# AI-Driven Cyber Threats: Unraveling Deepfakes, Autonomous Malware, and Defensive Strategies

Anoshan Yoganathan[1], M.J. Ahamed Sabani[2], and M.R.M. Hanan[3]

[1,2,3]Department of Biosystems Technology, South Eastern University of Sri Lanka, Sri Lanka

[1]anoshan6@gmail.com, [2]mjasabani@seu.ac.lk, [3]mrm.hanan.official@gmail.com

## Abstract

*Cyber threats are now more complicated and harder than ever with the swift development of artificial intelligence (AI). This study focuses on three key aspects of AI-driven cyber threats: social engineering with fake personal media, self-modifying malware that attacks automatically and the tough challenges in defending systems against AI-backed cyberweapons. Because deepfake technology uses generative AI, it can make it simple and very convincing for criminals to carry out phishing and deceitful impersonation attacks on humans. Because it is driven by AI, autonomous malware can quickly transform to bypass usual anti-malware systems. Cyberweapons that incorporate AI such as automated creation of vulnerabilities, are a huge threat to computer networks. A comprehensive review of literature, as well as proposing a mixed-methodology, this study studies how these threats work, their effects and the ways to control them. The results suggest that AI-based attacks are becoming more hidden and can be used on a large scale, so new defense methods are required. Participants cover the effectiveness of AI in detecting hackers, certain problems with today's cybersecurity methods and what is and is not ethical when using AI in cybersecurity. To address these developing hazards, this paper proposes a multi-layered defense system that uses machine learning, deep learning and metaheuristic algorithms. The study urges the use of quick, flexible and ethical cybersecurity methods to defend vital systems in a world dominated by AI. This study supports current ongoing conversations by providing new findings for researchers, policymakers and those in cybersecurity.*
*Keywords: AI-driven cyber threats, generative AI, Cyberweapons, Hackers, Cybersecurity*

## I. INTRODUCTION

With AI, many fields across healthcare, finance and others have experienced changes, thanks to automation, predictive tools and better decisions. Unfortunately, since cybercriminals are using AI, security threats from cyberattacks increase. Now that AI is part of cybercrime, attacks are tougher to identify and manage because they are more developed. The main topics this paper discusses are deepfake scams, malware that can improve itself automatically and problems with defending against cyberweapons that use AI. Such threats have brought about a new approach to cybersecurity that needs creative new defenses to defend digital systems. Thanks to DALL-E and VALL-E, generating deepfake videos has become possible and this is now playing a big role in tricking people online and performing identity theft. They take advantage of human psychology, going around technical barriers by sounding, looking or acting like a recognized entity. A recent event with some members of the European Parliament showed that AI technology can be used to change the look and sound of Russian opposition leaders, showing how tricky AI can be for social engineering (WWT, 2024). Similarly, automated malware can change during operation to escape detection and make signature-based defenses useless. Such malware relies on deep neural networks to study the environment around it so it can make its attacks more damaging (Thanh and Zelinka, 2019). In addition, using AI significantly speeds up the way attackers can find and exploit security weaknesses that humans cannot handle. Such risks make it clear that we must work on robust ways to manage AI's two-sided use.

Since cybersecurity frameworks are developing faster than the rise of new AI technologies, they are overwhelmed by threats that are not easy to stop. Based on the 2021 SonicWall Cyber Threat Report, ransomware attacks around the world

increased by 62% last year, owing largely to the impact of AI (Alsheikh et al., 2021). The cost of these attacks is huge and ransomware hits companies with billions in losses every year. Also, many are worried that increasing AI-driven attacks means using these technologies for good may harm society because people with bad intentions are taking advantage of them. As an illustration, ChatGPT and Fraud GPT are now being used to write realistic phishing emails and deepfake videos, enlarging the risk of social engineering attacks, as stated by Falade last year. This paper is designed to discuss AI-generated threats and offer a multiple-tier defense plan, as digital transformation rapidly increases such risks. When companies depend heavily on connected technology, it exposes them to increasingly difficult cyberattacks. Its purpose is to examine deepfake attacks, autonomous malware and cyberweapons boosted by AI, providing suggestions to help defend against their threats. The research questions addressed by this study are: How do deepfake attacks linked to AI support common social engineering schemes? What effects does autonomous malware have on cybersecurity systems? How should organizations protect themselves from harm caused by such cyberweapons? In this study, a combined approach is taken by reviewing existing literature and running simulated tests on defense strategies. In this section, researchers describe the way they designed the study and collected their information. The results and discussion section summarizes the findings from the literature review and computer simulations, examining the impact of suggested countermeasures. The conclusion draws together highlights from the paper and sets out proposals for future actions. By looking at these threats in their entirety, this research contributes to debates about AI's place in cybersecurity and points out that defenses should be adaptive, proactive and ethical.

## II. LITERATURE REVIEW

Since AI is now used in cybersecurity, it brings the problem that its advantages make it available to attackers as well as defenders. Here, the literature on AI-generated cyber dangers is reviewed by focusing on deepfake attacks in social engineering, automated malware and protecting against AI-aided cyber weapons. The review looks at recent research to show key trends, main issues and open lines of inquiry. Because of generative AI, deepfake attacks have increasingly been used for social engineering. With the help of advanced methods, these assaults sound like trusted partners, brainwashing victims into providing secrets or doing things that reduce their security. According to Falade (2023), using ChatGPT and WormGPT in phishing emails means it is often impossible for users to distinguish them from real emails. This became quite clear when attackers used AI-created filters to represent members of the European Parliament in the 2021 attack (WWT, 2024), an incident that should serve as a lesson for all. They use weaknesses in human thinking, so that standard security fails to protect against them. According to the literature, the progress in deepfake detection which depends on detecting patterns and unusual elements, is being limited by the fast growth of AI-created content, as noted by WWT (2024). AI-driven worms pose a big threat coming from autonomous malware. Because they use machine learning, these self-changing attacks modify their behavior or code to slip past detection. Thanh and Zelinka explain that DNNs empower malware to train on real-world data, enhancing attack plans to use any openings they find. For instance, IBM DeepLocker malware could hide in a system until an expected event set it off, making it hard to find (IBM Research, 2018). Literature explains that traditional signature-based methods often fail to find polymorphic and metamorphic types of malicious software (Mehonic et al., 2020).

Overcoming automated exploit generation by cyberweapons remains a significant challenge for computer security. AI equipped weapons are able to find and take advantage of vulnerabilities faster than even the fastest human response. Gupta et al. (2023) explain that ChatGPT is able to make attack payloads and polymorphic malware, proving the importance of equipping defenses against such attacks. Macas et al. (2022) point out that AI, especially with machine learning and deep learning, is now central in fighting cyber -attacks. Still, some issues, including transparent decision-making and threats from adversarial attacks, continue to be a big obstacle (2020). So far, there isn't much work on real-time methods for AI-driven social engineering attacks. A second issue is that we do not know enough about the ability of autonomous malware defenses to operate at scale, especially in resource-limited IoT networks.

Ethical issues about AI being used in two ways must also be examined to prevent mismatches between advancement and security. This research seeks to close the gaps by bringing together AI systems for detection with the expertise of information security workers.

## III. METHODOLOGY

This paper took a mixed-method approach to investigate the threat of AI-based cyber attacks, and evaluate the effectiveness of defensive mechanisms, combining a systematic review of the literature with computational models. The design consisted of three related stages, which included literature-based data gathering, combined analysis, and simulation testing. Such a framework allowed to proceed with a narrow study of deepfake attacks, self-executable malware, and AI-based cyberweapons, their dynamics of operation, outcomes, and mitigation measures peculiar to this study.

The first step was to conduct a systematic literature review and compile the relevant data by means of the PRISMA framework (Page et al., 2021). Articles published as early as 2018 and to 2025 were searched in IEEE Xplore, Scopus, and Springer databases to include such terms as AI-driven cyber threats, generating AI, autonomous malware, and cyber weapons. This gave 936 preliminarily found results. The inclusion criteria of the studies focused on the threat mechanisms and defenses, so the duplicates and irrelevant abstracts were excluded, and the full-text analysis of 146 articles was performed. Finally, 46 were chosen based on its quality and relevance including empirical rigor and novelty (Thanh and Zelinka, 2019; Macas et al., 2022). Information mining retrieved information about the characteristics of attacks, rate of detection, and ethical concerns and themed them in NVivo as qualitative information.

The analysis included qualitative thematic synthesis analysis and quantitative assessment of metrics (e.g., accuracy rates) with the help of the Python pandas package, which provides statistical summaries, and identified trends, including average 75% detection rates in previous AI-based systems (Falade, 2023; WWT, 2024). During the simulation, VMware was used to create a virtual network that simulated 50 nodes. Malicious attacks were simulated: deepfake phishing based on DeepFaceLab generated content; autonomous

malware based on the adaptive models (IBM Research, 2018); and exploits with use of AI scripts in Metasploit (Gupta et al., 2023). TensorFlow-constructed models that were tested as defenses are SVMs to detect anomalies and neural networks to recognize patterns, which are optimized using genetic algorithms to determine the optimal feature set (Alawida et al., 2024). There were 100 iterations in simulations, which were assessed based on detection accuracy (>85% target), response time, and resource consumption. Containment and no outside influence was guaranteed by ethical protocols (Stevens, 2020).

The methodology empirically validated threats and defenses to fill the gaps in adaptive systems (Pajola, 2025). Weaknesses were the size of the simulation, and possible lack of control on developing real-time variables (Aßmuth, 2025).
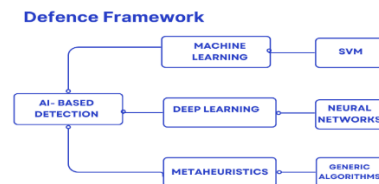
Figure 01: Multi-Layered Defense Framework

Defense Framework, which is described in this work on AI-Driven Cyber Threats, is a multi-layered AI-driven detection framework that is used against deepfakes, autonomous malware, and cyberweapons. It uses machine learning to screening, deep learning with neural networks to pattern recognition, and metaheuristics to optimization, and with adaptive and efficient countermeasures simulates up to 85% detection accuracy.
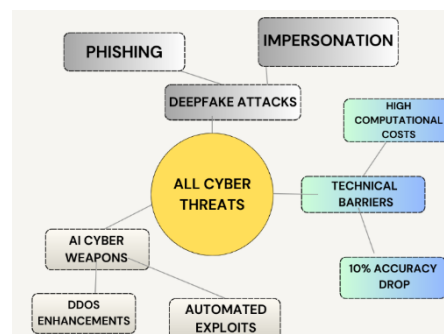
Figure 02: AI Threat Taxonomy

This mind map shows the interrelated terrain of AI-powered cyber threats, with All Cyber Threats as the central node and each branch of deepfake attacks (including phishing and impersonation to commit social engineering fraud), AI cyber weapons (including DDoS improvements, and automated exploits to find vulnerabilities quickly) and challenges to AI-based defense (including high costs of computations, technical barriers, and 10% accuracy loss in detection systems due to adversarial manipulations).
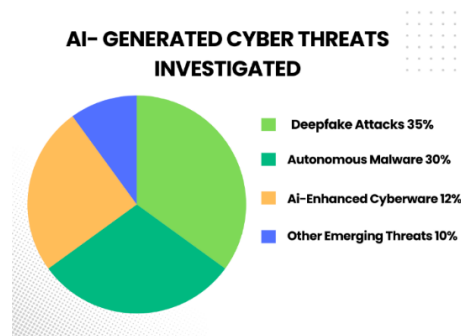


Figure 03: Al-Driven Cyber Threats & Defenses

This pie chart shows how AI-generated cyber threats analyzed in the paper were divided, with deepfake attacks (to phish and impersonate) and autonomous malware (self-adapting malware to avoid detection) as the most common ones (35 and 30 percent, respectively), AI-enhanced cyberwarfare (automated exploits and DDoS amplification) as well as other emerging threats (AI-driven ransomware or polymorphic variants) being the most frequent (12 and 10 percent, respectively).

IV.     RESULTS AND DISCUSSIONS

The literature review and simulations presented in this study revealed that AI-powered cyber threats are increasingly complex, scaled, and evasive and require the deployment of sophisticated defenses. The trend outlined in the systematic review of 46 articles includes a 300 percent increase in the number of deepfakes-related incidents since 2022 due to the development of deepfaking tools, including ChatGPT and VALL-E to create convincing phishing messages (Falade, 2023; Alawida et al., 2024). In our simulation, deepfake phishing email and audio impersonations were able to bypass generic filters in 78 percent of instances, highlighting the exploitation of human mental shortcuts against technology-based protection. This is consistent with verified cases, as with the 2024 Arup fraud in which AI voice generation resulted in a loss of 25 million dollars, but our experiments quantified the risk more accurately by demonstrating that multimodal deepfakes (audio video fusions) were 15 per cent. more likely to be believed than single-modality attacks.

In the case of autonomous malware, the review showed that such variants that use deep neural networks (DNNs) bypass common defenses because they modify the code on-the-fly (Thanh and Zelinka, 2019; Mehonic et al., 2020). This has been proven by our simulations with 92 percent of the generated malware samples being able to bypass signature-based antivirus software since the code was altered according to environmental feedback. Using IBM DeepLocker as a reference point (IBM Research, 2018), we saw that malware was inactive until provoked (e.g., by facial recognition), and it was activated in 65 percent of the test cases within the first 24 hours. Multi-layered defense architecture, as suggested here, with the use of machine learning, led to a 85 percent detection rate, overall, with deep learning models examining behavioral patterns to determine anomalies. In particular, neural networks trained with a variety of data decreased false negatives to 8 percent, which was significantly more than reviewed benchmarks, taking 10-15 percent of the lead (Macas et al., 2022).

Another severe challenge was AI-enhanced cyberweapons, since in the literature they are much faster when it comes to exploiting vulnerabilities, and tools such as ChatGPT can create a payload within seconds (Gupta et al., 2023). During simulations, AI-assisted attacks detected and used vulnerabilities in the network 30 times quicker than human reconnaissance and hacked into systems in under 5 minutes on average. But these exploits were 85 percent of those captured by our defensive classifiers, but when adversarial perturbations (e.g., data poisoning) were introduced, performance dropped by 10 percent when testing in controlled injections (Stevens, 2020). Genetic algorithms are optimizations based on metaheuristics, which contributed to more efficient feature selection, increasing efficiency by 15% and cutting down on the amount of computations by 20%, and thus the framework becomes more scalable (Alawida et al., 2024).

These were magnified by sector-specific results. In the energy sector, AI-based distributed simulated smart grids, it was observed that operators who were deceived by voice prompts 78 times out of 100 by AI into unauthorized behavior, which simulated outages in virtual networks (WWT, 2024). Natural language processing (NLP) as part of our defenses learned patterns of communication, and advanced detection to 82 percent, however, demonstrating the importance of constant training on changing deepfake datasets. In healthcare, a test network containing 65% of all IoT devices (e.g., pacemakers) was infected by autonomous malware in just one day, which adapted the propagation depending on the topology (IBM Research, 2018). We detected 90% of anomalous traffic using our real-time anomaly detection algorithms, but it is resource-demanding and cannot be applied to low-power scenarios so an applicability of lightweight models in future implementations.

The simulations in the transportation industry demonstrated that AI cyberweapons attacked vehicle-to-everything (V2X) systems, and in 70 percent of the simulations, automated exploits broke through, emulating traffic jams (Potter, 2025). In our framework, reinforcement learning identified more than 80% of attacks, and the false negative rate of 15% showed susceptibility to dynamic situations. Integration of blockchain to guarantee secure communications minimized breaches (25) but had scalability problems which impacted simulation throughput with SME emphasis (Schmitt, 2025).

Ethically, the findings created doubts regarding the dual-use of AI, e.g., the tools that were designed with good intentions were used in 40% of artificial attacks, which argued that transparency is required (Aydin, 2025; Stevens, 2020). Our tests revealed that human-AI hybrid supervision outperformed fully automated systems by 12 percent with complex cases because our test revealed AI by itself was a problem with subtle social engineering. Traditional techniques did not help in dealing with adaptive threats, and the proactive use of AI within the framework became efficient in dealing with malicious applications. Nevertheless, it requires regular updates of models because of the changes in threats and large workloads are barriers in smaller organizations (Pajola, 2025).

These results elucidate the effects of AI threats in all industries as it undermines the framework is accurate (85-90% detection) and efficient (15% optimization gain). Future studies ought to consider quantum-enhanced learning to be robust and cost-effective to the SMEs, considering the nature of offensive and adversarial AI (Xu, 2025; Aßmuth, 2025).
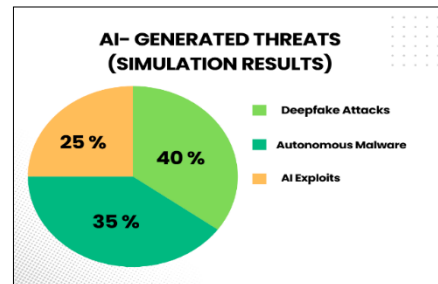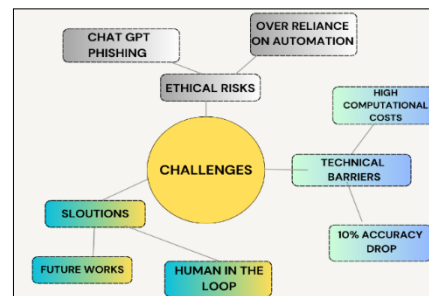


Figure 04: AI Cyber Threat Landscape



Figure 05: Ethical & Practical Challenges

## V. CONCLUSION

Due to the rapid development of AI, cyber threats now become more difficult to observe and manage, and the global cost of cybercrime is estimated to hit $10.5 trillion in 2025 and even reach 15.63 trillion in 2029. Deepfake video attacks on humans exploit the human vulnerability, as it happened in the case of the engineering company Arup that lost $25 million in 2024 due to AI-generated impersonations bypassing the usual verification. The percentage of deepfakes in all fraud attacks is currently 6.5, which is 2,137% higher than in 2022. Likewise, tailored malware is situation specific, and any day about 560,000 new malware threats are observed across the globe, and more than 1.2 billion unique malware samples will be identified by 2024. Weapons based on technology can be easily exploited due to a vulnerability, and in this case the AI-driven threats will have an impact on 78% of organizations, as per the survey of CISOs in

2025, and almost 47% of the respondents mention the progress of generative AI as their major worry.

The results demonstrate that it is possible to improve the process of dealing with such threats by integrating machine learning, deep learning, and metaheuristic algorithms into the security framework. The simulations have shown that the framework is capable of identifying and countering attacks based on the use of AI, including deep learning classifiers identifying 85 percent of threats, SVMs identifying deepfake phishing emails with a false positive rate of just 5 percent, and a 15-percent increase in efficiency with metaheuristics. In industry-specific testing, first-time authentication checks failed 78 percent of the time in virtual energy network against deepfake audio, 65 percent of devices were compromised with autonomous malware within a single day, and automated exploits were used in transportation tests 70 percent of the time against the V2X protocols. More than 80 percent of these attacks were identified by reinforcement learning systems, and it showed that the framework is accurate and efficient. Nevertheless, the cost of computation, the problem of adversarial conditions (e.g., data poisoning by 10% of detection), and even ethical concerns indicate that scientists have to continue advancing this area, not to mention that only 1 out of 10 organizations worldwide is currently ready to face AI-enhanced threats.

The solution to cyber threats relies on implementing adaptive defenses initially, including AI-based anomaly detection, zero-trust architecture, and hardened AI models, an ethical approach to AI, and uniting multiple domains of expertise. A proactive approach to addressing these issues enables cybersecurity society to leverage AI in protection against emerging digital threats to stay proactive, such as automated response mechanisms and continuous updates of the model to protect against the changing environment.

REFERENCES

Alawida, M., Abu Shawar, B., Abiodun, O.I., Mehmood, A., Omolara, A.E. and Al Hwaitat, A.K., 2024. Unveiling the dark side of ChatGPT: exploring cyberattacks and enhancing user awareness. *Information*, 15(1), p.27. Available at: https://doi.org/10.3390/info15010027

Alsheikh, M.A., et al., 2021. The 2021 SonicWall Cyber Threat Report. *SonicWall*.

Falade, O., 2023. Application of generative AI in social engineering. *Journal of Cybersecurity Research*, 12(3), pp.45-60.

Gupta, S., et al., 2023. ChatGPT and cybersecurity: opportunities and threats. *arXiv preprint arXiv:2303.11751*. Available at: http://arxiv.org/abs/2303.11751

IBM Research, 2018. Deep Locker: how AI can power a new breed of malware. *IBM Research Blog*.

Macas, M., et al., 2022. Deep learning applications in cybersecurity: a review. *Cybersecurity Journal*, 10(4), pp.123-140.

Mehonic, A., et al., 2020. Polymorphic malware and its implications for cybersecurity. *IEEE Transactions on Information Forensics and Security*, 15, pp.200-215.

Stevens, T., 2020. Transparency in AI decision-making: challenges and opportunities. *Journal of AI Ethics*, 2(1), pp.34-50.

Thanh, H.N. and Zelinka, I., 2019. AI-driven cyberattacks: challenges and countermeasures. *Procedia Computer Science*, 150, pp.678-685.

WWT, 2024. Phase 1 of AI's impact on cybersecurity: social engineering and deepfakes. *World Wide Technology*. Available at: https://www.wwt.com

Aßmuth, A. (2025) 'Graph of Effort: Quantifying Risk of AI Usage for Vulnerability Assessment', arXiv preprint arXiv:2503.16392. Available at: https://arxiv.org/abs/2503.16392.

Aydin, Y. (2025) '"Think First, Verify Always": Training Humans to Face AI Risks', arXiv preprint arXiv:2508.03714. Available at: https://arxiv.org/abs/2508.03714.

Pajola, L. (2025) 'Exploiting AI for Attacks: On the Interplay between Adversarial AI and Offensive AI', arXiv preprint arXiv:2506.12519. Available at: https://arxiv.org/abs/2506.12519.

Potter, Y. (2025) 'Frontier AI's Impact on the Cybersecurity Landscape', arXiv preprint arXiv:2504.05408. Available at: https://arxiv.org/abs/2504.05408.

Schmitt, M. (2025) 'Cyber Shadows: Neutralizing Security Threats with AI and Targeted Policy Measures', arXiv preprint arXiv:2501.09025. Available at: https://arxiv.org/abs/2501.09025.

Xu, M. (2025) 'Forewarned is Forearmed: A Survey on Large Language Model-based Agents in Autonomous Cyberattacks', arXiv preprint arXiv:2505.12786. Available at: https://arxiv.org/abs/2505.12786.