67GHz and 8GB RAM. A 10-fold cross validation was used in the validating, and the comparative summary is presented in Table 1. It is clear that our proposed approach outperforms the proposals of Hashemi et al. and Santos et al.

The approach of Santos et al. is a basic and commonly-known OpCode based malware detection algorithm and the approach of Hashemi et al. is the most similar in terms of using eigenspace as the basis. Accuracy is a general criteria for evaluating performance of an algorithm for both malware and benign class identification. The proposed approach achieves a high accuracy of 99.68%, while the approaches of Hashemi et al. and Santos et al. respectively achieve 98.59% and 95.91% accuracy. Recall or detection rate is an important criteria and the proposed approach achieves 98.37%, in comparison to 81.55% and 77.70% for the other two approaches. Our proposed approach also outperforms the approaches of Hashemi et al. and Santos et al., in terms of precision rate and F-Measure. Utilizing class-wise feature selection appears to result in beneficial features of minor class to be more effective during classification phase. Also, using Formulation to calculate OpCode's distance leads to the ability to represent more OpCode sequence patterns in the sample's graph. It also appears that employing deep neural networks for classification leads to a better classifier.



**MODULES:**

There are three modules can be divided here for this project they are listed as below

- User Apps
- DDOS Attack Deduction
- Classifications of DDOS attack
- Graphical analysis

From the above four modules, project is implemented. Bag of discriminative words are achieved

**1. User Apps**

User handling for some various times of smart phones , desktops laptops and tablets .If any kind of devices attacks for some unauthorized Malware softwares. In this Malware on threats for user personal dates includes for personal contact, bank account numbers and any kind of personal documents are hacking is possible.

**2. DDOS Attack Deduction**

User search the any link Notably, not all network traffic data generated by malicious apps correspond to malicious traffic. Many malware take the form of repackaged benign apps; thus, Malware can also contain the basic functions of a benign app. Subsequently, the network traffic they generate can be characterized by mixed benign and malicious









**CONCLUSION:**

In this paper we propose and demonstrate an effective DDoS attack detection and mitigation system based on the observation of the bi-directional nature of the Internet traffic. We discuss an agent based technique which enables each edge router to validate the bi-directionality of the incoming traffic passing through them. Also, we propose and demonstrate a packet marking scheme called XORID, which can be used to defend against spoofing based DDoS attack. As our future work, we are investigating techniques which can defend application layer DDoS attacks at the edge of a network based on informations received from the internal servers of the network.

Although this method is capable of detecting DDoS attack with a high accuracy and low detection time, however the method will be ineffective if the attacker generates the attack traffic with well-proportioned.

**REFERENCES:**

[1] C. Rossow, “Amplification hell: revisiting network protocols for DDoS abuse,” in Symposium on Network and Distributed System Security (NDSS), Feb. 2014.

[2] F.-Y. Lee and S. Shieh, “Defending against spoofed DDoS attacks with path fingerprint,” *Comput. Sec.*, Vol. 24, no. 7, pp. 571–586, Oct. 2005.

[3] . Bhuyan MH, Bhattacharyya DK, Kalita JK (2015) An empirical evaluation of information metrics for low-rate and high-rate ddos attack detection.

[4] Pattern Recogn Lett 51:1–7 2. Lin S-C, Tseng S-S (2004) Constructing detection knowledge for ddos intrusion tolerance. *Exp Syst Appl* 27(3):379–390 3.

[5] Chang RKC (2002) Defending against flooding-based distributed denial-of-service attacks: a tutorial. *IEEE Commun Mag* 40(10):42–51

[6] Yu S (2014) Distributed denial of service attack and defense. Springer, Berlin