(REVIEW ARTICLE)

Check for updates

# The role of AI and machine learning in cybersecurity: Advancements in threat detection, anomaly detection and automated response

Aminat Bolaji Bello [1], Akeem Olakunle Ogundipe [2, *], Awobelem A. George [3] and Olabode Anifowose [4]

[1] Department of Mathematical Science, Adekunle Ajasin University, Ondo, Nigeria.
[2] Department of Management Information Systems, Lamar University, Texas, USA.
[3] Jack H. Brown College of Business & Public Administration, California State University, California, USA
[4] Department of Mechanical Engineering, Georgia Southern University, Georgia, USA

## Abstract

The increasing complexity and frequency of cyber threats have prompted organizations to seek more sophisticated defense mechanisms. Traditional signature-based methods and manual threat-hunting processes often fall short against evolving malware, zero-day exploits, and social engineering techniques. Artificial Intelligence (AI) and Machine Learning (ML) have emerged as pivotal tools, enabling automated threat detection, real-time anomaly analysis, and proactive incident response. This review synthesizes current research and practices related to AI-driven cybersecurity, examining supervised and unsupervised learning for threat detection, AI-powered anomaly detection, and real-world industrial applications. The discussion also explores ethical considerations such as adversarial AI and bias, concluding with future directions that include quantum-safe cryptography, AI-augmented security operations centers, and the integration of blockchain for enhanced cybersecurity.

**Keywords:** Cybersecurity; Artificial Intelligence; Machine Learning; Deep Learning; Phishing

## 1. Introduction

The exponential growth of digital services and the global shift toward cloud-based computing have reshaped the cybersecurity landscape, exposing modern organizations to a host of increasingly sophisticated and persistent cyber threats. Historically, many security solutions relied on signature-based detection matching known malicious patterns or hash signatures of malware binaries [1]. While this method proved effective against common threats, the emergence of more advanced techniques such as polymorphic malware, zero-day exploits, and stealthy Advanced Persistent Threats (APTs) has revealed limitations in static or rule-based defenses [2]. Attackers have also begun to exploit the interconnectedness of supply chains and the vulnerabilities inherent in remote work environments, further amplifying the attack surface.

Against this backdrop, Artificial Intelligence (AI) and Machine Learning (ML) have gained prominence as potentially transformative tools in cybersecurity. By ingesting large volumes of data ranging from network logs to endpoint telemetry machine learning models can discern intricate patterns, detect anomalies, and generate predictive insights at a scale and speed beyond the capacity of human analysts [3]. This shift mirrors a broader trend in the technology industry, where AI-driven automation is increasingly valued for its ability to reduce operational overhead, adapt to novel threats, and minimize the mean time to detect (MTTD) and respond (MTTR) to incidents. Successful implementations span multiple domains, from user and entity behavior analytics (UEBA) to automated threat hunting and intelligent firewalls that dynamically adjust to changing network conditions [4].

---

* Corresponding author: Akeem Olakunle Ogundipe.

Beyond enhancing detection and response, AI also offers capabilities for proactive risk management. AI-powered threat intelligence platforms can aggregate open-source intelligence, dark web data, and social media signals to forecast emerging trends or identify sophisticated adversarial campaigns before they fully materialize [5]. These proactive measures align well with the growing emphasis on cyber resilience the notion that merely reacting to breaches is insufficient in a climate where attacks have become unrelenting and ever-evolving. However, integrating AI tools into a cybersecurity framework does not come without challenges. Adversarial AI, where threat actors intentionally manipulate inputs to fool machine learning models, highlights a new frontier in the arms race between defenders and attackers [6]. Biases within training datasets can also lead to false positives that disrupt legitimate user activities or false negatives that allow stealthy threats to slip through undetected [7].

In parallel to these technical considerations, the ethical dimensions of AI-driven defense strategies are increasingly scrutinized. Questions arise about how much autonomy these systems should have in executing mitigation actions, especially when human intervention may be required to avoid unintended consequences in critical infrastructure or sensitive healthcare setting [8]. Moreover, transparent explainability of AI decisions is critical for fostering trust among stakeholders, including security analysts, IT teams, and upper management. If an AI system erroneously blocks vital services or misclassifies user behavior, organizations must be able to investigate and rectify the underlying logic.

Given these developments, this review aims to illuminate the current state-of-the-art in AI-driven cybersecurity by examining foundational ML techniques, discussing approaches for threat detection and anomaly identification, and exploring real-world industry applications. The discussion further delves into the ethical and technical challenges adversarial AI, bias, and explainability offering insights into how researchers and practitioners can navigate these pitfalls. Finally, the review highlights emerging directions, such as quantum-safe cryptography and AI-augmented Security Operations Centers (SOCs), as potential avenues for advancing cybersecurity defenses in an era marked by constant technological disruption. Ultimately, the strategic deployment of AI and ML in cybersecurity may be pivotal in turning the tide against rapidly evolving adversaries, bolstering digital trust, and safeguarding critical assets and infrastructure.

## 2. AI and Machine Learning Techniques in Cybersecurity

### 2.1. Supervised vs. Unsupervised Learning in Threat Detection

The distinction between supervised and unsupervised learning is often considered a foundational concept in machine learning, particularly in the cybersecurity domain. Supervised learning relies on labeled datasets, where examples of benign and malicious activities are known [9, 10]. This approach is effective in environments where historical records provide a wealth of accurately tagged attack signatures or malware samples. Decision Trees, Random Forests, and Support Vector Machines (SVM) frequently serve as the backbone of supervised threat detection systems [11]. For instance, an SVM can be trained on known malicious executable files and their benign counterparts, learning to distinguish between them based on extracted features such as opcode frequency, file metadata, or API call sequences. Once trained, the model can identify suspicious files or network traffic patterns with high accuracy assuming the data distribution remains stable and the adversaries do not drastically alter their attack tactics [12].

However, a key limitation of supervised models lies in their reliance on the quality and breadth of labeled training data. In practice, labeling large-scale security datasets is resource-intensive, and newly emerging attacks (e.g., advanced persistent threats or zero-day exploits) may not appear in historical records [13, 14]. As a result, supervised models can underperform against novel or obfuscated threats, driving the need for more flexible or adaptive strategies.

In contrast, unsupervised learning does not rely on pre-labeled data. Instead, it identifies deviations from a learned baseline of "normal" activity, flagging outliers as potential intrusions. Techniques such as autoencoders, clustering (k-means, DBSCAN), or principal component analysis (PCA) have proved valuable for anomaly detection in network logs, system call traces, and user behavior analytics [15, 16]. By modeling the "typical" state of a system, unsupervised methods can uncover unknown attack vectors often critical in detecting zero-day exploits or stealthy attacks that evolve over time. Nonetheless, the high sensitivity of these models can lead to false positives, generating an influx of alerts that overwhelm security analysts [17]. Tuning hyperparameters to balance detection accuracy and alert volume remains an ongoing challenge, particularly in large-scale production environments where data can shift rapidly [18]. Notably, hybrid approaches are emerging that combine supervised and unsupervised learning. for example, an unsupervised anomaly detection model might first flag unusual activity, which is then classified by a downstream supervised model (or verified by human experts) to confirm its malicious or benign nature [19]. For such hybrid frameworks seek to optimize detection coverage while minimizing false positives, bridging the gap between purely supervised and purely unsupervised paradigms in threat detection.

## 2.2. Deep Learning for Cybersecurity

The rise of deep learning techniques has further transformed machine learning applications in cybersecurity, enabling automated feature extraction and often outperforming traditional approaches in complex tasks. Convolutional Neural Networks (CNNs) have demonstrated success in analyzing structured and semi-structured data ranging from network packet payloads to log files by detecting spatial patterns or repetitive signatures that might be missed by manual feature engineering [20, 21]. For instance, a CNN can treat raw byte sequences of network traffic as "images," allowing the model to learn local patterns that distinguish benign traffic from attacks such as DDoS, port scans, or data exfiltration.

Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) variants, excel at capturing temporal patterns, making them well-suited for intrusion detection in sequential data like netflow records, command-line arguments, or system calls [22]. By retaining information over time, LSTMs can model how malware or attack scripts evolve during an ongoing intrusion, potentially catching advanced persistent threats that unfold in multiple stages. As an example, a well-trained LSTM might detect subtle signs of lateral movement within a network, issuing alerts before an attacker can escalate privileges or exfiltrate large volumes of data.

More recently, transformer-based models have emerged to address cybersecurity tasks involving textual analysis, such as phishing email detection or malicious URL classification [23]. By leveraging attention mechanisms, transformer models contextualize relationships within input text, outperforming earlier architectures on a wide range of natural language processing tasks. In the cybersecurity context, these models can parse the semantics of email content, spot domain spoofing, and identify language patterns indicative of social engineering [ 24].

Despite the promise of deep learning, deploying these models in real-world settings requires addressing computational overhead, potential overfitting, and data availability. Training large neural networks demands substantial GPU or TPU resources, while real-time inference may call for dedicated hardware accelerators or optimized model compression techniques [25]. Additionally, the model's performance can degrade if the production data distribution changes significantly over time, emphasizing the importance of continuous retraining or incremental learning strategies [26].

## 2.3. Reinforcement Learning in Cybersecurity

Whereas supervised and unsupervised models strive to classify or detect anomalies within static datasets, reinforcement learning (RL) takes a more dynamic approach. Here, an RL agent learns by interacting with an environment such as a simulated network or honeypot and receives rewards based on the efficacy of its actions [27, 28]. For example, quarantining an infected host or blocking a malicious IP might yield a reward if it halts an attack, whereas erroneously blocking legitimate business traffic might incur a penalty.

Early studies on RL in cybersecurity suggest that agents can autonomously discover optimal strategies for intrusion detection and incident response. For instance, an agent could be trained to sequentially deploy security tools like sandboxing suspicious executables or segmenting a network zone depending on the threat level [29]. By continually updating its policy, the RL system can adapt to emerging attack tactics or changes in the network topology. This adaptive feature is particularly appealing for zero-day threats, where static or signature-based solutions might fail.

However, practical deployment of RL-based systems remains limited. Training in a live environment poses substantial risks if experimental actions cause network disruptions or inadvertently overlook critical threats [30]. Creating realistic simulations or cyber ranges is thus a prerequisite for RL research, but building and maintaining these complex environments can be time-consuming and costly. Additionally, RL agents often require extensive trial-and-error interactions to converge on optimal policies, which may be impractical for rapidly shifting or mission-critical operations [31]. Overcoming these limitations necessitates refining RL frameworks potentially by incorporating offline learning from historical data, combining RL with supervised components for safer exploration, or adopting hierarchical RL architectures that abstract high-level policies.

# 3. AI-Powered Anomaly Detection and Threat Identification

## 3.1. Behavioral Analytics for Threat Detection

Behavioral analytics has become a cornerstone in modern cybersecurity, offering a proactive means of identifying potential attacks before they escalate into severe breaches. Under the umbrella of User and Entity Behavior Analytics (UEBA), AI-driven systems collect and interpret a wide range of activity data such as user login times, data access volumes, file transfers, and cross-domain communications [32]. These insights help establish a dynamic baseline of "normal" behavior within an organization. Any significant deviation from this baseline, such as an employee

downloading an unusually large volume of files during non-business hours, prompts an alert that may warrant further investigation.

What sets behavioral analytics apart from traditional rule-based systems is its ability to adapt to evolving organizational patterns. For instance, if an enterprise adopts new workflows or integrates remote employees in different time zones, UEBA models can gradually recalibrate normal activity thresholds, minimizing false positives [33]. Additionally, these models can detect complex insider threats, where malicious activity may blend in with legitimate user actions. By continuously learning and correlating multiple data points from physical badge swipes to network file shares, behavioral analytics provides a high-fidelity view of potential risk factors. However, successful implementation demands robust data pipelines, consistent labeling for normal vs. abnormal activities, and careful consideration of privacy implications, particularly when monitoring employee actions in regions with strict data protection laws [34, 35].

## 3.2. Network Intrusion Detection Systems (NIDS)

Network Intrusion Detection Systems (NIDS) serve as the first line of defense against external threats, monitoring ingress and egress traffic in real time. As network architectures grow more complex encompassing on-premises data centers, public clouds, and remote endpoints AI-enabled NIDS have become indispensable. These systems parse massive volumes of packet-level data, applying advanced machine learning or deep learning algorithms to highlight anomalies indicative of malicious behavior [36]. By scrutinizing traffic patterns, AI-based NIDS can detect subtle shifts that might indicate an evolving Distributed Denial of Service (DDoS) attack, a stealthy APT infiltration, or data exfiltration attempts.

Organizations like IBM have integrated AI modules into their Security Operations Center (SOC) offerings, exemplified by solutions such as IBM Watson for Cybersecurity. By automating the analysis of logs and threat intelligence feeds, Watson can correlate alerts from disparate sources firewalls, intrusion detection systems, and endpoint protection platforms offering holistic insights into potential breaches [37]. Some NIDS solutions also incorporate Intrusion Prevention Systems (IPS), providing an automated, proactive defense mechanism that can block or divert malicious traffic as soon as it is flagged [38]. While these capabilities reduce the time to containment, they also introduce risks if the IPS accidentally terminates legitimate connections, emphasizing the need for finely tuned policies and periodic reviews of false positive rates.

## 3.3. Phishing and Social Engineering Detection

Despite advancements in perimeter defenses, phishing remains a leading cause of security incidents, with attackers capitalizing on human error through deceptive emails, websites, or messages [39]. AI-driven solutions address this persistent threat by employing Natural Language Processing (NLP) and sophisticated patt ern recognition to detect nuanced indicators of malicious intent. Transformer-based models, such as BERT or GPT variants, excel at contextual analysis, allowing them to parse email text for subtle cues like domain spoofing, suspicious URLs, or grammatical inconsistencies that signature-based methods might miss [40].

Beyond text analysis, image recognition technologies add another layer of protection by identifying visual elements such as logos or brand imagery that may be replicated to mislead recipients [41]. These tools can detect common tactics like pixel manipulation or slight color variations used to circumvent basic phishing filters. AI-driven phishing detection systems can also learn from historical user interactions; for example, if a certain address has frequently sent benign attachments in the past but suddenly exhibits unusual spikes in outbound phishing messages, the system can issue an

early alert. However, adversaries continuously refine their strategies to evade detection, underscoring the need for regular retraining of these models. Integrating user education with AI-based tools ensures a more comprehensive defense, as employees learn to spot and report suspicious communications even if automated systems fail to catch every threat [42, 43].

## 4. Challenges and Ethical Considerations

### 4.1. Adversarial AI and Model Poisoning

The integration of AI into cybersecurity tools has undoubtedly bolstered defense capabilities, but it has also opened avenues for attackers to weaponize AI techniques [44]. Adversarial AI occurs when malicious actors introduce deceptive inputs such as subtly altered images, network traffic patterns, or system logs to mislead or "confuse" machine learning models. These manipulations may be imperceptible to human observers yet significantly degrade a model's detection

performance [45]. In scenarios where an Intrusion Detection System (IDS) relies on neural networks to flag unusual traffic, even minor perturbations in packet metadata could cause the system to misclassify malicious activity as benign.

A closely related threat, model poisoning, targets the data pipelines that feed supervised or semi-supervised models. Attackers with access to the training process can manipulate labeled examples, thus skewing classification boundaries [46, 47]. This could allow new malware variants to blend in with normal traffic or files, evading detection altogether. Over time, a poisoned model systematically misclassifies attacks, defeating the very defenses it was designed to enhance. To counteract these threats, organizations must adopt rigorous data handling practices, including secure data pipelines that validate inputs before they reach a training environment. Adversarial training where models are exposed to potentially deceptive inputs during development can also harden defenses but implementing this requires additional computational resources and domain expertise. Moreover, verifying model integrity through checksums or secure enclaves ensures that unauthorized modifications do not go undetected, thereby bolstering the reliability of AI-driven security solutions [48, 49].

## 4.2. Bias and Explainability in AI for Cybersecurity

While adversarial AI presents external threats, bias and lack of explainability represent internal challenges that can undermine the efficacy and trustworthiness of AI-driven defenses. Bias arises when models train on datasets that are not representative of real-world conditions, resulting in skewed or inaccurate outcomes. For example, a model trained predominantly on Windows-based malware samples may fail to detect newly emerging threats targeting Linux or IoT systems. Similarly, certain behaviors or user groups might be underrepresented in the training data, leading to disproportionate false positives for those demographics or under detection of attacks in less-sampled categories [50]. Continuous monitoring of model performance across different segments geographic regions, device types, or user roles can help identify and correct biases before they compromise security operations.

A second dimension of this challenge is explainability, or the ability to interpret how and why AI systems make specific decisions. In cybersecurity contexts, decisions like blocking critical network traffic or quarantining system processes can have high stakes, potentially impacting business continuity or user productivity. Explainable AI (XAI) techniques such as LIME (Local Interpretable Model-Agnostic Explanations) or SHAP (SHapley Additive exPlanations) aim to generate human-readable explanations for a model's predictions [51, 52]. By highlighting the features or rule sets that influenced a decision, XAI fosters transparency and accountability, vital for both regulatory compliance and organizational trust. Security analysts can review these explanations to confirm whether an alert is valid, helping to reduce the frequency of false positives and potentially mitigating disruptions to legitimate workflow. As AI becomes further entrenched in security operations, striking the right balance between automated decision-making and human oversight will hinge on investments in explainability, bias mitigation, and continuous model governance.

## 5. Case Studies and Industry Applications

### 5.1. Google's Chronicle Security AI

Chronicle, developed by Alphabet's cybersecurity subsidiary, leverages AI-driven analytics to process and correlate massive log datasets from enterprise networks. By ingesting historical logs alongside real-time threat intelligence, Chronicle surfaces hidden Indicators of Compromise (IoCs) that might otherwise remain obscured in large and diverse environments [53]. This cloud-native platform typically handles trillions of security events per second, enabling it to identify malicious patterns and behavioral anomalies at scale. Security teams benefit from contextualized alerts and visualized attack timelines, making it easier to pinpoint the root cause of suspicious activities.

Unlike traditional Security Information and Event Management (SIEM) solutions that rely heavily on rule-based correlation, Chronicle applies advanced machine learning techniques to adapt continuously to new threat vectors [54]. For instance, if attackers modify Tactics, Techniques, and Procedures (TTPs), Chronicle's models can dynamically update detection criteria without requiring manual rule creation [55]. Moreover, tight integration with Google's broader ecosystem encompassing VirusTotal intelligence and Google Cloud logs allows for rapid ingestion of global threat data, further enhancing the platform's capacity to detect zero-day exploits. As organizations move toward hybrid and multi-cloud infrastructures, Chronicle's ability to scale horizontally while maintaining low latency becomes a critical advantage for large enterprises [56].

### 5.2. Microsoft's AI-Driven Threat Intelligence Platform

Microsoft's approach to AI-enhanced cybersecurity extends across its suite of products, including Microsoft Defender, Sentinel, and the Azure cloud services [57]. By analyzing billions of signals daily originating from Windows endpoints,

Office 365 accounts, Azure subscriptions, and third-party connectors Microsoft's machine learning models identify and correlate events that could indicate malicious intrusions. The system automatically triages alerts, aggregating related events into cohesive incidents to minimize alert fatigue for security analysts.

A notable feature of Microsoft's platform is automated threat hunting, which proactively seeks out compromised accounts, suspicious file transfers, or anomalous network connections. For example, if a legitimate user credential is used in an unusual location or at an atypical time, the system raises a flag for deeper investigation. This continuous data-driven scrutiny helps reduce the mean time to detect (MTTD) and mean time to respond (MTTR), key metrics in cybersecurity incident management. Additionally, the platform's AI capabilities are integrated with Microsoft's broader ecosystem of tools such as Active Directory or Intune facilitating rapid remediation or containment actions, like automatically revoking compromised credentials or isolating infected hosts [58].

## 5.3. Darktrace's AI-Driven Network Security

Darktrace's Enterprise Immune System employs an unsupervised machine learning model that is conceptually inspired by biological immune systems. Instead of relying on predefined signatures or rules, Darktrace's algorithms autonomously learn what constitutes "normal" behavior for each specific network gauging typical connection patterns, data flows, and user activities over time. By establishing a dynamic baseline, the platform identifies deviations that might indicate malicious behavior. This can range from an insider threat exfiltrating sensitive files to an external attacker attempting lateral movement across the network [59].

A hallmark of Darktrace's approach is its autonomous response mechanism, often referred to as "digital antibodies." Once a threat is flagged, the system can take immediate, proportionate action like throttling suspicious data transfers or temporarily suspending access to a compromised account while allowing legitimate operations to continue without disruption [60, 61]. This capability significantly reduces the time window in which attackers can inflict damage, ultimately minimizing the scope and cost of potential breaches. Though powerful, such autonomous responses require careful tuning to avoid false positives that could disrupt critical processes, making ongoing collaboration between security teams and Darktrace's machine learning models essential.

## 5.4. AI in Endpoint Protection: CrowdStrike and Cylance

Traditional endpoint protection often relies on signature-based detection, requiring frequent updates to maintain defenses against evolving malware. AI-driven endpoint security solutions like CrowdStrike Falcon and CylancePROTECT shift the paradigm by analyzing behavioral signatures rather than static file hashes [62, 63]. By pre-executing or emulating code in a secure environment often referred to as a "sandbox" these platforms assess whether the code exhibits characteristics consistent with malicious software. This approach helps detect polymorphic or fileless malware that disguises its signature to evade conventional antivirus scans.

Continuous learning is a key element in these endpoint solutions. As new threats emerge, AI models update to recognize novel exploit techniques, rendering them more resilient against zero-day vulnerabilities. Additionally, cloud-based threat intelligence informs the entire customer base once a new threat is detected in one environment [64]. Both CrowdStrike and Cylance tout real-time analysis capabilities that can automatically quarantine files or isolate endpoints, mitigating damage before it spreads laterally across a network. However, because these platforms rely on complex algorithms, maintaining explainability and transparency around automated decisions can be a challenge especially in scenarios where essential business processes may be abruptly halted to prevent infection [65].

# 6. Future Research Directions

## 6.1. AI-Augmented SOCs (Security Operations Centers)

As threats grow in frequency and complexity, Security Operations Centers (SOCs) have embraced AI solutions to streamline incident detection and response. Modern SOC platforms generate thousands of daily alerts, many of which require manual triage [66]. Integrating AI-driven analytics, automated correlation of events, and adaptive response mechanisms can reduce this alert deluge, enabling security analysts to concentrate on strategic tasks like threat hunting, deep-dive forensics, and adversary simulation (red teaming). Looking ahead, organizations may develop fully AI-augmented SOCs that operate with minimal human intervention, relying on unsupervised or reinforcement learning algorithms to detect novel intrusions and initiate containment protocols automatically. While these systems promise rapid response times, questions remain about how to manage false positives or complex attack chains that demand nuanced human judgment. Striking a balance between autonomous detection and human oversight will be critical,

particularly in regulated industries such as finance or healthcare where unintended shutdowns of essential services can have serious repercussions [67].

## 6.2. AI for Proactive Cyber Threat Intelligence

Contemporary AI-based security solutions often excel at identifying known threats, but advanced adversaries frequently develop bespoke attack vectors. Proactive AI-driven threat intelligence goes beyond real-time detection by predicting emerging trends or uncovering advanced persistent threat (APT) campaigns in their early stages. This predictive capability stems from combining open-source intelligence (OSINT) including social media sentiment, vulnerability disclosures, and dark web chatter with advanced ML or deep learning algorithms to spot early indicators of malicious coordination [68]. For instance, an AI model might detect anomalous patterns in newly registered domain names or repeated chatter about specific software exploits, prompting security teams to patch relevant systems preemptively. Overcoming challenges in data quality, language translation, and swift data ingestion will be key to fully realizing these proactive capabilities. Additionally, privacy and ethical considerations must be addressed when gathering large-scale data from public or semi-private channels, ensuring compliance with regulations like GDPR and preserving users' civil liberties [69].

## 6.3. Quantum AI for Next-Generation Encryption

The advent of quantum computing introduces both an existential threat to current cryptographic algorithms like RSA and ECC and a potential accelerator for computational tasks, including machine learning. Should quantum computing become readily available to malicious actors, existing public-key encryption standards might be rendered obsolete by quantum-based decryption techniques [70]. Consequently, quantum-safe encryption algorithms are under active development, and AI can facilitate this transition by automating key generation, managing encryption lifecycle updates, and dynamically selecting suitable protocols as quantum capabilities evolve. Simultaneously, research into Quantum AI the fusion of quantum computing and ML holds promise for dramatically faster threat analysis and real-time defense adaptations. Yet, this fusion also raises new vulnerabilities, such as the risk of adversaries employing quantum AI to refine evasive attack strategies. Achieving next-generation encryption and computational security will require interdisciplinary collaboration, spanning computer science, cryptography, quantum physics, and cybersecurity policy [71].

## 6.4. Integration of AI with Blockchain for Enhanced Cybersecurity

Blockchain technologies, known for their decentralized and tamper-evident ledgers, have the potential to complement AI-based threat detection. By storing cryptographic hashes, verified threat intelligence, or critical event logs on a blockchain, organizations can maintain an immutable audit trail that attackers cannot easily alter or erase [72, 73]. This transparency is especially valuable when multiple stakeholders such as incident response teams, government agencies, and private sector partners must collaborate on complex threat investigations. AI systems could query the blockchain to confirm the integrity of shared data, detect inconsistencies in real time, or coordinate a response across a consortium of organizations. However, integrating AI with blockchain also entails challenges related to scalability, privacy, and the overhead of consensus mechanisms. As blockchain-based deployments expand, solutions will need to address storage constraints, avoid unnecessary duplication of large-scale datasets, and implement selective disclosure techniques to protect sensitive information. If properly orchestrated, the synergy between AI's advanced analytics and blockchain's trustless architecture could forge a more robust and collaborative defense against cyberattacks in an interconnected world [74, 75, 76].

## 7. Conclusion

Artificial Intelligence (AI) and Machine Learning (ML) have fundamentally redefined the landscape of cybersecurity, allowing defenders to stay ahead of increasingly sophisticated adversaries. Techniques such as supervised learning have proven effective for signature-based malware detection, while unsupervised anomaly detection has made it possible to spot zero-day exploits and stealthy intrusion attempts. Reinforcement learning adds another dimension, offering the possibility of autonomous defense strategies that adapt to ever-changing threat environments in real time. Meanwhile, AI-augmented Security Operations Centers (SOCs) herald an era of semi- or fully automated security workflows, freeing human analysts to focus on nuanced threat hunting and strategic oversight.

Despite these advancements, critical challenges persist. Explainable AI (XAI) is crucial for establishing trust in automated decision-making processes, especially when blocking traffic or quarantining systems can have real-world operational consequences. Adversarial AI tactics, including model poisoning and the creation of malicious inputs, demonstrate that attackers can exploit the very tools meant to defend against them. Bias in training data also continues

to pose a significant threat to the accuracy and fairness of AI-driven security solutions. Addressing these concerns will require rigorous research, collaborative information sharing, and the integration of robust governance frameworks.

Looking ahead, emerging technologies offer both opportunities and challenges. Quantum computing stands to outpace current cryptographic standards, prompting the development of quantum-safe algorithms and possibly ushering in a new wave of AI-driven cryptographic systems. Concurrently, blockchain integration could create tamper-proof repositories of threat intelligence, bolstering data integrity and fostering inter-organizational cooperation. These innovations highlight the continued need for multidisciplinary collaboration across industry, academia, and government entities. By uniting expertise in cryptography, data science, legal and regulatory frameworks, and ethical AI practices, the cybersecurity community can develop solutions that remain resilient in the face of relentless threat evolution.

Ultimately, the future of AI-driven cybersecurity hinges on striking a balance between technological sophistication and ethical responsibility. Sustaining public trust, protecting critical infrastructure, and safeguarding sensitive data will depend on how effectively stakeholders coordinate to refine algorithms, share threat intelligence, and establish transparent standards for accountability and privacy. As digital transformation accelerates across all sectors, AI's role in cybersecurity will undoubtedly expand, making continuous innovation, rigorous oversight, and collaborative partnership indispensable for defending against the threats of tomorrow.

## Compliance with ethical standards

*Disclosure of conflict of interest*

No conflict of interest to be disclosed.

## References

[1] LaRocque A, Gross G, Lindholm F, Greco P, Dupont B, Kruger J. Effective ransomware detection using autonomous patternbased signature extraction.

[2] Mbah GO, Evelyn AN. AI-powered cybersecurity: Strategic approaches to mitigate risk and safeguard data privacy.

[3] Hassan A, Mhmood AH. Optimizing network performance, automation, and intelligent decision-making through real-time big data analytics. International Journal of Responsible Artificial Intelligence. 2021 Aug 8;11(8):12-22.

[4] Kountardas N. Big data real-time security analytics.

[5] Obioha-Val OA, Lawal TI, Olaniyi OO, Gbadebo MO, Olisa AO. Investigating the Feasibility and Risks of Leveraging Artificial Intelligence and Open Source Intelligence to Manage Predictive Cyber Threat Models. Journal of Engineering Research and Reports. 2025 Jan 23;27(2):10-28.

[6] Malik J, Muthalagu R, Pawar PM. A systematic review of adversarial machine learning attacks, defensive controls and technologies. IEEE Access. 2024 Jul 4.

[7] Sagar R, Jhaveri R, Borrego C. Applications in security and evasions in machine learning: a survey. Electronics. 2020 Jan 3;9(1):97.

[8] Vegesna VV. Comprehensive analysis of AI-enhanced defense systems in cyberspace. International Numeric Journal of Machine Learning and Robots. 2023 Dec 21;7(7).

[9] Mbona I, Eloff JH. Classifying social media bots as malicious or benign using semi-supervised machine learning. Journal of Cybersecurity. 2023 Jan 1;9(1):tyac015.

[10] Mbona, I. and Eloff, J.H., 2023. Classifying social media bots as malicious or benign using semi-supervised machine learning. Journal of Cybersecurity, 9(1), p.tyac015.

[11] Das S, Nayak SP, Sahoo B, Nayak SC. Machine Learning in Healthcare Analytics: A State-of-the-Art Review. Archives of Computational Methods in Engineering. 2024 Apr 4:1-40.

[12] He K, Kim DD, Asghar MR. Adversarial machine learning for network intrusion detection systems: A comprehensive survey. IEEE Communications Surveys & Tutorials. 2023 Jan 3;25(1):538-66.

[13] El Husseini F, Noura H, Salman O, Chehab A. Advanced Machine Learning Approaches for Zero-Day Attack Detection: A Review. In2024 8th Cyber Security in Networking Conference (CSNet) 2024 Dec 4 (pp. 297-304). IEEE.

[14] El Husseini, Fatema, Hassan Noura, Ola Salman, and Ali Chehab. "Advanced Machine Learning Approaches for Zero-Day Attack Detection: A Review." In 2024 8th Cyber Security in Networking Conference (CSNet), pp. 297-304. IEEE, 2024.

[15] Sharma N, Arora B, Ziyad S, Singh PK, Singh Y. A Holistic review and performance evaluation of unsupervised learning methods for network anomaly detection. International Journal on Smart Sensing and Intelligent Systems.;17(1).

[16] Sharma, N., Arora, B., Ziyad, S., Singh, P. K., & Singh, Y. A Holistic review and performance evaluation of unsupervised learning methods for network anomaly detection. International Journal on Smart Sensing and Intelligent Systems, 17(1).

[17] Alahmadi BA, Axon L, Martinovic I. 99% false positives: A qualitative study of {SOC} analysts' perspectives on security alarms. In31st USENIX Security Symposium (USENIX Security 22) 2022 (pp. 2783-2800).

[18] Sharief F, Ijaz H, Shojafar M, Naeem MA. Multi-Class Imbalanced Data Handling with Concept Drift in Fog Computing: A Taxonomy, Review, and Future Directions. ACM Computing Surveys. 2024 Oct 7;57(1):1-48.

[19] Wang C, Zhu H. Wrongdoing monitor: A graph-based behavioral anomaly detection in cyber security. IEEE Transactions on Information Forensics and Security. 2022 Jul 15;17:2703-18.

[20] Madhu AS, Rapolu S. Anomaly Detection in Wait Reports and its Relation with Apache Cassandra Statistics.

[21] Madhu, Abheyraj Singh, and Sreemayi Rapolu. "Anomaly Detection in Wait Reports and its Relation with Apache Cassandra Statistics." (2021).

[22] Kelani AM. Feature Based Transfer Learning Intrusion Detection System (Doctoral dissertation, University of Guelph).

[23] Liu Z. A review of advancements and applications of pre-trained language models in cybersecurity. In2024 12th International Symposium on Digital Forensics and Security (ISDFS) 2024 Apr 29 (pp. 1-10). IEEE.

[24] Τσίγγανος N. Utilizing deep learning and natural language processing to recognise chat-based social engineering attacks for cyber security situational awareness.

[25] Shuvo MM, Islam SK, Cheng J, Morshed BI. Efficient acceleration of deep learning inference on resource-constrained edge devices: A review. Proceedings of the IEEE. 2022 Dec 14;111(1):42-91.

[26] Zhou DW, Wang QW, Qi ZH, Ye HJ, Zhan DC, Liu Z. Class-incremental learning: A survey. IEEE Transactions on Pattern Analysis and Machine Intelligence. 2024 Jul 16.

[27] Javadpour A, Ja'fari F, Taleb T, Shojafar M, Benzaïd C. A comprehensive survey on cyber deception techniques to improve honeypot performance. Computers & Security. 2024 Mar 1:103792.

[28] Javadpour, A., Ja'fari, F., Taleb, T., Shojafar, M. and Benzaïd, C., 2024. A comprehensive survey on cyber deception techniques to improve honeypot performance. Computers & Security, p.103792.

[29] Geng J, Wang J, Fang Z, Zhou Y, Wu D, Ge W. A Survey of strategy-driven evasion methods for PE malware: transformation, concealment, and attack. Computers & Security. 2024 Feb 1;137:103595.

[30] Neumann PG. Computer-related risks. Addison-Wesley Professional; 1994 Oct 18.

[31] Yeşiltepe D. A heuristic solution approach for dynamic mission abort problem based on deep reinforcement learning (Master's thesis, Middle East Technical University).

[32] Hassan, A., Nizam-Uddin, N., Quddus, A., Hassan, S.R., Rehman, A.U. and Bharany, S., 2024. Navigating IoT Security: Insights into Architecture, Key Security Features, Attacks, Current Challenges and AI-Driven Solutions Shaping the Future of Connectivity. Computers, Materials & Continua, 81(3).

[33] Loukkaanhuhta M. Transforming technical IT security architecture to a cloud era.

[34] Mattila R. Data pipeline monitoring solution and data quality in manufacturing company.

[35] Herath, H.M.S.S., Herath, H.M.K.K.M.B., Madhusanka, B.G.D.A. and Guruge, L.G.P.K., 2024. Data protection challenges in the processing of sensitive data. In Data Protection: The Wake of AI and Machine Learning (pp. 155-179). Cham: Springer Nature Switzerland.

[36] Vela E. Anomaly Detection in IoT Devices using LogBERT (Doctoral dissertation, Concordia University).

[37] Jain S. Advancing cybersecurity with artificial intelligence and machine learning: Architectures, algorithms, and future directions in threat detection and mitigation.

[38] Gupta N, Jindal V, Bedi P. A Survey on Intrusion Detection and Prevention Systems. SN Computer Science. 2023 Jun 10;4(5):439.

[39] Ayeni RK, Adebiyi AA, Okesola JO, Igbekele E. Phishing attacks and detection techniques: A systematic review. In2024 International Conference on Science, Engineering and Business for Driving Sustainable Development Goals (SEB4SDG) 2024 Apr 2 (pp. 1-17). IEEE.

[40] Mwaruwa MC. Long Short Term Memory Based Detection Of Web Based Sql Injection Attacks (Doctoral dissertation, UoN).

[41] Simonson A, Schmitt BH. Marketing aesthetics: The strategic management of brands, identity, and image. Simon and Schuster; 1997 Aug 30.

[42] Sarker IH. AI-driven cybersecurity and threat intelligence: cyber automation, intelligent decision-making and explainability. Springer Nature; 2024.

[43] Fakhouri HN, Alhadidi B, Omar K, Makhadmeh SN, Hamad F, Halalsheh NZ. AI-Driven Solutions for Social Engineering Attacks: Detection, Prevention, and Response. In2024 2nd International Conference on Cyber Resilience (ICCR) 2024 Feb 26 (pp. 1-8). IEEE.

[44] Heng L. Strategic Overview of Applying Artificial Intelligence on the Future Battlefield.

[45] Elsayed G, Shankar S, Cheung B, Papernot N, Kurakin A, Goodfellow I, Sohl-Dickstein J. Adversarial examples that fool both computer vision and time-limited humans. Advances in neural information processing systems. 2018;31.

[46] Miller DJ, Xiang Z, Kesidis G. Adversarial learning targeting deep neural network classification: A comprehensive review of defenses against attacks. Proceedings of the IEEE. 2020 Feb 26;108(3):402-33.

[47] Xiang Z, Miller DJ, Kesidis G. Detection of backdoors in trained classifiers without access to the training set. IEEE Transactions on Neural Networks and Learning Systems. 2020 Dec 16;33(3):1177-91.

[48] Srinivas MB, Konguvel E. Era of Sentinel Tech: Charting Hardware Security Landscapes through Post-Silicon Innovation, Threat Mitigation and Future Trajectories. IEEE Access. 2024 May 13.

[49] Srinivas, M.B. and Konguvel, E., 2024. Era of Sentinel Tech: Charting Hardware Security Landscapes through Post-Silicon Innovation, Threat Mitigation and Future Trajectories. IEEE Access.

[50] Fang GH, Lin ZM, Xie CZ, Han QZ, Hong MY, Zhao XY. Optimized Machine Learning Model for Predicting Compressive Strength of Alkali-Activated Concrete Through Multi-Faceted Comparative Analysis. Materials. 2024 Oct 18;17(20):5086.

[51] Rane N, Choudhary S, Rane J. Explainable Artificial Intelligence (XAI) in healthcare: Interpretable Models for Clinical Decision Support. Available at SSRN 4637897. 2023 Nov 15.

[52] Rane, N., Choudhary, S., & Rane, J. (2023). Explainable Artificial Intelligence (XAI) in healthcare: Interpretable Models for Clinical Decision Support. Available at SSRN 4637897.

[53] Bhardwaj A. Insecure digital frontiers: Navigating the global cybersecurity landscape. CRC Press; 2024 Oct 30.

[54] Farooq U. Cyber-physical security: AI methods for malware/cyber-attacks detection on embedded/IoT applications (Doctoral dissertation, Politecnico di Torino).

[55] Mennuni M. An Analysis of SOC Monitoring Systems (Doctoral dissertation, Politecnico di Torino).

[56] Ali SA. Designing Secure and Robust E-Commerce Plaform for Public Cloud. The Asian Bulletin of Big Data Management. 2023 Nov 25;3(1):164-89.

[57] Morić Z, Dakić V, Kapulica A, Regvart D. Forensic Investigation Capabilities of Microsoft Azure: A Comprehensive Analysis and Its Significance in Advancing Cloud Cyber Forensics. Electronics. 2024 Nov 19;13(22):4546.

[58] Loukasmäki H. Cyber Incident Response in Public Cloud: implications of modern cloud computing characteristics for cyber incident response.

[59] Schlicher BG, MacIntyre LP, Abercrombie RK. Towards reducing the data exfiltration surface for the insider threat. In2016 49th Hawaii International Conference on System Sciences (HICSS) 2016 Jan 5 (pp. 2749-2758). IEEE.

[60] Redwood WO. Apecs: A dynamic framework for preventing and mitigating theft, loss, and leakage of mission critical information in trust management networks.

[61] Kaul D. Blockchain-Powered Cyber-Resilient Microservices: AI-Driven Intrusion Prevention with Zero-Trust Policy Enforcement.

[62] Afianian A, Niksefat S, Sadeghiyan B, Baptiste D. Malware dynamic analysis evasion techniques: A survey. ACM Computing Surveys (CSUR). 2019 Nov 14;52(6):1-28.

[63] George AS. Riding the AI Waves: An Analysis of Artificial Intelligence's Evolving Role in Combating Cyber Threats. Partners Universal International Innovation Journal. 2024 Feb 25;2(1):39-50.

[64] Chadwick DW, Fan W, Costantino G, De Lemos R, Di Cerbo F, Herwono I, Manea M, Mori P, Sajjad A, Wang XS. A cloud-edge based data security architecture for sharing and analysing cyber threat information. Future generation computer systems. 2020 Jan 1;102:710-22.

[65] Dwivedi YK, Hughes DL, Coombs C, Constantiou I, Duan Y, Edwards JS, Gupta B, Lal B, Misra S, Prashant P, Raman R. Impact of COVID-19 pandemic on information management research and practice: Transforming education, work and life. International journal of information management. 2020 Dec 1;55:102211.

[66] Basta A, Basta N, Anwar W, Essar MI. Open-Source Security Operations Center (SOC): A Complete Guide to Establishing, Managing, and Maintaining a Modern SOC. John Wiley & Sons; 2024 Sep 23.

[67] Kabadayi S, O'Connor GE, Tuzovic S. The impact of coronavirus on service ecosystems as service mega-disruptions. Journal of Services Marketing. 2020 Nov 7;34(6):809-17.

[68] Alzaabi FR, Mehmood A. A review of recent advances, challenges, and opportunities in malicious insider threat detection using machine learning methods. IEEE Access. 2024 Feb 26;12:30907-27.

[69] Shehu VP, Shehu V. Human rights in the technology era–Protection of data rights. European Journal of Economics, Law and Social Sciences. 2023;7(2):1-0.

[70] Lior A. A Quantum of Privacy. Nevada Law Journal. 2024 Mar 19;25.

[71] Bilal H. The Impact of Quantum Computing on Cybersecurity: Challenges and Opportunities.

[72] Duncan B, Whittington M. Creating and Configuring an Immutable Database for Secure Cloud Audit Trail and System Logging. International Journal On Advances in Security. 2017;10(3&4):155-66.

[73] Oluwabanke A.S., Oyindamola M.O., Selina A.O., Adeniyi P.P. & Akeem O.O. Transforming corporate finance and advisory services with machine learning applications in risk management. GSC Advanced Research and Reviews, 2025, 22(02), 094-103

[74] Chatziamanetoglou D, Rantos K. Cyber Threat Intelligence on Blockchain: A Systematic Literature Review. Computers. 2024 Feb 26;13(3):60.

[75] Rafique W, Qadir J. Internet of everything meets the metaverse: Bridging physical and virtual worlds with blockchain. Computer Science Review. 2024 Nov 1;54:100678.

[76] Oyindamola M.O., Oluwabanke A.S., Adeniyi P.P. & Selina A.O. Unlocking new opportunities for strategic advisory and innovation with digital twin technology in corporate finance. World Journal of Advanced Research and Reviews, 2025, 25(02), 733-744.