Scientific Research Publishing

# Navigating AI Cybersecurity: Evolving Landscape and Challenges

**Maryam Roshanaei, Mahir R. Khan, Natalie N. Sylvester**

IST Department, Pennsylvania State University Abington, Abington, USA
Email: mur45@psu.edu

## Abstract

The rapid integration of artificial intelligence (AI) into critical sectors has revealed a complex landscape of cybersecurity challenges that are unique to these advanced technologies. AI systems, with their extensive data dependencies and algorithmic complexities, are susceptible to a broad spectrum of cyber threats that can undermine their functionality and compromise their integrity. This paper provides a detailed analysis of these threats, which include data poisoning, adversarial attacks, and systemic vulnerabilities that arise from the AI's operational and infrastructural frameworks. This paper critically examines the effectiveness of existing defensive mechanisms, such as adversarial training and threat modeling, that aim to fortify AI systems against such vulnerabilities. In response to the limitations of current approaches, this paper explores a comprehensive framework for the design and implementation of robust AI systems. This framework emphasizes the development of dynamic, adaptive security measures that can evolve in response to new and emerging cyber threats, thereby enhancing the resilience of AI systems. Furthermore, the paper addresses the ethical dimensions of AI cybersecurity, highlighting the need for strategies that not only protect systems but also preserve user privacy and ensure fairness across all operations. In addition to current strategies and ethical concerns, this paper explores future directions in AI cybersecurity.

## Keywords

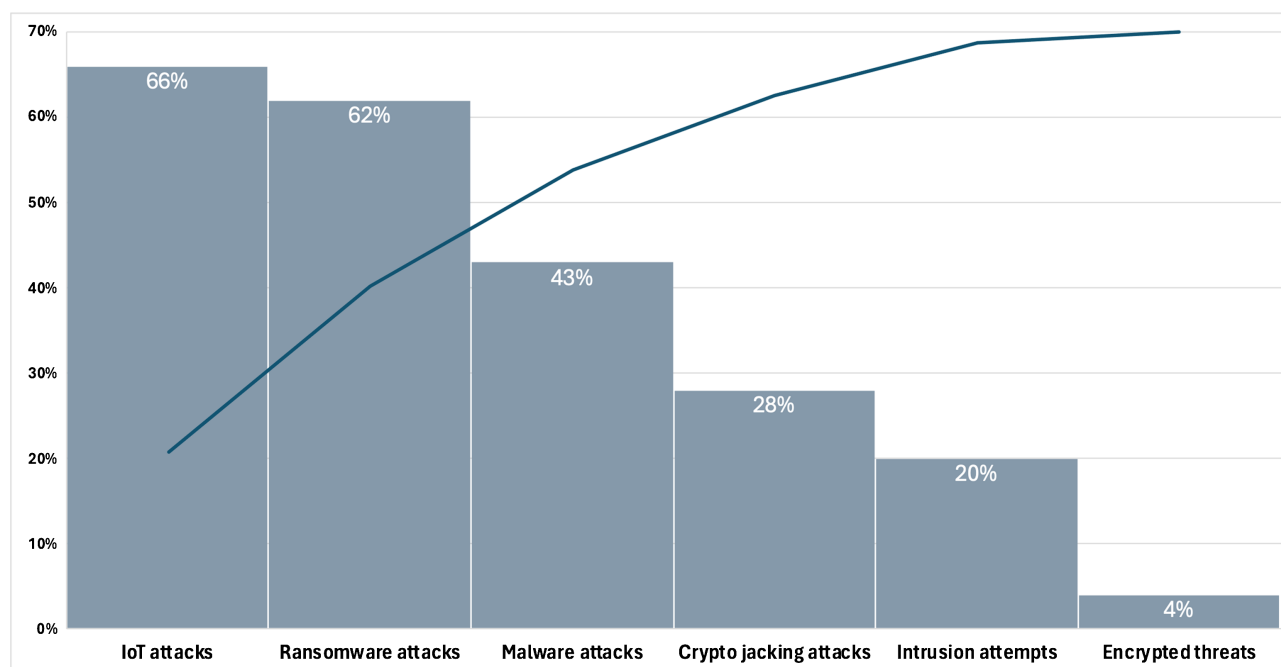AI Cybersecurity, Adversarial Attacks, Defensive Strategies, Ethical AI

## 1. Introduction

Artificial Intelligence (AI) [1] is transforming numerous industries by automating complex processes, enhancing data analytics, and creating new technological

capabilities. Its applications range from autonomous vehicles navigating city streets, to sophisticated trading algorithms in financial markets, to diagnostic tools in healthcare that predict patient outcomes. This widespread adoption of AI underscores its potential to significantly benefit society. However, as AI systems become more integral to critical infrastructures, the urgency to address their security vulnerabilities increases dramatically. AI systems [2] are not merely passive targets of cyber threats common to traditional IT environments; they are also susceptible to specialized attacks that exploit the unique characteristics of machine learning (ML) models and data-driven decision-making processes. For instance, adversarial attacks can subtly alter input data in ways that are imperceptible to humans but cause AI to make erroneous decisions, potentially leading to catastrophic outcomes. Similarly, data poisoning attacks can corrupt the training data of an AI model, leading to flawed learning and compromised operational performance. Furthermore, the complexity and opacity of many AI models [3]—often referred to as "black boxes"—pose significant challenges for cybersecurity. This lack of transparency can obscure vulnerabilities, making it difficult to detect when and how an AI system has been compromised. This issue is exacerbated in systems that continually evolve based on new data, as ongoing learning can introduce new vulnerabilities or amplify existing ones. Therefore, a robust approach to AI cybersecurity must address not only traditional cyber threats but also those that are uniquely enabled or enhanced by AI technologies. It involves ensuring the integrity and security of the data used by AI, protecting the AI models themselves, and securing the infrastructures that support them. The objective is not only to defend against attacks but also to design AI systems that maintain their intended functionality and continue to operate safely, even when under threat. Addressing these challenges requires a collaborative effort among cybersecurity experts, AI developers, industry stakeholders, and policymakers. Together, they must develop and implement comprehensive cybersecurity strategies [4] that integrate advanced security technologies, adhere to rigorous standards, and evolve in response to the dynamic nature of threats in the AI landscape. This proactive and comprehensive approach is crucial for maintaining the integrity, trustworthiness, and reliability of AI applications, thereby ensuring that AI continues to drive positive outcomes across all sectors of society. AI systems face unique challenges that traditional cybersecurity measures are not fully equipped to handle [5]. Table 1 provides a detailed overview of the challenges faced by AI Systems in Cybersecurity.

Figure 1 illustrates the historical progression of cybersecurity threats and responses over the past few decades [6]. It collectively demonstrates the escalating nature of AI-specific threats and the critical importance of innovative and robust defense strategies. As AI technologies continue to permeate critical sectors, the collaboration between cybersecurity experts, developers, and policymakers becomes increasingly vital. They must work together to implement these strategies, ensuring that AI systems can resist and recover from cyber-attacks, thereby maintaining their functionality and reliability in a challenging security landscape.

**Table 1.** Challenges faced by AI systems in cybersecurity.

| Attack Type | Description |
|---|---|
| Adversarial Attacks | Subtle alterations to input data that deceive AI models. |
| Data Poisoning | Deliberate corruption of the training dataset, leading to flawed learning outcomes. |
| Model Theft | Unauthorized access and duplication of AI models. |
| Infrastructure Attacks | Targeting the physical and virtual environments supporting AI systems. |



**Figure 1.** Increase in AI-Specific Cybersecurity Incidents (2010-2023).

**Table 2.** Defensive strategies for AI systems and their effectiveness.

| Strategy | Description | Application Example | Effectiveness (%) |
|---|---|---|---|
| Robust AI Training | Training models with adversarial examples to detect attacks | Used in facial recognition systems | 75% |
| Data Integrity Measures | Techniques to ensure data is not altered or corrupted | Applied in healthcare AI for patient data | 85% |
| Enhanced Model Security | Encryption and rigorous access controls | Financial AI systems for fraud detection | 90% |
| Infrastructure Security | Advanced protocols to secure physical and virtual components | Autonomous vehicles | 80% |

Effective defense mechanisms are essential to protect AI systems from these evolving threats. Table 2 outlines these strategies and their effectiveness of Defensive Strategies for AI Systems [7].

## 2. AI Vulnerabilities and Attack Vectors

As AI technologies continue to evolve and integrate into various industries, they

become prime targets for sophisticated cyber threats. These threats [8] exploit unique vulnerabilities inherent to AI systems, such as their reliance on data integrity, the transparency of their algorithms, and the security of their supporting infrastructure. Understanding the nature of these vulnerabilities and the corresponding attack vectors is crucial for developing robust countermeasures that can protect AI systems from potential cyber-attacks. Data poisoning [9] represents a critical threat to AI systems, particularly because these systems rely heavily on data integrity for training and operation. In a data poisoning attack, adversaries deliberately manipulate the training data to compromise the model's learning process, leading to flawed decision-making or predictive abilities. This type of attack is particularly dangerous because it can be difficult to detect and can have far-reaching effects once a model is deployed. Techniques include Injection Attacks and Modification Attacks. Injection Attacks [10] involve inserting malicious data points into the training dataset. The AI system, unaware of the tampering, learns from this corrupted data, which can lead to significant deviations in its behavior. For instance, an AI model used for financial forecasting could be taught incorrect associations, leading to erroneous investment recommendations. Modification Attacks, [11] on the other hand, alter existing data within the dataset rather than adding new data points. Even minor changes to critical data points can retrain the model with false information, resulting in incorrect outputs. Such attacks might be used to manipulate systems like automated surveillance, where altering image data could prevent the recognition of specific individuals or objects. Concrete examples of these techniques' real-world impacts include attackers compromising a facial recognition system by introducing subtly altered images into its training set. These alterations, imperceptible to humans, were significant enough to fool the system, resulting in the failure to identify or misidentification of individuals. This vulnerability was exploited to manipulate facial recognition, leading to incorrect tagging, or ignoring of faces, which could potentially bypass security protocols. Another instance involved a traffic control AI in an urban smart city system. Attackers injected faulty data representing fake traffic conditions, such as non-existent traffic jams or accidents. The AI, trained with these false data points, generated incorrect traffic flow predictions, causing chaos in city traffic management and emergency response services.

Defensive measures against attacks on AI systems [12] include robust data validation, data provenance, and redundancy in data sources. Implementing stringent data validation processes is crucial for verifying the integrity of training data before use. This can involve statistical analyses to detect outliers or inconsistencies that may indicate tampering. Maintaining a secure and transparent chain of custody for training data, known as data provenance, helps in identifying potential points of compromise by tracking the source and history of the data. Additionally, using multiple, independent data sources to train AI models can reduce the risk of poisoning. Discrepancies between sources [13] can be flagged

and investigated, ensuring greater reliability and security in the training process. Data poisoning not only highlights the vulnerabilities in AI systems but also underscores the need for comprehensive security practices in AI training environments. By understanding and implementing rigorous defensive strategies, organizations can better protect their AI applications from these insidious attacks. Table 3 summarizes the types, techniques, case studies, and defensive measures in data poisoning attacks.

## 3. Model Stealing and Inversion

Model stealing and model inversion [14] represent significant threats in the domain of AI cybersecurity, targeting the core intellectual properties and sensitive data integral to AI systems. These methods are particularly insidious because they can undermine the competitive advantage of technology companies and violate privacy laws. Model stealing [15], or model extraction attacks, occur when attackers aim to replicate an AI system's model without direct access to the underlying architecture or training data. Typically, this is done by observing the outputs of a model in response to various inputs, and then using this information to train a new model that exhibits similar behavior. Techniques used to compromise AI models include query attacks and side-channel attacks. Query attacks involve attackers feeding inputs into the model and gathering the outputs, which they then use to train a duplicate model. Side-channel attacks, on the other hand, involve inferring model properties from indirect information such as computational time or power consumption. These attacks pose significant risks and have real-world impacts. Economic damage is a major concern, as intellectual property theft through model stealing can lead to substantial financial losses for the original developers. This is because they lose the competitive edge gained through their investment in research and development. Additionally, there are regulatory compliance issues to consider. In sectors where AI models are part of regulated activities, unauthorized duplication can lead to compliance violations, affecting the legality of the copied model's deployment. Model inversion attacks [16] aim to extract sensitive or proprietary information

**Table 3.** Overview of data poisoning in AI systems.

| Category | Technique | Description | Case Study | Defensive Measure |
|---|---|---|---|---|
| Injection Attacks | Insertion of malicious data | Deliberately adding harmful data to the training set, leading to erroneous model training. | Financial Forecasting Misguidance | Robust Data Validation, Anomaly Detection |
| Modification Attacks | Alteration of existing data | Subtly changing critical data points within the dataset, thereby corrupting the model's output. | Manipulation of Facial Recognition | Data Provenance, Statistical Data Integrity Checks |

about the dataset used for training from a trained model. These attacks exploit the model's ability to recall parts of its training data, potentially exposing personal data even if the model is only intended to make generalized inferences. Techniques used to compromise AI models include reconstruction and inference. Reconstruction involves attackers carefully crafting input queries and analyzing the model's responses to infer characteristics of the training data. Inference, on the other hand, uses statistical techniques to deduce information about the training data's distribution or specific data entries. These techniques pose significant risks and have real-world impacts. One major concern is privacy violations [17]. If the training data includes personal information, model inversion can lead to significant privacy breaches, potentially exposing sensitive personal data such as medical records or financial information. Another critical impact is the loss of public trust. Organizations suffering from such breaches risk losing the trust of their customers and partners, as well as facing potential legal ramifications. Defensive measures against AI model attacks include differential privacy, limited model access, and homomorphic encryption. Integrating differential privacy during the training phase helps protect against both model stealing and model inversion by adding noise to the model's outputs, making it difficult for attackers to gain useful information from querying the model. Limiting the number of queries an individual can make and monitoring for unusual access patterns can also reduce the risk of these attacks. Additionally, using homomorphic encryption [18] allows AI models to operate on encrypted data, preventing attackers from accessing usable data directly from model queries. These measures collectively enhance the security and privacy of AI systems.

## 4. Evasion Attacks

Evasion attacks [19] represent a significant security threat to AI systems, particularly in environments where decisions are made based on input data that could be manipulated by an adversary. These attacks involve deliberately crafting input data that an AI model misinterprets, leading to incorrect decisions or actions. This type of cyberattack exploits specific weaknesses in the model's design or its learning process, typically by using knowledge about the model's algorithms and sensitivities.

Adversaries can bypass AI systems [20] using adversarial examples, which are carefully modified inputs that appear normal to human observers but deceive AI systems into making errors. This manipulation involves adding small, often imperceptible perturbations to input data that cause the AI to misclassify it. For example [21], in image recognition, altering pixel values in an image of a stop sign might lead an autonomous vehicle's AI to recognize it as a yield sign instead. Attackers often use gradient-based techniques to determine the most effective alterations to the input. By understanding how slight changes to input affect the outputs, attackers can calculate the minimum perturbation needed to fool the model. An intriguing aspect of adversarial examples is their transferabil-

ity; an adversarial example crafted to deceive one model often works against another model, even if they have different architectures or were trained on different data sets. This makes defending against evasion attacks particularly challenging. The misinterpretation of traffic signs by autonomous vehicles can lead to incorrect driving decisions, posing safety risks to both the vehicle's occupants and other road users. In security applications, such as facial recognition or anomaly detection, evasion attacks can allow unauthorized individuals to bypass biometric security measures or evade detection by surveillance systems [22].

Defensive measures [23] include adversarial training, which involves including adversarial examples during the training phase of the AI model. By learning from these perturbed inputs, the model becomes more robust against similar attacks in the future. Implementing input sanitization steps to detect and mitigate suspicious alterations in input data before it is fed into the model can help reduce the risk of evasion attacks. Regularly updating AI models to recognize new adversarial tactics and patching discovered vulnerabilities is crucial for maintaining security.

## 5. Infrastructure Attacks

Evasion attacks on AI systems [24] are a type of cyber threat that exploits specific vulnerabilities in how these systems process inputs. The primary objective of evasion attacks is to cause the AI to make incorrect decisions or take inappropriate actions by subtly manipulating the input data. These attacks are particularly challenging because they often involve minimal changes that are difficult to detect with standard security measures, yet they can significantly alter the behavior of an AI system.

Gradient-based techniques [25] leverage the gradient of the AI model's loss function (which measures the error of the model's predictions) to determine how to modify inputs in a way that maximizes error. By calculating the gradient with respect to the input data, attackers can identify the most effective changes to the input that will lead the model to an incorrect output. This approach is commonly used in creating adversarial examples for image recognition systems, where slight adjustments to the pixel values of an image can cause the model to mislabel the image. The transferability of attacks implies that an adversarial example designed for one model can often deceive another model, even if the two models have different architectures or were trained on different datasets. This characteristic makes evasion attacks particularly menacing as it allows for scalable and effective attacks across different systems. An adversarial example crafted [26] to fool a model used in one type of autonomous car could potentially be effective against another brand's model, raising significant concerns for industry-wide standards and defense strategies. Real-world impacts of evasion attacks are significant. In the context of autonomous vehicles [27], incorrect interpretations of road signs or unexpected obstacles due to manipulated input data can lead to accidents, endangering human lives and property. In security settings,

such as facial recognition systems used at airports or in smartphones, evasion attacks [28] can allow unauthorized access, compromising personal security and organizational safety. Defensive measures against evasion attacks include adversarial training, input sanitization, and regular model updates and patching. Adversarial training involves training the AI model on a mixture of regular and adversarial examples. The goal is to make the model less sensitive to the small perturbations typically used in evasion attacks, thereby improving its ability to generalize from manipulated inputs. While this method enhances robustness, it can also reduce the model's accuracy on non-adversarial examples, creating a trade-off between security and performance. Input sanitization involves applying rigorous checks and transformations to the inputs before they are processed by the AI model. Techniques such as noise reduction, anomaly detection, or even reformatting [29] can help mitigate the impact of subtle input manipulations. Sanitization can prevent many types of evasion attacks, especially when combined with other monitoring techniques that flag unusual input patterns. Regular model updates and patching involve continuously updating the AI models to detect new patterns of adversarial examples and patching any identified vulnerabilities in the model's design or its data processing pipeline. Regular updates ensure that the AI system evolves in response to emerging threats, maintaining its defensive capabilities against novel attack strategies. Table 4 summarizes the key aspects of evasion attack mechanisms, their impacts, and the defensive measures [30] that can be employed to mitigate these risks.

## 6. Defensive Mechanisms and Frameworks

Ensuring the security of AI systems requires a multi-layered approach that encompasses the integrity of the data, the robustness of the models, and the

**Table 4.** Summary of evasion attack mechanisms, impacts, and defensive measures.

| Aspect | Description |
| --- | --- |
| Gradient-Based Techniques | Leverage model's gradient to modify inputs for maximizing error in outputs. |
| Transferability of Attacks | Adversarial examples for one model often work against different models. |
| Autonomous Vehicles | Manipulated data leads to incorrect driving decisions, posing safety risks. |
| Security System Breaches | Allow unauthorized access, compromising personal and organizational security. |
| Adversarial Training | Training on both regular and adversarial examples to improve model robustness. |
| Input Sanitization | Rigorous checks and transformations to detect and mitigate suspicious inputs. |
| Regular Model Updates | Continuous updates to recognize new adversarial tactics and patch vulnerabilities. |

adaptability of the security measures. Securing the training data for AI systems [31] is paramount, as the quality and integrity of this data directly impact the performance and reliability of the AI models. Effective strategies for ensuring data integrity include both preventive measures to protect data from being compromised and reactive measures to identify and mitigate any integrity issues that occur.

Data provenance [32] involves maintaining a comprehensive, traceable record of the data's origin, where it has been, who has handled it, and how it has been modified over time. This can be achieved through the integration of provenance-tracking technologies that automatically log every interaction with the data. By having detailed logs of data handling and modifications, organizations can quickly trace the source of any corrupted data, understand the nature of the corruption, and determine the appropriate remedial actions. This traceability is crucial for maintaining not only the integrity of the data but also for regulatory compliance, especially in industries governed by strict data protection regulations. Implementations might include blockchain for immutable data logs, metadata management systems, and automated auditing tools. Anomaly detection in the context of AI training involves continuously monitoring the data for any signs of tampering or abnormalities. This is typically done using sophisticated algorithms that can learn the normal patterns of the data and detect deviations that may indicate attempts to manipulate the model's learning process. These can range from outright corrupted data entries to subtle statistical outliers that could skew the model's understanding and predictions. The system can be set to alert administrators or take automated corrective actions when anomalies are detected. Statistical methods, machine learning models such as neural networks or clustering algorithms, and rule-based systems [33] are commonly employed. Advanced implementations might use unsupervised learning to adapt to new data without needing pre-labeled examples of anomalies. Effective anomaly detection [34] can preemptively prevent data poisoning attacks by identifying and isolating corrupt data before it is used in training. This not only protects the AI from learning incorrect information but also minimizes the risk of deploying flawed AI models. The integration of data provenance and anomaly detection systems provides a robust framework for data integrity. Provenance systems provide a clear audit trail, which is invaluable when investigating and resolving issues identified by anomaly detection systems. Conversely, anomaly detection can initially flag issues that, upon further investigation using provenance data, can be quickly resolved.

AI systems face several challenges and considerations [35] that must be addressed to ensure their effectiveness and reliability. One significant challenge is scalability. As data volumes grow, maintaining extensive provenance records and performing real-time anomaly detection can become computationally expensive, requiring substantial resources and sophisticated infrastructure. Another critical issue is managing false positives and negatives. Balancing the sensitivity of anomaly detection systems to minimize false positives (benign anomalies

flagged as threats) and false negatives (genuine threats missed) is crucial for maintaining the accuracy and reliability of the system. Additionally, data privacy is a paramount concern. Implementing these systems must not compromise the privacy of the data, especially when personal or sensitive information is involved, ensuring that robust privacy measures are integrated into the design and operation of AI systems. By focusing on these strategies [36], organizations can significantly enhance the security and reliability of their AI systems. Ensuring data integrity is not just about protecting data from external threats, but also about maintaining the trustworthiness and dependability of the AI systems that rely on this data. Table 5 summarizes these strategies for securing AI training data.

## 7. Secure AI Architectures

Dynamic security [37] measures are essential in the rapidly evolving field of cybersecurity, particularly with AI-driven systems that face constantly changing threat landscapes. These adaptive systems not only need to recognize threats but must also evolve to handle new types of attacks. Implementing advanced AI-driven threat prediction involves using predictive analytics and real-time monitoring to identify potential security threats before they materialize. Predictive analytics leverages machine learning algorithms to analyze historical data and identify patterns that could indicate future attacks. By training models on vast datasets that include instances of security breaches, the system learns to recognize precursors to an attack. Real-time monitoring continuously scans system operations to detect anomalies that deviate from established patterns, which is critical in environments like network security where traffic patterns can be monitored for signs of intrusion or data exfiltration. The benefits of AI-driven threat prediction [38] are substantial. Proactive security measures allow organizations to shift from a reactive to a proactive security posture, potentially stopping attacks before they start. Additionally, predictive systems enable more

**Table 5.** Strategies for securing AI training data.

| Strategy | Implementation Details | Benefits | Tools and Technologies |
|---|---|---|---|
| Data Provenance | Maintaining comprehensive records of data origin, handling, and modifications. | Quick traceability of corrupted data, understanding corruption nature, regulatory compliance. | Blockchain, metadata management systems, automated auditing tools |
| Anomaly Detection | Continuously monitoring data for tampering or abnormalities using sophisticated algorithms. | Preemptively prevents data poisoning, protecting AI from learning incorrect information, minimizing risk of deploying flawed AI models. | Statistical methods, machine learning models (neural networks, clustering algorithms), rule-based systems, unsupervised learning |

efficient allocation of security resources by focusing efforts where an attack is deemed most likely. Continuous learning systems, incorporating adaptive learning technologies, are essential for maintaining up-to-date security measures. Online learning models adapt to new data as it comes in without needing to retrain from scratch, which is crucial for keeping the system current with the latest threats. Reinforcement learning allows the system to learn optimal actions in different scenarios through trial and error, using feedback from its own actions to improve over time.

The benefits of these systems include adaptability, as they enable the system to adjust to new strategies employed by attackers that might not have been in the original training data. Moreover, they ensure long-term relevance by continually integrating new insights and data, maintaining the effectiveness of security measures. Automated response mechanisms enhance the speed and consistency of security responses. Automated patching applies [40] software patches automatically when vulnerabilities are detected, swiftly closing security gaps. Smart isolation can automatically isolate affected nodes or network segments in the event of a detected breach, preventing the spread of the attack. Dynamic access control adjusts user or system access privileges in real-time based on detected threats, restricting access when abnormal behavior is identified until further investigation can occur. The primary benefits of automated response mechanisms [41] are speed and consistency. Automated responses are typically faster than those orchestrated by humans, which is crucial for mitigating fast-moving cyber threats. Consistency ensures that every incident is responded to uniformly, reducing the likelihood of errors or oversights. While dynamic security measures offer significant advantages, they also come with challenges. Implementing sophisticated AI-driven security systems can be complex and costly, requiring significant upfront investment in technology and expertise. Overly sensitive systems might flag normal activities as suspicious, leading to disruptions in user operations and potential dissatisfaction. Additionally, the use of AI in monitoring and response mechanisms must be balanced with ethical considerations and compliance with privacy laws and regulations. By addressing these challenges [42], organizations can harness the full potential of dynamic security measures to protect against a wide array of cyber threats, ensuring their AI systems are robust and adaptable to the ever-changing digital threat landscape. Table 6 summarizes these aspects of AI-driven threat prediction, continuous learning systems, and automated response mechanisms.

## 8. Anomaly Detection

Dynamic security measures are pivotal in adapting [43] to and countering the rapidly evolving threats in cybersecurity, particularly when dealing with AI systems. These measures, characterized by their adaptability and automation, play crucial roles in not just detecting threats but also responding to them in real time. AI-driven threat prediction systems [44] analyze vast amounts of data

**Table 6.** AI-driven threat prediction, continuous learning systems, and automated response mechanisms.

| Aspect | Advanced Implementation | Benefits |
|---|---|---|
| AI-Driven Threat Prediction | Predictive Analytics: Utilizing machine learning to identify potential future attacks. Real-Time Monitoring: Continuously scanning for anomalies. | Proactivity: Shifts to a proactive security posture, stopping attacks before they start. Resource Optimization: Efficient allocation of security resources. |
| Continuous Learning Systems | Online Learning Models: Adapt to new data without retraining. Reinforcement Learning: Learns optimal actions through trial and error. | Adaptability: Adjusts to new attacker strategies. Long-Term Relevance: Continually integrates new insights and data. |
| Automated Response Mechanisms | Automated Patching: Automatically applies software patches. Smart Isolation: Automatically isolates affected nodes or segments. Dynamic Access Control: Adjusts access privileges based on threat levels. | Speed: Faster responses to cyber threats. Consistency: Uniform response to incidents. |

from past security incidents and current system interactions to identify patterns that might indicate potential threats. This data could include log files, network traffic data, or even patterns of user behavior. Using techniques [45] like machine learning, these systems develop models that can predict where and how the next cyber-attack might occur. This modeling often employs complex algorithms capable of handling large-scale data and delivering predictions with high accuracy. Potential threats can be scored based on their likelihood and potential impact, allowing security teams to prioritize their responses or preventive measures effectively. Additionally, AI-driven threat prediction [46] is often integrated with anomaly detection systems to refine the accuracy of alerts and reduce false positives, ensuring that security resources are focused where they are most needed.

Continuous learning systems [47] are built on technology foundations such as online machine learning, which continuously ingests and learns from new data as it is generated, adapting their models in real time. This is crucial in environments where threat vectors evolve rapidly, as it allows the AI to stay current without manual retraining. Feedback from the results of previous security actions (such as the success or failure of a mitigation strategy) is used to refine the AI models, improving their predictive capabilities and decision-making processes over time. The AI's ability to adapt [48] based on new data helps it anticipate shifts in cyberattack strategies, potentially identifying threats before they become active issues. Continuous learning contributes to the resilience of security systems by enabling them to evolve with the threats, rather than requiring periodic overhauls or updates. Automated response mechanisms offer functional capabil-

ities such as incident response automation, where automated scripts or AI-driven systems can take immediate action in response to detected threats, such as isolating affected systems, shutting down compromised operations, or deploying patches. In addition to responding to detected threats, automated systems [49] can also take proactive measures based on predicted risks, such as adjusting firewall rules or changing access permissions dynamically. However, the complexity of deploying these systems requires sophisticated programming and deep integration with existing IT infrastructure, which can be complex and costly. It's crucial to balance automation with oversight to prevent disruptive responses to false alarms, which can affect system reliability and user trust. Table 7 summarizes [50] these aspects of AI-driven threat prediction, continuous learning systems, and automated response mechanisms.

## 9. Ethical Considerations in AI Cybersecurity

Ethical considerations [51] are central to the deployment and operation of AI systems, particularly as these systems are integrated into critical and sensitive areas of our lives. Ethical challenges in AI cybersecurity not only influence public trust but also impact the effectiveness of the AI systems in their intended functions. In the context of AI cybersecurity, addressing bias and ensuring fairness are critical challenges. Biases can originate from various sources, including historical inequalities captured in data, subjective human decisions during data

**Table 7.** AI-driven threat prediction, continuous learning systems, and automated response mechanisms.

| Aspect | Advanced Implementation | Benefits |
| --- | --- | --- |
| AI-Driven Threat Prediction | Data Utilization: Analyzing data from past incidents and current interactions.<br>Predictive Modeling: Using machine learning to predict future attacks.<br>Threat Scoring: Prioritizing responses based on likelihood and impact.<br>Anomaly Detection Integration: Refining alerts accuracy and reducing false positives. | Proactivity: Shifts to a proactive security posture, stopping attacks before they start.<br>Resource Optimization: Efficient allocation of security resources. |
| Continuous Learning Systems | Online Machine Learning: Adapts models in real time with new data.<br>Feedback Mechanisms: Refines models based on results of previous actions. | Adaptability: Adjusts to new attacker strategies.<br>Long-Term Relevance: Continually integrates new insights and data. |
| Automated Response Mechanisms | Incident Response Automation: Takes immediate action in response to detected threats.<br>Proactive Measures: Takes preventive actions based on predicted risks.<br>Complexity in Deployment: Requires sophisticated programming and deep integration.<br>Balancing Act: Preventing disruptive responses to false alarms. | Speed: Faster responses to cyber threats.<br>Consistency: Uniform response to incidents. |

collection, and selective or incomplete data gathering processes. These biases [52] can lead to differential security outcomes and increased susceptibility to manipulation, necessitating robust strategies to mitigate bias and promote fairness.

In the context of AI cybersecurity, addressing bias and ensuring fairness are critical challenges. Biases can originate from various sources, including historical inequalities captured in data, subjective human decisions during data collection, and selective or incomplete data-gathering processes. These biases can lead to differential security outcomes and increased susceptibility to manipulation, necessitating robust strategies to mitigate bias and promote fairness. Inherent biases in data stem from historical inequalities, subjective human decisions, and selective data gathering. For instance [53], an AI system trained predominantly on security data from urban areas may not perform well in rural settings, potentially leading to inadequate security measures in less represented areas. These biases can become self-reinforcing. A predictive policing system targeting certain communities due to biased historical arrest data may generate more data from those same communities, further entrenching the bias. Biases in AI-driven cybersecurity tools can lead to differential security outcomes. Anomaly detection systems might flag legitimate behaviors in certain demographic groups [54] as suspicious, resulting in higher rates of false positives. Additionally, biased AI systems might be more susceptible to cyber-attacks. Attackers could exploit known biases to craft inputs that are more likely to be misclassified, bypassing security measures. To enhance strategies for mitigating bias, diverse training datasets are crucial. Techniques such as synthetic data generation can simulate underrepresented groups or scenarios, enhancing dataset diversity without compromising privacy. Continuous evaluation systems [55] are also necessary to ensure that as new data is incorporated or as the model evolves, it remains free of emerging biases. Bias audits are essential, with the frequency and scope of audits being crucial. High-risk AI applications may require frequent audits that cover training data, model outputs, algorithms, and deployment environments. Conducting these audits by independent third parties can maintain objectivity and boost public trust in the audit process.

Transparent AI models are another key strategy. Implementing explainable AI (XAI) techniques [56], such as feature importance scoring and decision trees, can help illuminate the decision-making process of AI models, aiding users, and regulators in understanding how decisions are made. Engaging a broad spectrum of stakeholders in the design and review process of AI models can ensure diverse perspectives, helping to identify and mitigate potential biases that may not be apparent to the original design team. By understanding the sources and impacts of bias in AI systems, and by implementing strategies such as diverse training datasets, bias audits, and transparent AI models, organizations can enhance the fairness and effectiveness of AI in cybersecurity. This holistic approach addresses [57] the challenges, impacts, and mitigation strategies for bias in AI cybersecurity, as detailed in Table 8.

**Table 8.** Challenges, impacts, and mitigation strategies for bias in AI cybersecurity.

| Challenge | Description | Impact on Cybersecurity | Mitigation Strategies |
|---|---|---|---|
| Sources of Bias | Biases from historical inequalities, human decisions, and data gathering | Inadequate security measures in underrepresented areas | Diverse training datasets, synthetic data generation |
| Compounding Effects | Self-reinforcing biases due to feedback loops | Entrenched bias in AI systems | Continuous evaluation, bias audits |
| Differential Security Outcomes | Different protection levels for different groups | Higher rates of false positives for certain demographics | Transparent AI models, explainable AI |
| Risk of Manipulation | Exploitation of biases by attackers | Increased susceptibility to cyber-attacks | Stakeholder engagement, independent review |

## 10. Conclusion and Future Work

AI introduces unique vulnerabilities distinct from traditional cybersecurity challenges—such as data poisoning, model theft, and adversarial attacks—which significantly complicate the security landscape. In response, a set of robust defense mechanisms tailored to these advanced threats, including adversarial training, enhanced encryption, and sophisticated anomaly detection systems, has been examined. These defenses represent the lead of a dynamic and continuously evolving toolkit designed to not only address current security challenges but also anticipate future threats. Moreover, the development of AI-specific solutions like self-healing systems and decentralized threat detection models showcases the potential of AI to revolutionize cybersecurity practices by enabling more autonomous and adaptive security measures. The ethical considerations surrounding AI cybersecurity are equally complex. Issues such as bias, privacy, and accountability in AI systems are of paramount concern, as these technologies have the potential to significantly impact society. Ethical AI deployment in cybersecurity is crucial for maintaining public trust and ensuring that these technologies are used responsibly. This involves integrating comprehensive ethical guidelines into every phase of AI development and deployment, from initial design to real-world implementation.

Addressing these challenges requires a multidisciplinary approach. Collaboration across different fields—cybersecurity experts, AI developers, ethicists, and policymakers—is essential to develop solutions that are not only effective but also ethically sound and socially responsible. Moreover, as AI technologies rapidly evolve, so too must the policies and regulations that govern their use. This may involve regular updates to laws and regulations, proactive international cooperation to standardize AI security practices, and the establishment of global norms and protocols that address both the opportunities and challenges presented by AI in cybersecurity. Navigating the regulatory landscape involves sev-

eral key elements. Adaptive regulatory frameworks are essential to staying relevant in the face of technological advancements and cybersecurity threats. These frameworks should include mechanisms for continuous review and rapid updates, possibly through dedicated regulatory bodies or task forces. Looking ahead, the advancement of AI in cybersecurity hinges on several critical areas of future work that are essential for addressing the dynamic and increasingly sophisticated landscape of cyber threats. A primary focus is the development of advanced predictive models. Future research should aim to enhance the predictive capabilities of AI systems, enabling them to foresee potential cyber threats before they materialize. This involves not only refining the accuracy of predictions but also the speed at which these AI systems can detect and respond to threats, ensuring they can operate ahead of potential breaches. Another significant area of future work is the refinement of AI ethics and privacy protections. As AI technologies become more embedded in cybersecurity efforts, continuous efforts are needed to ensure that the ethical frameworks governing their use evolve in line with new technological developments and societal expectations. This includes addressing concerns related to data privacy, the transparency of AI decision-making processes, and the equitable treatment of all individuals affected by AI-driven actions.

## Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

## References

[1] Mallick, M.A.I. and Nath, R. (2024) Navigating the Cyber Security Landscape: A Comprehensive Review of Cyber-Attacks, Emerging Trends, and Recent Developments. *World Scientific News*, **190**, 1-69.

[2] Aldoseri, A., Al-Khalifa, K.N. and Hamouda, A.M. (2023) Re-Thinking Data Strategy and Integration for Artificial Intelligence: Concepts, Opportunities, and Challenges. *Applied Sciences*, **13**, Article 7082. https://doi.org/10.3390/app13127082

[3] Goni, A., Jahangir, M.U.F. and Chowdhury, R.R. (2024) A Study on Cyber Security: Analyzing Current Threats, Navigating Complexities, and Implementing Prevention Strategies. *International Journal of Research and Scientific Innovation*, **10**, 507-522. https://doi.org/10.51244/ijrsi.2023.1012039

[4] Thakur, M. (2024) Cyber Security Threats and Countermeasures in Digital Age. *Journal of Applied Science and Education* (*JASE*), **4**, 1-20.

[5] Camacho, N.G. (2024) The Role of AI in Cybersecurity: Addressing Threats in the Digital Age. *Journal of Artificial Intelligence General Science*, **3**, 143-154. https://doi.org/10.60087/jaigs.v3i1.75

[6] Mohamed, N. (2023) Current Trends in AI and ML for Cybersecurity: A State-of-the-Art Survey. *Cogent Engineering*, **10**, Article 2272358. https://doi.org/10.1080/23311916.2023.2272358

[7] Kaloudi, N. and Li, J. (2020) The AI-Based Cyber Threat Landscape: A Survey. *ACM Computing Surveys*, **53**, Article No. 20. https://doi.org/10.1145/3372823

[8] Guembe, B., Azeta, A., Misra, S., Osamor, V.C., Fernandez-Sanz, L. and Pospelova, V. (2022) The Emerging Threat of AI-Driven Cyber Attacks: A Review. *Applied Artificial Intelligence*, **36**, Article 2037254. https://doi.org/10.1080/08839514.2022.2037254

[9] Sun, G., Cong, Y., Dong, J., Wang, Q., Lyu, L. and Liu, J. (2022) Data Poisoning Attacks on Federated Machine Learning. *IEEE Internet of Things Journal*, **9**, 11365-11375. https://doi.org/10.1109/jiot.2021.3128646

[10] Tufail, S., Batool, S. and Sarwat, A.I. (2021) False Data Injection Impact Analysis in AI-Based Smart Grid. *SoutheastCon* 2021, Atlanta, 10-13 March 2021, 1-7. https://doi.org/10.1109/southeastcon45413.2021.9401940

[11] De Mello, F.L. (2020) A Survey on Machine Learning Adversarial Attacks. *Journal of Information Security and Cryptography* (*Enigma*), **7**, 1-7. https://doi.org/10.17648/jisc.v7i1.76

[12] Sadeghi, K., Banerjee, A. and Gupta, S.K.S. (2020) A System-Driven Taxonomy of Attacks and Defenses in Adversarial Machine Learning. *IEEE Transactions on Emerging Topics in Computational Intelligence*, **4**, 450-467. https://doi.org/10.1109/tetci.2020.2968933

[13] Ramirez, M.A., Kim, S.K., Hamadi, H.A., Damiani, E., Byon, Y.J., Kim, T.Y., Yeun, C.Y., *et al.* (2022) Poisoning Attacks and Defenses on Artificial Intelligence: A Survey. arXiv: 2202.10276. https://doi.org/10.48550/arXiv.2202.10276

[14] Zhou, S., Zhu, T., Ye, D., Zhou, W. and Zhao, W. (2024) Inversion-Guided Defense: Detecting Model Stealing Attacks by Output Inverting. *IEEE Transactions on Information Forensics and Security*, **19**, 4130-4145. https://doi.org/10.1109/tifs.2024.3376190

[15] Chesterman, S. (2024) Good Models Borrow, Great Models Steal: Intellectual Property Rights and Generative AI. *Policy and Society*, puae006. https://doi.org/10.1093/polsoc/puae006

[16] Zhao, X., Zhang, W., Xiao, X. and Lim, B. (2021) Exploiting Explanations for Model Inversion Attacks. 2021 *IEEE/CVF International Conference on Computer Vision* (*ICCV*), Montreal, 10-17 October 2021, 662-672. https://doi.org/10.1109/iccv48922.2021.00072

[17] Thuraisingham, B.M. (2020) Can AI Be for Good in the Midst of Cyber Attacks and Privacy Violations? A Position Paper. *Proceedings of the Tenth ACM Conference on Data and Application Security and Privacy*, New Orleans, 16-18 March 2020, 1-4. https://doi.org/10.1145/3374664.3379334

[18] Su, G., Wang, J., Xu, X., Wang, Y. and Wang, C. (2024) The Utilization of Homomorphic Encryption Technology Grounded on Artificial Intelligence for Privacy Preservation. *International Journal of Computer Science and Information Technology*, **2**, 52-58. https://doi.org/10.62051/ijcsit.v2n1.07

[19] Roshanaei, M. (2024) Enhancing Mobile Security through Comprehensive Penetration Testing. *Journal of Information Security*, **15**, 63-86. https://doi.org/10.4236/jis.2024.152006

[20] Rawal, A., Rawat, D.B. and Sadler, B. (2021) Recent Advances in Adversarial Machine Learning: Status, Challenges and Perspectives. *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications III*, 117462Q. https://doi.org/10.1117/12.2583970

[21] Khayyam, H., Javadi, B., Jalili, M. and Jazar, R.N. (2020) Artificial Intelligence and Internet of Things for Autonomous Vehicles. In: Jazar, R. and Dai, L., Eds., *Nonlinear Approaches in Engineering Applications*, Springer, 39-68.

https://doi.org/10.1007/978-3-030-18963-1_2

[22] Sagar, R., Jhaveri, R. and Borrego, C. (2020) Applications in Security and Evasions in Machine Learning: A Survey. *Electronics*, **9**, Article 97. https://doi.org/10.3390/electronics9010097

[23] Gupta, R., Tanwar, S., Tyagi, S. and Kumar, N. (2020) Machine Learning Models for Secure Data Analytics: A Taxonomy and Threat Model. *Computer Communications*, **153**, 406-440. https://doi.org/10.1016/j.comcom.2020.02.008

[24] Roshanaei, M. and Duan, Q. (2021) International Telecommunication Union Standardization for Trust Provisioning in Information, Communication and Technology Infrastructure toward Achieving United Nation's Sustainable Development Goals. *Journal of Computer and Communications*, **9**, 44-59. https://doi.org/10.4236/jcc.2021.910004

[25] Wang, J., Tuyls, J., Wallace, E. and Singh, S. (2020) Gradient-Based Analysis of NLP Models Is Manipulable. *Findings of the Association for Computational Linguistics: EMNLP* 2020, 16-20 November 2020, 247-258. https://doi.org/10.18653/v1/2020.findings-emnlp.24

[26] Aldahdooh, A., Hamidouche, W., Fezza, S.A. and Déforges, O. (2022) Adversarial Example Detection for DNN Models: A Review and Experimental Comparison. *Artificial Intelligence Review*, **55**, 4403-4462. https://doi.org/10.1007/s10462-021-10125-w

[27] Bathla, G., Bhadane, K., Singh, R.K., Kumar, R., Aluvalu, R., Krishnamurthi, R., *et al.* (2022) Autonomous Vehicles and Intelligent Automation: Applications, Challenges, and Opportunities. *Mobile Information Systems*, **2022**, Article 7632892. https://doi.org/10.1155/2022/7632892

[28] Awad, A.I., Babu, A., Barka, E. and Shuaib, K. (2024) AI-Powered Biometrics for Internet of Things Security: A Review and Future Vision. *Journal of Information Security and Applications*, **82**, Article 103748. https://doi.org/10.1016/j.jisa.2024.103748

[29] Agrawal, S. (2022) Enhancing Payment Security through AI-Driven Anomaly Detection and Predictive Analytics. *International Journal of Sustainable Infrastructure for Cities and Societies*, **7**, 1-14.

[30] Hossain, M.T., Afrin, R. and Biswas, M.A.A. (2024) A Review on Attacks against Artificial Intelligence (AI) and Their Defence Image Recognition and Generation Machine Learning, Artificial Intelligence. *Control Systems and Optimization Letters*, **2**, 52-59.

[31] Habbal, A., Ali, M.K. and Abuzaraida, M.A. (2024) Artificial Intelligence Trust, Risk and Security Management (AI TRiSM): Frameworks, Applications, Challenges and Future Research Directions. *Expert Systems with Applications*, **240**, Article 122442. https://doi.org/10.1016/j.eswa.2023.122442

[32] Kim, J.J.H., Um, R.S., Lee, J.W.Y. and Ajilore, O. (2024) Generative AI Can Fabricate Advanced Scientific Visualizations: Ethical Implications and Strategic Mitigation Framework. *AI and Ethics*. https://doi.org/10.1007/s43681-024-00439-0

[33] Barbierato, E. and Gatti, A. (2024) The Challenges of Machine Learning: A Critical Review. *Electronics*, **13**, Article 416. https://doi.org/10.3390/electronics13020416

[34] Trilles, S., Hammad, S.S. and Iskandaryan, D. (2024) Anomaly Detection Based on Artificial Intelligence of Things: A Systematic Literature Mapping. *Internet of Things*, **25**, Article 101063. https://doi.org/10.1016/j.iot.2024.101063

[35] Adelani, F.A., Okafor, E.S., Jacks, B.S. and Ajala, O.A. (2024) Theoretical Frameworks for the Role of AI and Machine Learning in Water Cybersecurity: Insights

from African and U.S. Applications. *Computer Science & IT Research Journal*, **5**, 681-692. https://doi.org/10.51594/csitrj.v5i3.928

[36] Cinà, A.E., Grosse, K., Demontis, A., Biggio, B., Roli, F. and Pelillo, M. (2024) Machine Learning Security against Data Poisoning: Are We There Yet? *Computer*, **57**, 26-34. https://doi.org/10.1109/mc.2023.3299572

[37] Abdi, A.H., Audah, L., Salh, A., Alhartomi, M.A., Rasheed, H., Ahmed, S. and Tahir, A. (2024) Security Control and Data Planes of SDN: A Comprehensive Review of Traditional, AI and MTD Approaches to Security Solutions. *IEEE Access*, **12**, 69941-69980. https://doi.org/10.1109/ACCESS.2024.3393548

[38] Sontan, A.D. and Samuel, S.V. (2024) The Intersection of Artificial Intelligence and Cybersecurity: Challenges and Opportunities. *World Journal of Advanced Research and Reviews*, **21**, 1720-1736. https://doi.org/10.30574/wjarr.2024.21.2.0607

[39] Zhang, Z., Hamadi, H.A., Damiani, E., Yeun, C.Y. and Taher, F. (2022) Explainable Artificial Intelligence Applications in Cyber Security: State-of-the-Art in Research. *IEEE Access*, **10**, 93104-93139. https://doi.org/10.1109/access.2022.3204051

[40] Shwedeh, F., Malaka, S. and Rwashdeh, B. (2023) The Moderation Effect of Artificial Intelligent Hackers on the Relationship between Cyber Security Conducts and the Sustainability of Software Protection: A Comprehensive Review. *Migration Letters*, **20**, 1066-1072.

[41] Enholm, I.M., Papagiannidis, E., Mikalef, P. and Krogstie, J. (2021) Artificial Intelligence and Business Value: A Literature Review. *Information Systems Frontiers*, **24**, 1709-1734. https://doi.org/10.1007/s10796-021-10186-w

[42] Ahmad, K., Iqbal, W., El-Hassan, A., Qadir, J., Benhaddou, D., Ayyash, M., *et al.* (2024) Data-Driven Artificial Intelligence in Education: A Comprehensive Review. *IEEE Transactions on Learning Technologies*, **17**, 12-31. https://doi.org/10.1109/tlt.2023.3314610

[43] Sikder, M.N.K. and Batarseh, F.A. (2023) Outlier Detection Using AI: A Survey. In: Batarseh, F.A. and Freeman, L.J., Eds., *AI Assurance*, Academic Press, 231-291. https://doi.org/10.1016/B978-0-32-391919-7.00020-2

[44] Sarker, I.H. (2022) AI-Based Modeling: Techniques, Applications and Research Issues towards Automation, Intelligent and Smart Systems. *SN Computer Science*, **3**, Article No. 158. https://doi.org/10.1007/s42979-022-01043-x

[45] Lau, P.L., Nandy, M. and Chakraborty, S. (2023) Accelerating UN Sustainable Development Goals with AI-Driven Technologies: A Systematic Literature Review of Women's Healthcare. *Healthcare*, **11**, Article 401. https://doi.org/10.3390/healthcare11030401

[46] Naik, B., Mehta, A., Yagnik, H. and Shah, M. (2021) The Impacts of Artificial Intelligence Techniques in Augmentation of Cybersecurity: A Comprehensive Review. *Complex & Intelligent Systems*, **8**, 1763-1780. https://doi.org/10.1007/s40747-021-00494-8

[47] Wiafe, I., Koranteng, F.N., Obeng, E.N., Assyne, N., Wiafe, A. and Gulliver, S.R. (2020) Artificial Intelligence for Cybersecurity: A Systematic Mapping of Literature. *IEEE Access*, **8**, 146598-146612. https://doi.org/10.1109/access.2020.3013145

[48] Zhang, Z., Hamadi, H.A., Damiani, E., Yeun, C.Y. and Taher, F. (2022) Explainable Artificial Intelligence Applications in Cyber Security: State-of-the-Art in Research. *IEEE Access*, **10**, 93104-93139. https://doi.org/10.1109/access.2022.3204051

[49] Vegesna, V.V. (2023) Comprehensive Analysis of AI-Enhanced Defense Systems in Cyberspace. *International Numeric Journal of Machine Learning and Robots*, **7**.

https://injmr.com/index.php/fewfewf/article/view/21

[50] Sewak, M., Sahay, S.K. and Rathore, H. (2022) Deep Reinforcement Learning in the Advanced Cybersecurity Threat Detection and Protection. *Information Systems Frontiers*, **25**, 589-611. https://doi.org/10.1007/s10796-022-10333-x

[51] Ansari, M.F., Dash, B., Sharma, P. and Yathiraju, N. (2022) The Impact and Limitations of Artificial Intelligence in Cybersecurity: A Literature Review. *International Journal of Advanced Research in Computer and Communication Engineering*, **11**, 81-90. https://doi.org/10.17148/ijarcce.2022.11912

[52] Kapoor, P. (2023) Machine Learning for Cyber Threat Detection: Advancements and Challenges. *International Journal of Machine Learning for Sustainable Development*, **5**. https://ijsdcs.com/index.php/IJMLSD/article/view/420

[53] Sinha, A.R., Singla, K. and Victor, T.M.M. (2023) Artificial Intelligence and Machine Learning for Cybersecurity Applications and Challenges. In: Kumar, R. and Pattnaik, P.K., Eds., *Risk Detection and Cyber Security for the Success of Contemporary Computing*, IGI Global, 109-146.
https://doi.org/10.4018/978-1-6684-9317-5.ch007

[54] Taddeo, M., Jones, P., Abbas, R., Vogel, K. and Michael, K. (2023) Socio-Technical Ecosystem Considerations: An Emergent Research Agenda for AI in Cybersecurity. *IEEE Transactions on Technology and Society*, **4**, 112-118.
https://doi.org/10.1109/tts.2023.3278908

[55] Dhirani, L.L., Mukhtiar, N., Chowdhry, B.S. and Newe, T. (2023) Ethical Dilemmas and Privacy Issues in Emerging Technologies: A Review. *Sensors*, **23**, Article 1151.
https://doi.org/10.3390/s23031151

[56] Dwivedi, R., Dave, D., Naik, H., Singhal, S., Omer, R., Patel, P., *et al.* (2023) Explainable AI (XAI): Core Ideas, Techniques, and Solutions. *ACM Computing Surveys*, **55**, Article No. 194. https://doi.org/10.1145/3561048

[57] Taskeen, and Garai, S. (2024) Emerging Trends in Cybersecurity: A Holistic View on Current Threats, Assessing Solutions, and Pioneering New Frontiers. *Blockchain in Healthcare Today*, **7**, Article 302. https://doi.org/10.30953/bhty.v7.302