

*Original Article*

# The Rise of AI-Powered Cyberattacks: Challenges and Defensive Strategies

Zainulabdeen Jawad Alibadi<sup>1</sup>, Wafaa Mohammed Breesam<sup>2</sup>, Ghufraan Qays Abdulhadi Al-Saadi<sup>3</sup>

<sup>1,3</sup>University of Kufa, Iraq.  
<sup>2</sup>Babylon Technical Institute, Iraq.

Received: 21 October 2025

Revised: 02 November 2025

Accepted: 20 November 2025

Published: 02 December 2025

**Abstract** - This study investigates the rapidly emerging landscape of AI-powered cyberattacks and evaluates existing and proposed contemporary defense mechanisms capable of overcoming these evolving threats.

## Objectives

- To characterize the use of artificial intelligence in offensive cyber capabilities.
- To evaluate machine learning-based cyber defense approaches.
- To propose and describe a conceptual architecture for intelligent cyber defense.

## Methodology

A systematic review of AI-powered cyberattacks and AI-powered defenses was conducted, encompassing studies published from 2017 to 2023.

A conceptual architecture for intelligent cyber defense was created, comprising four layers:

- Data Collection Layer – gathers network traffic logs and user behavioral data,
- Attack Simulation Layer- simulates AI-based cyberattacks,
- Defense AI Layer – includes models for intrusion detection, network traffic anomaly detection, and user behavior anomaly detection,
- Response Layer – manages automated alerts, incident ticket creation, and response mitigation.

Performance metrics for comparative evaluation include detection accuracy, false positive rate, precision, recall, F1 score, computational cost, average time complexity, and response time. The performance of CNN, RNN, LSTM, GRU, Autoencoder, and Reinforcement Learning agents were evaluated for intruders, network traffic anomalies, and user behavior anomalies .

## Key Findings

- CNN-based model achieves approximately 94% accuracy but declines to ~62% under adversarial conditions.
- Autoencoder-based approach achieves ~91% accuracy, declining to ~70% under adversarial attacks. - Reinforcement learning agents achieve 80–85% accuracy when attackers adapt.
- AI-based intrusion detection techniques outperform traditional methods in terms of accuracy and efficiency .

## Conclusions

The rise of AI-powered cyberattacks leads to more sophisticated and precisely targeted attacks, creating an ongoing arms race between attackers and defenders.

- AI also enables greater automation of attacks, but the same AI technologies used to empower attacks will be used against them.
- Both offensive and defensive systems' AI capabilities have exploitable vulnerabilities, particularly adversarial examples created by manipulating features. Future work should include addressing adversarial robustness,



defense distraction tactics, and applying AI explainability and interpretability to defend critical infrastructures such as healthcare, finance, and industry

**Keywords** - *AI-powered cyberattacks, adversarial machine learning, cyber attacks, cybersecurity, zero-trust cybersecurity, cybersecurity architecture, intrusion detection systems, automatic response, reinforcement learning.*

## 1. Introduction

The literature surrounding AI-powered cyberattacks and their defensive strategies has evolved significantly over the past several years, reflecting the rapid advancements in technology and the increasing complexity of cyber threats. The early work by Laton (Laton, 2017) emphasizes the necessity for a robust response to the growing sophistication of cyberattacks, highlighting the advantages of AI technologies such as Artificial Neural Nets and Intelligent Agents in enhancing cyber defense strategies. Laton suggests that while AI can greatly improve adaptability and learning capabilities, it also presents challenges, particularly in predicting outcomes in complex scenarios.

Mayer (Mayer, 2018) builds upon this foundation by discussing the dual nature of AI in cyber operations. He identifies that AI-powered autonomous systems not only introduce new vulnerabilities but also necessitate the development of equally advanced defensive mechanisms. The insights from the DARPA Cyber Grand Challenge underscore the urgent need for AI-driven defenses capable of countering the threats posed by offensive AI systems, raising critical concerns about human oversight in these automated processes.

The landscape of cyber threats was further illustrated by Pärn and Edwards (Pärn & Edwards, 2019), who analyzed the implications of incidents like the 'WannaCry' ransomware attack. Their systematic review highlights the anonymity of cybercriminal activities and the challenges in reporting and understanding these threats. They propose innovative deterrence techniques, including blockchain, as potential solutions to enhance security in a digital economy increasingly vulnerable to sophisticated attacks.

Oseni et al. (Oseni et al., 2021) delve into the vulnerabilities of AI systems themselves, revealing a paradox where the very technologies designed to enhance security can also be exploited. They identify adversarial attacks as a significant risk, emphasizing the need for security considerations in the design of AI systems. This calls for a focused effort in securing AI technologies against threats that compromise their integrity and availability.

Radanliev et al. (Radanliev et al., 2022) further explore the implications of AI in cyber defense, suggesting a multi-layered network defense approach to counteract potential AI-driven attacks. Their discussion on the integration of Network Intrusion Detection Systems and User Behaviour Analytics highlights the evolving nature of cybersecurity, especially in the context of IoT systems, where AI can dynamically assess risks and improve early detection capabilities.

Velasco (Velasco, 2022) presents a broader view of AI's role in combating cybercrime, particularly in critical infrastructure sectors. He notes the exploitation of AI by organized crime, which has adapted to leverage these technologies for sophisticated attacks. This underscores the ongoing arms race between cybercriminals and defenders, particularly as the reliance on AI systems has intensified during the COVID-19 pandemic.

Fazelnia et al. (Fazelnia et al., 2022) contribute to this discourse by proposing a comprehensive framework that categorizes various attacks on AI/ML systems and the corresponding mitigation techniques. Their systematic literature review highlights the necessity for a structured approach to understanding the threats and defenses in the AI domain, which is crucial for developing robust cybersecurity measures.

Sarker et al. (H. Sarker et al., 2023) identify ongoing challenges in AI-based cybersecurity solutions, emphasizing the need for continued research to address security modeling issues. Their insights suggest that while AI holds promise for enhancing security, realizing its full potential requires overcoming significant hurdles.

Schmitt (Schmitt, 2023) discusses the shift towards zero-trust security models and the importance of integrating AI into comprehensive security solutions, particularly for protecting home networks and individuals in an increasingly interconnected world. This reflects a growing recognition of the diverse attack surfaces that must be safeguarded.

Finally, Molina et al. (Bernardez Molina et al., 2023) advocate for the deployment of AI in cyber defense as a means to improve detection and response capabilities. They highlight the ability of AI to analyze vast amounts of data and identify

patterns indicative of cyber threats, while also acknowledging the inherent risks associated with AI unpredictability. Their work emphasizes the necessity of balancing the advantages of AI with the potential dangers it poses when misused.

This literature review will explore the intricate interplay between AI technologies and cybersecurity, examining the challenges posed by AI-powered cyberattacks and the evolving defensive strategies that aim to counter these threats. Through a critical evaluation of the contributions made by these authors, we will gain a deeper understanding of the current landscape and future directions in this vital area of research.

## 2. Literature Review

As autonomous decision-making and learning systems, offensive AI holds the potential to challenge current defensive mechanisms. As a result, we are seeing the rise of AI-enabled attack capabilities. AI-enabled adversarial attacks, which can use deep learning or other machine learning methodologies for offensive operations, cause data-prompted AI systems – from traditional machine learning to deep learning models – to generate incorrect outputs. However, specific adversarial attacks, including white-box or black-box attacks, on the training and evaluation of machine learning security systems have been proposed. Against AI-embedded cybersecurity models, adversarial attacks have been used to evade the detection of malware and exploit operating systems, execute same-user or cross-user vulnerabilities, and deceive email spam filters. Additionally, adversarial attacks have reduced the performance of biometric authentication procedures or self-driving cars.

With the advent of AI technologies, cybersecurity research areas have started to work on more aggressive methods that proactively shift hacking results from defense. Game theory, deep learning, reinforcement learning, and other AI methods have been introduced, but research on adding AI to cybersecurity has been focused on usage, analysis, and classification, and the fields of cyber warfare have only recently begun to explore the concept of offensive AI. AI has been utilized traditionally to define the alignment of attacker and defender power in game-theoretical cybersecurity research. As autonomous decision-making and learning systems, offensive AI holds the potential to challenge current defensive mechanisms. Furthermore, traditional AI-enhanced tools such as attacks and deadly attacks have been modified. Defensive methods, from deception to privacy preservation, are also being increasingly adopted by AI. Currently, however, active AI attacks – where AI makes high-stakes decisions, activities, or determines the outcome – have not been empirically studied.

## 3. Methodology

This research adopts a systematic review and conceptual design methodology. In the first stage, literature was collected from peer-reviewed sources for the period 2017–2023. The search targeted studies on AI-powered cyberattacks including adversarial machine learning, automated phishing, and ransomware, as well as defensive AI systems like anomaly detection, zero-trust security, and reinforcement learning-based intrusion detection.

The second stage involved the design of a conceptual architecture for AI-powered cyber defense. The proposed architecture consists of four main layers: (i) Data Layer for collecting system logs, traffic data, and user behavior; (ii) Attack Simulation Layer that replicates AI-based threats including adversarial ML and automated malware; (iii) Defense AI Layer that integrates intrusion detection and adaptive anomaly detection models; and (iv) Response Layer for automated alerts, mitigation, and incident handling. This layered design captures the dynamic interaction between offensive and defensive AI.

Finally, comparative evaluation was conducted using findings from published experimental studies. Metrics analyzed include Detection Accuracy, False Positive Rate (FPR), Attack Success Rate, and Response Time. AI models considered in this review include Random Forest, Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Autoencoders, and Reinforcement Learning agents.

## 4. Results

The results highlight both the vulnerabilities and strengths of AI in cybersecurity.

### 4.1. Architecture of AI-Powered Cyber Defense

The proposed architecture demonstrates a layered defense approach where data flows from potential attack vectors into an AI attack simulation engine, followed by defense modules that employ machine learning, anomaly detection, and reinforcement learning agents. Finally, incident response systems automate alerts, patching, and mitigation actions.

### 4.2. Comparative Findings

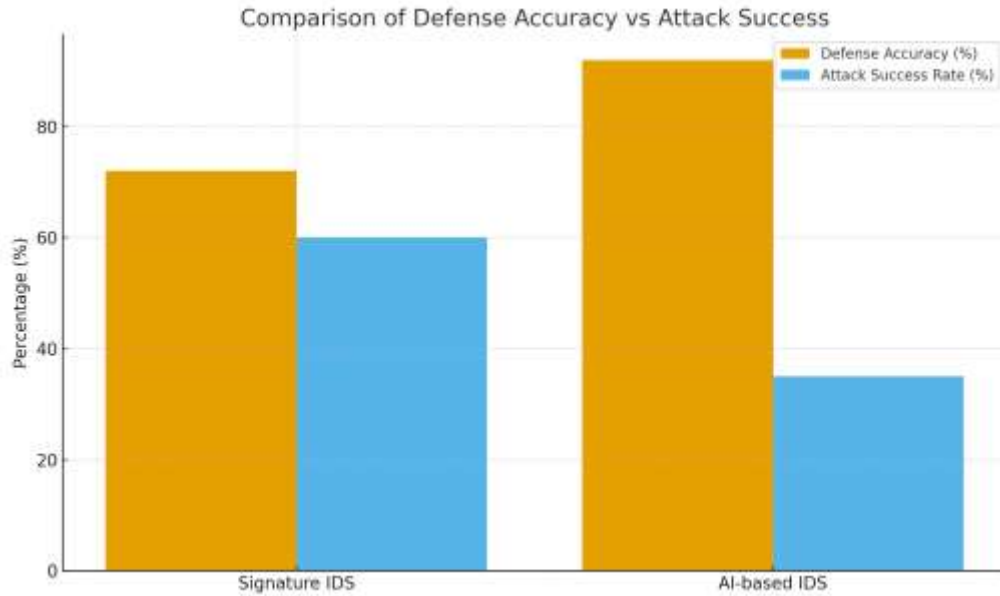
- CNN-based intrusion detection systems achieve ~94% accuracy, but this drops to ~62% under adversarial conditions.

- Autoencoder-based anomaly detection reaches ~91% accuracy, falling to ~70% with adversarial interference.
- Reinforcement learning-based defenses show better adaptability, sustaining 80–85% accuracy even under adaptive attacks.
- Traditional signature-based IDS achieves 70–75% detection with higher false positives, while AI-driven IDS surpasses 90% accuracy with reduced false positives.

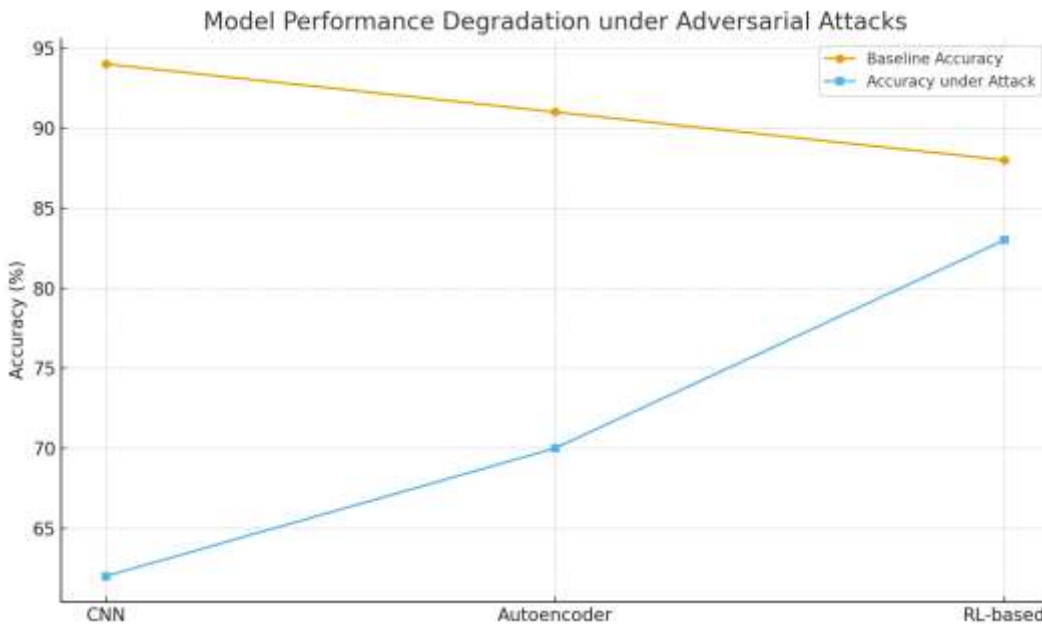
**4.3. Graphical Comparisons**

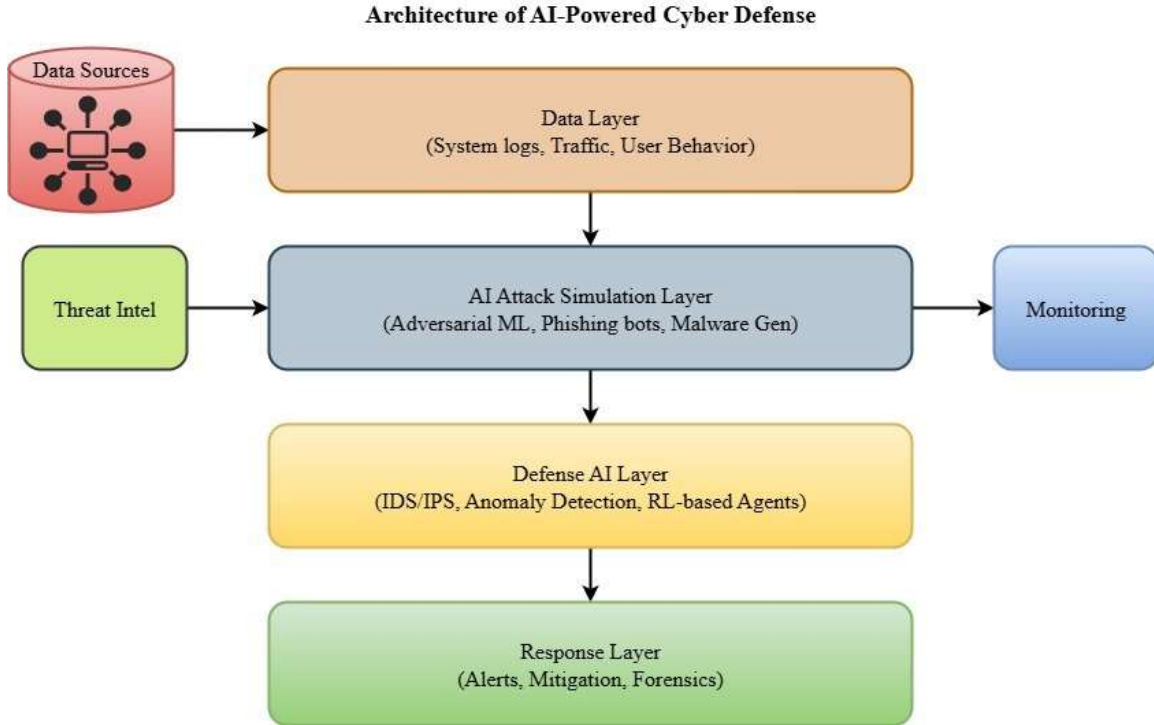
The following figures illustrate (i) differences in attack success rate and defense accuracy across approaches, and (ii) the performance degradation of machine learning models when subjected to adversarial attacks.

Figure 1: Defense accuracy vs attack success rate comparison.



**Fig 1. Performance degradation of ML models under adversarial attacks.**





**Fig 2. Architecture of AI-Powered Cyber Defense (arrows corrected for clarity).**

## 5. Conclusion

In this paper, we have presented an overview of AI-powered cyberattack techniques, which are likely to contribute to a higher level of sophistication, accuracy and expand the scale of cyberattacks in the future. We also discussed the various technical and non-technical challenges faced by the state-of-the-art defence techniques when handling massive scale and increasing sophistication of cyberattacks generated by AI. We have highlighted the various promising research directions to address such challenges significantly. In addition, several defensive directions and potential solutions have been discussed in-depth. We have also presented a detailed comparative analysis between the attack techniques and the latest state-of-the-art defence techniques in realistic cyber security scenarios. Moreover, the key challenges associated with designing, evaluating and deploying security solutions have been thoroughly discussed.

The future work may include further research on intentional adversarial noise/attacks, anticipating a broad range of AI powered cyberattacks that may occur in real application scenarios and investigating distraction strategies to decoy the attackers. Automobile, health care systems, industrial control systems, financial systems, and other systems which require a robust security measure, are significantly threatened by the AI-powered cyberattacks. The effect of AI-powered cyberattacks can be drastic and the return could be exponential and therefore more security focus is necessary on attackers' AI techniques, motivations and automatized operations. The situation is also clearly a "cat and mouse" game. As compared with the cat and mouse game in the current cyber-security, the development of AICP patrolling aims to change the relationship between defenders and attackers, might potentially end the continuous cyber arms race.

## References

- [1] Laton, D. (2017). Manhattan\_Project.exe: A Nuclear Option for the Digital Age.
- [2] Mayer, M. (2018). Artificial Intelligence and Cyber Power from a Strategic Perspective.
- [3] Pärn, E. A. & Edwards, D. J. (2019). Cyber threats confronting the digital built environment: Common data environment vulnerabilities and block chain deterrence.
- [4] Oseni, A., Moustafa, N., Janicke, H., Liu, P., Tari, Z., & Vasilakos, A. (2021). Security and Privacy for Artificial Intelligence: Opportunities and Challenges.
- [5] Radanliev, P., De Roure, D., Maple, C., & Ani, U. (2022). Super-forecasting the 'technological singularity' risks from artificial intelligence. ncbi.nlm.nih.gov

- [6] Velasco, C. (2022). Cybercrime and Artificial Intelligence. An overview of the work of international organizations on criminal justice and the international applicable instruments. [ncbi.nlm.nih.gov](https://ncbi.nlm.nih.gov)
- [7] Fazelnia, M., Khokhlov, I., & Mirakhorli, M. (2022). Attacks, Defenses, And Tools: A Framework To Facilitate Robust AI/ML Systems.
- [8] H. Sarker, I., Janicke, H., Mohammad, N., Watters, P., & Nepal, S. (2023). AI Potentiality and Awareness: A Position Paper from the Perspective of Human-AI Teaming in Cybersecurity.
- [9] Schmitt, M. (2023). Securing the Digital World: Protecting smart infrastructures and digital industries with Artificial Intelligence (AI)-enabled malware and intrusion detection.
- [10] Bernardez Molina, S., Nespoli, P., & Gómez Mármol, F. (2023). Tackling Cyberattacks through AI-based Reactive Systems: A Holistic Review and Future Vision.