

Indian Journal of Modern Research and Reviews

This Journal is a member of the 'Committee on Publication Ethics'

Online ISSN:2584-184X



Research Article

A Review of Machine Learning Algorithms for Malware Detection

Shinda Singh ^{1*}, Shalu Gupta ², Jaswinder Brar ³

¹ Student, Department of Computer Applications, Guru Kashi University, Talwandi Sabo, Punjab, India

² Associate Professor, Department of Computer Applications, Guru Kashi University, Talwandi Sabo, Punjab, India

³ Assistant Professor, Department of Computer Applications, Guru Kashi University, Talwandi Sabo, Punjab, India

Corresponding Author: *Shinda Singh

DOI: <https://doi.org/10.5281/zenodo.18309966>

Abstract

This study centres on dynamic malware detection, recognizing that malicious software evolves continuously and demands more adaptive security approaches. With new malware emerging almost every day and exploiting weaknesses across the Internet, traditional manual and heuristic-based analysis has become insufficient. To address this growing challenge, this research employs automated, behaviour-based detection supported by machine learning techniques. In this approach, malware samples are executed within a controlled environment, their behaviours are monitored, and detailed reports are generated. These reports are then transformed into sparse vector representations, which serve as input for various machine learning models. The classifiers applied in this study include kNN, DT, RF, AdaBoost, SGD, Extra Trees, and Gaussian NB. An evaluation of the experimental results shows that RF, SGD, Extra Trees, and Gaussian NB all reached 100% accuracy on the test set, along with perfect precision (1.00), recall (1.00), and f1-scores (1.00). These findings suggest that a proof-of-concept system combining autonomous behaviour analysis with machine learning can detect malware both effectively and efficiently.

Manuscript Information

- ISSN No: 2583-7397
- Received: 20-11-2025
- Accepted: 25-12-2025
- Published: 20-01-2026
- IJCRM:4(1); 2026: 128-133
- ©2025, All Rights Reserved
- Plagiarism Checked: Yes
- Peer Review Process: Yes

How to Cite this Article

Singh S, Gupta S, Brar J. A Review of Machine Learning Algorithms for Malware Detection. Indian J Mod Res Rev. 2026;4(1):128-133.

Access this Article Online



www.multiarticlesjournal.com

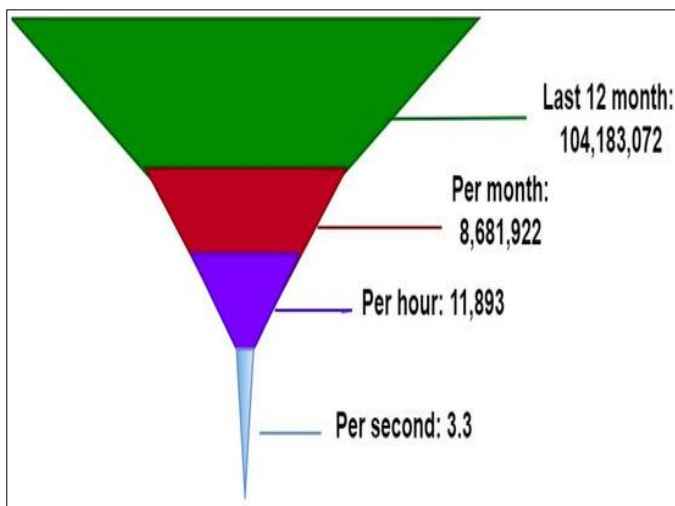
KEYWORDS: malware; cyberattacks; IoT; malicious threats; machine learning classifiers; RF; DT; cyber security; suspicious activity; SGD; extra trees; Gaussian NB

INTRODUCTION

Cyberattacks have become a major concern in the digital world, and traditional signature-based antivirus tools often fail to detect new or polymorphic malware. With malicious software spreading rapidly across the Internet, manual static analysis is no longer efficient or practical. As a result, researchers are turning toward automated techniques that combine dynamic malware analysis with machine learning methods to improve detection accuracy. The breakthrough came with deep convolutional neural networks, which reduced error rates dramatically by learning features directly from raw pixels. Object detection and recognition form an essential component of image processing and have emerged as a significant research area within the domains of image processing and pattern recognition [29, 30]. Edge detection techniques are widely used in various research domains, including computer vision, machine learning, and pattern recognition [31, 32].

Although many antivirus systems exist, malware incidents continue to rise, highlighting the need for more reliable solutions. Dynamic analysis offers clear advantages over static approaches, as it is harder for malware to hide its behaviour during execution. Machine learning has recently gained attention for predicting malicious patterns and identifying malware families, yet there is still no unified comparison of different algorithms. To address this, we conducted experiments evaluating multiple machine learning models for malware detection and classification. Figure 1 shows the number of new malware threats detected per second.

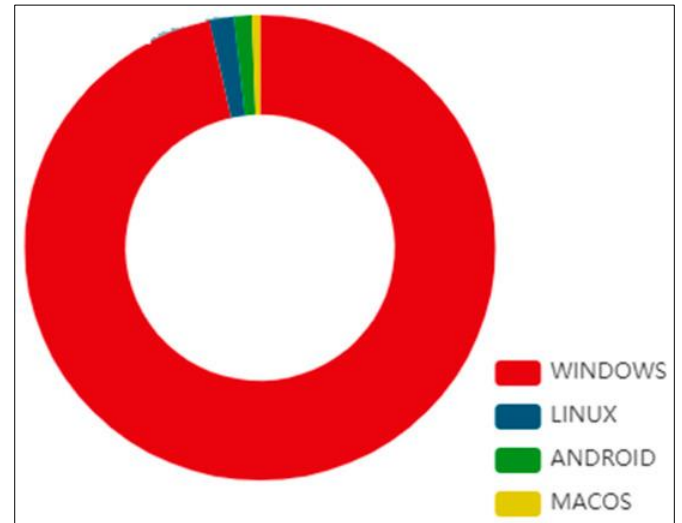
Figure 1. Threats of new malware per second



A dataset of real malware samples and benign programs collected from Virus Total was executed in a sandbox environment to capture behavioural information, which was later used to evaluate multiple machine learning models based on standard performance metrics. The execution data, stored as JSON reports, offered a rich set of features representing each sample's behaviour, allowing us to separate malicious files from harmless ones. This study is motivated by the fact that

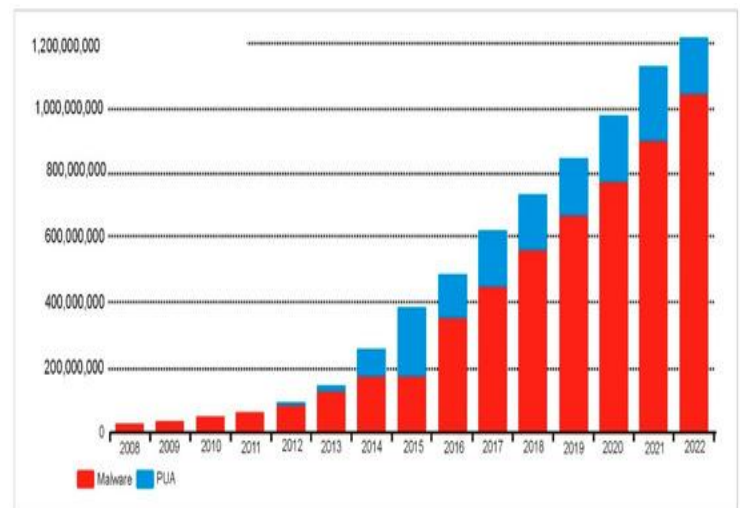
different detection approaches often behave inconsistently, even under similar conditions, due to their varying optimisation goals. To address this, we provide recommendations for researchers and discuss future directions for improving dynamic malware detection using machine learning. Figure 2 presents the classification of OS-based threats.

Figure 2. Classification of OS-based malware threats.



As Internet access expands across a wide range of devices—from desktop computers to embedded systems—more people rely on it for information and quick communication. With constant connectivity, users can access online services at any time. However, this growth has also created opportunities for cybercriminals, leading to a rapid rise in malware. As Internet use has increased, so has the appeal of distributing malicious software. Figure 3 shows how malware detections have grown exponentially in recent years.

Figure 3. Total amounts of malware and potentially unwanted applications (PUAs).

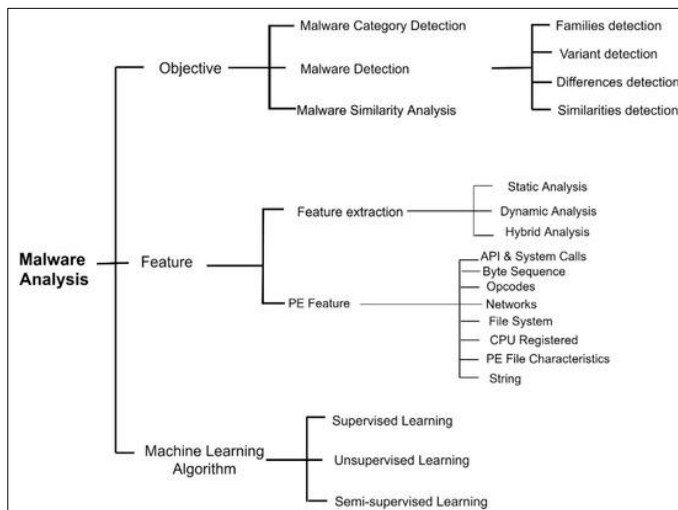


Anti-malware tools, intrusion detection systems, and other security measures have emerged in response to the growing damage caused by malware. Yet several critical challenges remain, especially as attackers constantly adopt new techniques and exploit vulnerabilities in widely used software. Many methods from various fields have been proposed to improve malware detection, but dynamic analysis has proven more effective than static approaches because it is harder for malware to hide its malicious actions during execution. As the advantages of automated and dynamic techniques have become clearer, researchers have increasingly moved away from traditional static detection methods.

2. LITERATURE REVIEW

Trinius (2016) introduced a new behaviour-tracking representation called MIST, designed to capture and analyse malicious program actions more effectively using data mining and machine learning. This representation can be collected automatically through behaviour-monitoring tools or generated manually from existing analysis reports. Rieck (2018) explored how similarities among malware samples can be used to group and categorise them. Patil (2020) further observed that different versions of malware often display consistent behavioural patterns that reveal the intentions of their authors. Their approach begins with monitoring malware activity in a sandbox environment, followed by using an antivirus-labelled dataset, and finally analysing the collected results as shown in Figure 4.

Figure 4. Malware analysis methods.



Learning-based methods are widely used to train malware behaviour classifiers, emphasising the most relevant behavioural features that explain classification decisions. Rieck (2018) introduced a machine-learning-driven framework that automatically analyses malware activity and assigns unknown samples to known behavioural classes. Christodorescu (2018) proposed another approach that compares the execution patterns of malicious programs against benign applications. By

identifying harmful features present only in malware, their method enables detectors to flag new threats.

Machine learning systems rely heavily on high-quality feature engineering, selection, and representation. Models are trained on labelled data to form a decision boundary separating malware from legitimate software. Domain knowledge remains essential for designing effective features. However, traditional ML-based malware detectors face challenges: adversaries can reverse-engineer models to evade detection, and high-quality public datasets are limited due to privacy and security concerns. As a result, many researchers build their own datasets using standard data-science procedures. Ye (2017) highlighted the scale and complexity of analysing such datasets, making real-time ML-based malware detection difficult.

Modern AI systems increasingly use deep learning, which captures complex feature representations through layered learning. Neelam (2020) reviewed several studies applying deep learning models to malware analysis. In 2015, Microsoft organised a Kaggle malware classification challenge with nearly 20,000 samples, prompting Ronen (2015) to examine published and emerging research in the field.

Souri (2020) conducted an extensive review of malware detection methods based solely on data analysis, dividing the literature into signature-based and behaviour-based approaches. Their findings suggest that hybrid methods—combining static and dynamic analysis—achieve higher accuracy than using either technique alone. Yanfeng (2018) similarly summarised cloud-based malware detection work, outlining feature extraction techniques, classification methods, and malware evolution trends; however, these studies only covered research up to 2017, indicating the need for further updates.

Ucci (2017) compiled a structured review of machine learning algorithms used to identify malicious PE files on Windows systems, organising prior studies by objectives, methodology, and dataset characteristics. The authors also discussed the broader “malware analysis economy,” noting persistent challenges and the need for continued research, especially as three years have passed since their initial publication.

A. Research Gap

Cybercriminals continue to develop and spread malicious software to infiltrate systems or cause damage. Organisations typically rely on antivirus tools, log analysis, and activity monitoring to detect suspicious patterns that signal known threats. Although signature-based detection works well for identifying previously documented malware, attackers can easily evade these systems by modifying or obfuscating their code. As a result, researchers have focused on improving detection accuracy, reducing false positives, and lowering processing time.

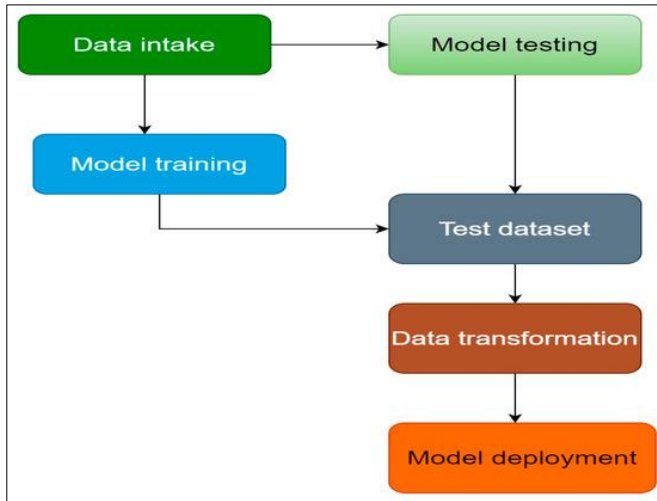
However, advancing malware detection remains difficult due to several challenges within the malware ecosystem. In this study, we review existing methods used to detect malware in previously released files and highlight areas that require further investigation. We also examine ongoing efforts to standardise how malware is measured, described, evaluated, and

architected. Identifying these factors can help make malware detection research more consistent, expandable, and accessible to future researchers.

4. Research Framework

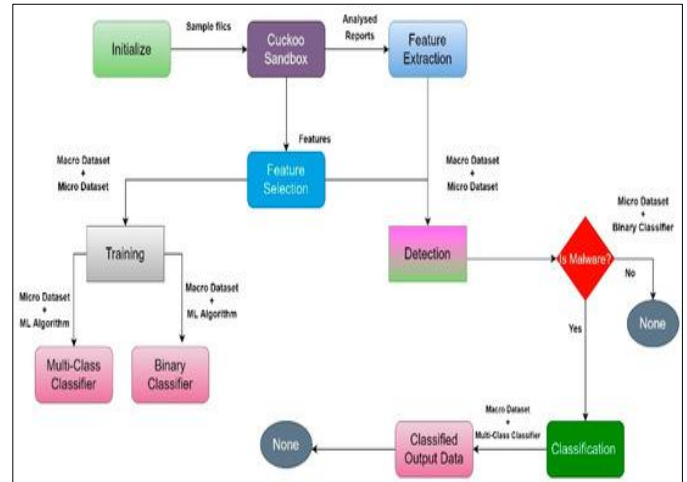
The rapid rise of advanced and sophisticated malware has become a significant threat to modern computing systems. Traditional signature-based detection methods are increasingly ineffective as the number of new malware samples grows at an exponential rate. Research shows that machine learning techniques can reliably detect and classify malicious files. Their performance can be further enhanced through feature-selection methods, which identify the most relevant attributes and reduce dataset size, leading to faster processing and improved accuracy. The overall research framework used in this study is illustrated in Figure 5.

Figure 5. Research framework.



In this study, we present a machine learning-driven approach to improve the efficiency and accuracy of malware detection and classification. For dynamic analysis, we used the Cuckoo sandbox, which runs malware in an isolated environment and produces detailed reports of its behaviour. We also developed a feature extraction and selection module that collects relevant attributes from these reports and identifies the most significant ones, ensuring high accuracy with minimal computational effort. A variety of machine learning algorithms were then applied for precise classification and detection. Our experimental evaluation showed that the proposed approach achieves higher accuracy than existing methods. The overall structure of the malware detection framework is illustrated in Figure 6.

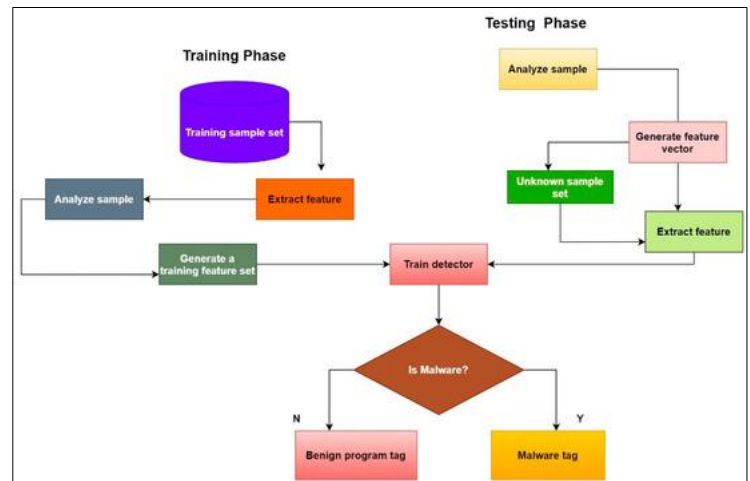
Figure 6. Malware detection framework structure.



5. RESEARCH METHODOLOGY

Figure 6 provides a high-level overview of our machine learning-based malware detection process. This workflow involves selecting suitable datasets for training the classifier, identifying advanced malware samples, and choosing the most relevant features for model development. Below, we provide a more detailed description of each step followed in the study. The complete proposed methodology is illustrated in Figure 7.

Figure 7. Proposed method of malware detection.



6. RESULTS AND DISCUSSION

To ensure the effectiveness of any classification method, both training and testing phases are essential. The system must be trained using a mix of malicious and benign samples so that the classifier can learn to distinguish between them. With machine learning, the model gradually improves as it processes more labelled data, enabling it to generate accurate predictions. In our study, classifiers such as Random Forest, SGD, Extra Trees, and Gaussian NB showed enhanced performance as the dataset

size increased. During validation, each model was tested on a separate set of unseen files—some malicious and some benign—and was required to classify them correctly.

Figure 9 provides a visual comparison of the RF, SGD, Extra Trees, and Gaussian NB models. A dropout layer is used in the final fully connected stage of these models. In practice, this dropout appears to function more as an additional architectural layer rather than a strict regularisation mechanism.

7. CONCLUSIONS

This study addresses the limitations of manual feature engineering and traditional learning methods by integrating RF, ASG, Extra Trees, and Gaussian NB into a novel ensemble deep neural network for malware detection. The ASG, Extra Trees and Gaussian NB models achieved perfect performance with 100% accuracy, precision, recall, and F1-score. The combined ensemble effectively captures sequential patterns, long-term dependencies, and spatial correlations, resulting in accuracy close to 1 in both training and testing phases. Overall, machine learning-based approaches significantly enhance malware detection by improving accuracy, reducing false positives, and enabling faster identification. Researchers typically evaluate these algorithms by dividing data into training and testing sets to assess real-world performance.

REFERENCE

1. Akhtar MS, Feng T. Malware analysis and detection using machine learning algorithms. *Symmetry*. 2022;14:2304.
2. Akhtar MS, Feng T. Detection of malware by deep learning as CNN-LSTM machine learning techniques in real time. *Symmetry*. 2022;14:2308.
3. Akhtar MS, Feng T. Deep learning-based framework for the detection of cyberattack using feature engineering. *Secur Commun Netw*. 2021;2021:6129210.
4. Baghirov E. Techniques of malware detection: Research review. In: *Proc IEEE 15th Int Conf Appl Inf Commun Technol (AICT)*; 2021; Baku, Azerbaijan. p. 1–6.
5. Akhtar M, Feng T. Comparison of classification model for the detection of cyber-attack using ensemble learning models. *EAI Endorsed Trans Scal Inf Syst*. 2022;9:e6.
6. Saad S, Briguglio W, Elmiligi H. The curious case of machine learning in malware detection. *arXiv*. 2019:1905.07573.
7. Muppalaneni N, Patgiri R. Malware detection using machine learning approach. In: *Int Conf Big Data Mach Learn Appl*; 2021; Vancouver, Canada. Singapore: Springer.
8. Baset M. *Machine learning for malware detection* [master's thesis]. Edinburgh: Heriot-Watt University; 2016.
9. Singhal P, Raul N. Malware detection module using machine learning algorithms to assist in centralized security in enterprise networks. *Int J Netw Secur Its Appl*. 2012;4:61–67.
10. Agarkar S, Ghosh S. Malware detection and classification using machine learning. In: *Proc IEEE Int Symp Sustain Energy Signal Process Cyber Secur (iSSSC)*; 2020; Gunupur, India. p. 1–6.
11. Cuan B, Damien A, Delaplace C, Valois M. Malware detection in PDF files using machine learning. In: *Proc 15th Int Conf Secur Cryptogr (SECRYPT)*; 2018; Porto, Portugal. p. 578–585.
12. Rimon SI, Haque MM. Malware detection and classification using hybrid machine learning algorithm. In: *Intell Comput Optim (ICO 2022)*. LNNS, vol. 569. Cham: Springer; 2023.
13. Hussain A, Asif M, Ahmad M, Mahmood T, Raza M. Malware detection using machine learning algorithms for Windows platform. In: *Int Conf Inf Technol Appl*; 2022; Lisbon, Portugal. Singapore: Springer.
14. Gavriluț D, Cimpoeșu M, Anton D, Ciortuz L. Malware detection using machine learning. In: *Proc Int Multiconf Comput Sci Inf Technol*; 2009; Mragowo, Poland. Vol 4. p. 735–741.
15. Ye Y, Li T, Adjero D, Iyengar S. A survey on malware detection using data mining techniques. *ACM Comput Surv*. 2017;50:1–40.
16. Neelam C, Singh A, Gaurav G. Android malware detection using improvised random forest algorithm. *Glob J Res Anal*. 2020;9(3).
17. Mazuera-Rozo A, Bautista-Mora J, Linares-Vásquez M, Rueda S, Bavota G. The Android OS stack and its vulnerabilities: An empirical study. *Empir Softw Eng*. 2019;24:2056–2101.
18. Azmoodeh A, Dehghantanha A, Choo KKR. Robust malware detection for Internet of battlefield things devices using deep eigenspace learning. *IEEE Trans Sustain Comput*. 2018;4:88–95.
19. Android malware dataset for machine learning. *Kaggle* [Internet]. Available from: <https://www.kaggle.com/shashwatwork/android-malwaredataset-for-machine-learning>
20. Jin X, Xing X. A malware detection approach using malware images and autoencoders. In: *Proc IEEE 17th Int Conf Mobile Ad Hoc Sensor Syst (MASS)*; 2020; Delhi, India.
21. Darem AA, Ghaleb FA. An adaptive behavioral-based incremental batch learning malware variants detection model using concept drift detection and sequential deep learning. *IEEE Access*. 2021;9.
22. Wu D, Guo P. Malware detection based on cascading XGBoost and cost sensitive. In: *Int Conf Comput Commun Netw Secur (CCNS)*; 2020; Xi'an, China.
23. McGiff J, Hatcher WG. Towards multimodal learning for Android malware detection. In: *Int Conf Comput Netw Commun (ICNC)*; 2019; Istanbul, Turkey. p. 432–436.
24. Anuar NA, Mas'ud MZ, Bahaman N, Ariff NAM. Analysis of machine learning classification in Android malware detection through opcode. In: *IEEE Conf Appl Inf Netw Secur (AINS)*; 2020; Kota Kinabalu, Malaysia.
25. Pedregosa F, et al. Scikit-learn: Machine learning in Python. *J Mach Learn Res*. 2011;12:2825–2830.

26. Introduction to Simple Imputer class. *Scikit-learn* [Internet]. Available from: <https://scikitlearn.org/stable/modules/generated/sklearn.impute.SimpleImputer.html>
27. Huang T, Zhao R, Bi L, Zhang D, Lu C. Neural embedding singular value decomposition for collaborative filtering. *IEEE Trans Neural Netw Learn Syst.* 2022;33:6021–6029.
28. Li Q, Zheng X, Wu X. Neural collaborative autoencoder. *arXiv.* 2017:1712.09043.
29. Gupta S, Singh YJ, Kumar M. Object detection using multiple shape-based features. In: *IEEE 4th Int Conf Parallel Distrib Grid Comput (PDGC)*; 2016. p. 433–437.
30. Gupta S, Singh YJ. Glowing window based feature extraction technique for object detection. In: *Int Conf Data Manag Anal Innov*; 2020; New Delhi, India.
31. Gupta S, Singh H, Singh YJ. Comprehensive study on edge detection. In: *Int Conf Commun Electron Digit Technol (NICE)*; 2023; India.

Creative Commons License
<p>This article is an open-access article distributed under the terms and conditions of the Creative Commons Attribution–NonCommercial–NoDerivatives 4.0 International (CC BY-NC-ND 4.0) License. This license permits users to copy and redistribute the material in any medium or format for non-commercial purposes only, provided that appropriate credit is given to the original author(s) and the source. No modifications, adaptations, or derivative works are permitted.</p>