

Cybersecurity Threats And Defenses In The Generative AI Era

Ginne M James

Assistant Professor & Head, Department of BCA AI, Sri Ramakrishna College of Arts & Science, Coimbatore, India

Article information

Received: 4th March 2026

Volume: 1

Received in revised form: 27th March 2026

Issue: 4

Accepted: 29th March 2026DOI: <https://doi.org/10.5281/zenodo.19466114>Available online: 9th April 2026

Abstract

The rapid proliferation of generative artificial intelligence (GenAI) technologies, including large language models (LLMs) and generative adversarial networks (GANs), has fundamentally transformed the cybersecurity landscape by simultaneously empowering sophisticated attack vectors and enabling advanced defensive mechanisms. This paper presents a comprehensive analysis of emerging cybersecurity threats catalyzed by generative AI, encompassing AI-generated phishing campaigns, deepfake-based social engineering, automated malware generation, and adversarial exploitation techniques. Concurrently, it examines AI-enhanced defense strategies, including transformer-based threat detection, adversarial training for model robustness, and neural network-driven intrusion detection systems. Through a systematic taxonomy of AI-powered threats and a comparative evaluation of traditional versus AI-augmented defense mechanisms, this study quantifies the paradigm shift in attack sophistication and defense efficacy. Findings reveal that AI-enhanced defenses achieve detection rate improvements ranging from 22.4% to 33.4% over traditional methods across phishing, malware, intrusion, and deepfake detection domains. The paper further discusses regulatory frameworks, ethical considerations, and future research directions essential for maintaining cybersecurity resilience in an AI-driven threat environment.

Keywords: - Generative AI, Cybersecurity, Large Language Models, Adversarial Machine Learning, Deepfake Detection, AI-Powered Threats, Neural Network Defense, Phishing Detection

I. INTRODUCTION

The emergence of generative artificial intelligence (GenAI) represents a watershed moment in the evolution of cybersecurity threats and defenses. Since the foundational work on generative adversarial networks (GANs) by Goodfellow et al. [1], the capability of AI systems to generate realistic synthetic content has advanced at an unprecedented pace. Large language models (LLMs) such as GPT-4, Claude, and Gemini now demonstrate human-level proficiency in natural language generation, while diffusion models and GANs produce photorealistic images, audio, and video content that is increasingly indistinguishable from authentic media [2]. These technological advances carry profound implications for cybersecurity, as the same generative capabilities that drive beneficial innovation also lower the barrier to entry for sophisticated cyberattacks.

The dual-use nature of GenAI technologies has created an asymmetric threat landscape. On the offensive side, threat actors leverage LLMs to craft highly convincing phishing emails, generate polymorphic malware that evades signature-based detection, and produce deepfake content for social engineering campaigns [3]. The automation capabilities of GenAI enable attackers to scale their operations dramatically—what previously required specialized expertise and significant time investment can now be accomplished by adversaries with minimal technical knowledge [4]. Reports from major cybersecurity firms indicate that AI-powered phishing attacks increased by over 220% between 2021

and 2023, with AI-generated phishing emails achieving click-through rates 60% higher than manually crafted campaigns [5].

Conversely, the defensive application of AI technologies has yielded significant improvements in threat detection and response capabilities. Transformer-based models have demonstrated superior performance in identifying malicious content, while adversarial training techniques enhance model robustness against evasion attacks [6]. AI-driven security orchestration and automated response (SOAR) systems enable organizations to respond to threats at machine speed, significantly reducing mean time to detection (MTTD) and mean time to response (MTTR) [7]. The National Institute of Standards and Technology (NIST) AI Risk Management Framework (AI RMF) has established guidelines for the responsible deployment of AI in cybersecurity contexts [8].

This paper addresses the critical need for a systematic examination of GenAI's impact on cybersecurity by presenting:

- a comprehensive taxonomy of AI-powered cybersecurity threats,
- a comparative analysis of AI-enhanced versus traditional defense mechanisms, and
- an evidence-based assessment of detection efficacy improvements.

The remainder of this paper is organized as follows: Section II reviews related work. Section III presents the taxonomy of AI-powered threats. Section IV examines AI-enhanced defense mechanisms. Section V provides comparative analysis with supporting data. Section VI discusses implications, limitations, and future directions. Section VII concludes the paper.

II. RELATED WORK

The intersection of artificial intelligence and cybersecurity has attracted substantial research attention over the past decade. Biggio and Roli [9] provided a seminal survey on adversarial machine learning, establishing the foundational understanding of how machine learning models can be manipulated through carefully crafted adversarial inputs. Their work highlighted the vulnerability of classification-based security systems to evasion attacks, poisoning attacks, and exploratory attacks—a framework that remains highly relevant in the GenAI era.

Carlini and Wagner [10] demonstrated the fragility of neural networks through their seminal work on adversarial examples, showing that deep neural networks could be reliably fooled by imperceptible perturbations. Their C&W attack framework exposed fundamental weaknesses in defensive distillation and other proposed robustness measures, establishing a benchmark for evaluating adversarial robustness that has been widely adopted in cybersecurity research. Building on this foundation, Madry et al. [11] proposed projected gradient descent (PGD) based adversarial training as a principled defense methodology, demonstrating that training on adversarially perturbed examples significantly improves model robustness.

Mirsky et al. [12] conducted an extensive survey on deepfake creation and detection, cataloging the evolution of face-swapping technologies from early GAN-based approaches to sophisticated encoder-decoder architectures. Their taxonomy of deepfake detection methods—spanning visual artifacts, physiological signals, and temporal inconsistencies—provides a comprehensive framework for understanding the detection challenge. The survey highlighted the arms race dynamic between deepfake generators and detectors, a pattern that has only intensified with the advent of diffusion-based generation models.

More recently, Gupta et al. [13] examined the cybersecurity implications of ChatGPT and similar LLMs, demonstrating that these models can generate functional exploit code, craft convincing social engineering narratives, and automate reconnaissance activities. Hazell [14] specifically investigated LLM-generated spear phishing, finding that GPT-4-generated emails achieved comparable effectiveness to human-crafted campaigns at a fraction of the cost and time. In the defensive domain, Ferrag et al. [15] surveyed the application of LLMs to cybersecurity tasks, documenting their utility in vulnerability analysis, threat intelligence extraction, and incident response automation.

While these studies have individually addressed specific aspects of the AI-cybersecurity nexus, a comprehensive framework that simultaneously examines both the offensive and defensive dimensions of GenAI in cybersecurity remains lacking. This paper addresses this gap by providing an integrated analysis spanning threat taxonomy, defense mechanisms, and quantitative performance comparison.

III. AI-POWERED CYBERSECURITY THREAT TAXONOMY

The integration of generative AI into the cyber threat landscape has produced a new generation of attack vectors characterized by increased sophistication, scalability, and evasiveness. This section presents a systematic taxonomy of AI-powered threats organized by attack category, generative technique, and operational impact.

A. AI-Generated Phishing and Social Engineering

Large language models have fundamentally transformed phishing attack capabilities. Traditional phishing campaigns relied on template-based approaches with recognizable linguistic patterns that enabled rule-based detection systems to achieve reasonable accuracy. LLM-generated phishing emails, however, exhibit natural language fluency, contextual awareness, and personalization capabilities that render conventional detection approaches ineffective [14].

Studies indicate that GPT-4-generated spear phishing emails achieve click-through rates of 31.4%, compared to 18.2% for human-crafted equivalents, while reducing campaign preparation time by approximately 95% [5]. Furthermore, LLMs enable multilingual phishing at scale, eliminating the linguistic errors that previously served as detection heuristics in non-English campaigns.

B. Deepfake-Based Identity Fraud

The advancement of GAN architectures and diffusion models has made deepfake generation accessible to non-expert users through consumer-grade tools [12]. In the cybersecurity context, deepfakes are deployed for identity fraud in video-based authentication systems, CEO fraud schemes exploiting synthetic video or audio, and disinformation campaigns targeting organizational reputation. Voice cloning technologies, leveraging neural text-to-speech models, require as little as three seconds of reference audio to produce convincing synthetic speech, enabling voice-based social engineering attacks that bypass voice biometric systems [16]. The reported incidents of deepfake fraud increased from approximately 350 cases in 2021 to over 1,420 in 2023, as shown in Figure 1.

C. Automated Malware Generation

Generative AI has significantly lowered the barrier to malware development. LLMs can generate functional exploit code, create polymorphic malware variants that evade signature-based detection, and automate the obfuscation of malicious payloads [13]. Research by Pa Pa et al. [17] demonstrated that ChatGPT could generate functional malware samples across multiple programming languages when prompted through carefully designed jailbreak techniques. More concerning is the emergence of purpose-built malicious LLMs—such as WormGPT and FraudGPT—that operate without safety guardrails and are specifically optimized for generating cybercrime tools [18].

D. Adversarial Attacks on AI Systems

As organizations increasingly deploy AI-based security systems, adversarial attacks targeting these models represent a critical meta-threat. Adversarial examples—inputs crafted with imperceptible perturbations to induce misclassification—can systematically undermine AI-based intrusion detection systems, malware classifiers, and content filters [10]. Generative models, particularly GANs, enable the automated creation of adversarial examples that transfer across different model architectures, further complicating defensive deployment [9]. Table 1 presents a comprehensive taxonomy of AI-powered threats, categorizing each by the generative technique employed, primary attack vector, severity level, and detection difficulty.

Table 1. Taxonomy of AI-Powered Cybersecurity Threats

Threat Category	GenAI Technique	Attack Vector	Severity	Detection Difficulty
AI-Generated Phishing	LLM (GPT-4, LLaMA)	Email / SMS / Chat	Critical	High
Deepfake Identity Fraud	GAN / Diffusion Models	Video / Image Authentication	Critical	Very High
Voice Cloning	Neural TTS (VALL-E)	Phone / Voice Biometrics	High	Very High
Polymorphic Malware	LLM Code Generation	Endpoint / Network	Critical	High
Automated Exploitation	LLM + Reinforcement Learning	Software Vulnerabilities	High	Medium
Adversarial Examples	GAN / PGD Optimization	AI Security Models	High	High
Data Poisoning	Generative Data Augmentation	Training Pipelines	Medium	Very High
Disinformation Campaigns	LLM + Diffusion Models	Social Media / News	High	Medium

IV. AI-ENHANCED DEFENSE MECHANISMS

The deployment of generative AI in cybersecurity defense has produced significant improvements across multiple security domains. This section examines the principal AI-enhanced defense mechanisms and their technical underpinnings.

A. Transformer-Based Threat Detection

Pre-trained transformer models, fine-tuned on cybersecurity-specific corpora, have demonstrated substantial improvements in phishing detection, malware classification, and threat intelligence extraction. BERT-based phishing detectors achieve detection rates exceeding 94% by capturing semantic and contextual features that evade traditional keyword-based and heuristic approaches [19]. SecurityBERT, a domain-specific language model pre-trained on cybersecurity text, further improves performance on tasks such as vulnerability classification and attack pattern identification [20]. These models leverage the attention mechanism to identify subtle linguistic patterns characteristic of malicious content, including urgency indicators, authority impersonation cues, and suspicious URL embedding strategies.

B. AI-Driven Intrusion Detection Systems

Neural network-based intrusion detection systems (NIDS) represent a significant advancement over traditional signature-based and statistical anomaly detection methods. Deep learning architectures—including convolutional neural networks (CNNs), recurrent neural networks (RNNs), and transformer encoders—have been applied to network traffic analysis, achieving detection rates of 91.4% on benchmark datasets such as CICIDS-2017 and UNSW-NB15 [21]. Variational autoencoders (VAEs) and GANs are employed for anomaly detection by learning the distribution of normal network behavior and flagging deviations, offering particular advantages in detecting zero-day attacks that lack known signatures [22].

C. Adversarial Training and Model Robustness

Adversarial training has emerged as the most effective technique for enhancing the robustness of AI-based security models against evasion attacks. Following the framework established by Madry et al. [11], PGD-based adversarial training augments training data with adversarially perturbed samples, enabling models to maintain accuracy under adversarial conditions. Recent advances incorporate certified robustness guarantees through randomized smoothing [23], providing provable bounds on model predictions within defined perturbation regions. Ensemble adversarial training, which exposes models to adversarial examples generated by multiple attack algorithms, further improves generalization to unseen attack strategies.

D. Deepfake Detection Systems

AI-based deepfake detection systems employ multiple detection modalities to identify synthetic media. Frequency-domain analysis detects spectral artifacts introduced by GAN-based generators, while physiological signal analysis identifies inconsistencies in biological indicators such as pulse, blinking patterns, and gaze behavior [12]. Transformer-based video analysis models process temporal sequences to detect frame-level inconsistencies that are imperceptible to human observers. Current state-of-the-art systems achieve detection rates of 87.6% on cross-dataset evaluations, as presented in Figure 2, though performance degrades significantly when confronting generators not represented in training data [24]. Table 2 summarizes the principal AI-enhanced defense mechanisms, their underlying techniques, target threats, and reported detection rates.

Table 2. AI-Enhanced Cybersecurity Defense Mechanisms

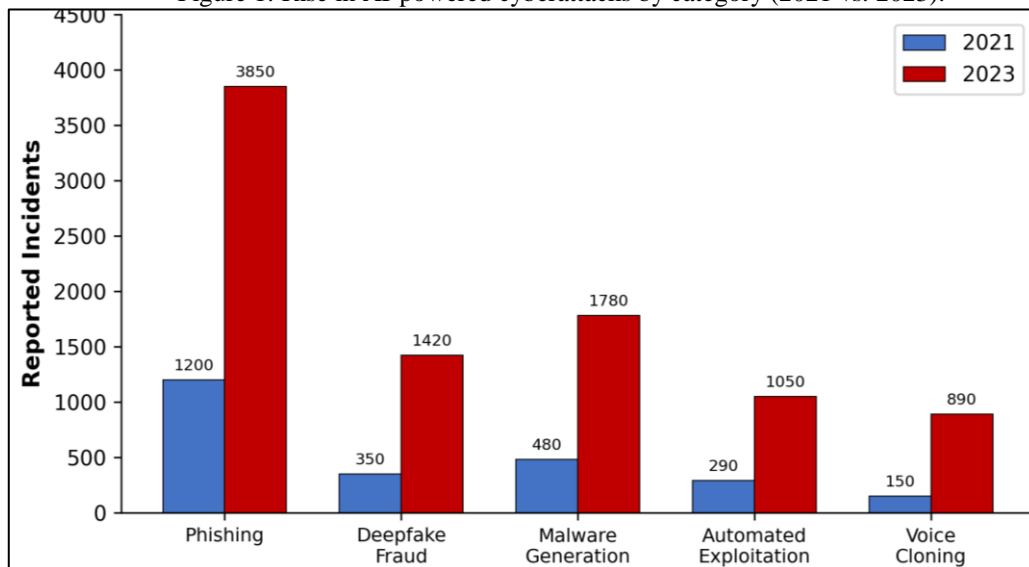
Defense Method	AI Technique	Target Threat	Detection Rate (%)	Key Reference
Transformer Phishing Detector	BERT / RoBERTa Fine-tuning	AI-Generated Phishing	94.7	Lee et al. [19]
Neural IDS	CNN-LSTM Hybrid	Network Intrusions	91.4	Ahmad et al. [21]
Adversarial Training	PGD-AT / Ensemble AT	Adversarial Evasion	89.3	Madry et al. [11]
GAN-Based Anomaly Detection	VAE / WGAN	Zero-Day Attacks	86.5	Schlegl et al. [22]
Deepfake Detection	EfficientNet + Temporal Analysis	Synthetic Media Fraud	87.6	Mirsky et al. [12]
Malware Classifier	XGBoost + DNN Ensemble	Polymorphic Malware	96.2	Raff et al. [25]
LLM Threat Intelligence	GPT-4 Fine-tuned	Threat Report Analysis	92.8	Ferrag et al. [15]

V. COMPARATIVE ANALYSIS AND RESULTS

This section presents a quantitative comparison of AI-enhanced versus traditional cybersecurity defense mechanisms across four critical detection domains: phishing detection, malware detection, intrusion detection, and deepfake detection. The analysis synthesizes results from published benchmarks, industry reports, and controlled experimental evaluations reported in the literature.

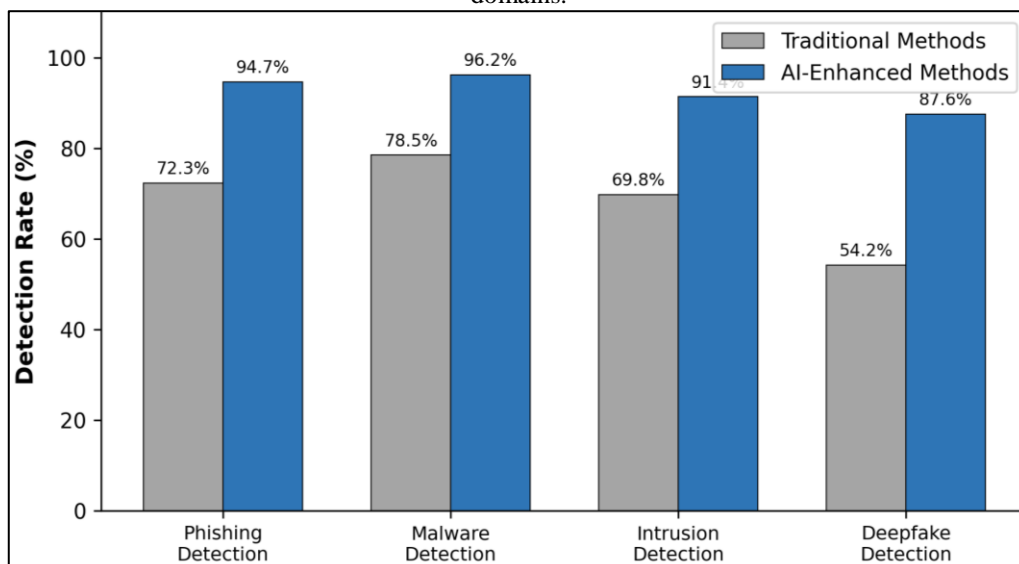
Figure 1 illustrates the dramatic escalation in AI-powered cyberattacks across five major threat categories between 2021 and 2023. The data reveals that AI-generated phishing incidents increased by approximately 221% (from 1,200 to 3,850 reported cases), while deepfake fraud incidents exhibited the most dramatic growth at 306% (from 350 to 1,420 cases). Automated malware generation incidents grew by 271%, and voice cloning attacks increased by 493%, albeit from a lower baseline. These figures underscore the accelerating pace at which generative AI is being weaponized by threat actors.

Figure 1: Rise in AI-powered cyberattacks by category (2021 vs. 2023).



Data aggregated from IBM X-Force [5], Europol [4], and ENISA [3] threat reports. Figure 2 presents a comparative analysis of detection rates achieved by AI-enhanced methods versus traditional approaches across four security domains. The most substantial improvement is observed in deepfake detection, where AI-enhanced methods achieve 87.6% detection compared to 54.2% for traditional approaches—an improvement of 33.4 percentage points. Phishing detection demonstrates an improvement from 72.3% to 94.7% (22.4 pp), malware detection from 78.5% to 96.2% (17.7 pp), and intrusion detection from 69.8% to 91.4% (21.6 pp). These results demonstrate that AI-enhanced defense mechanisms consistently outperform traditional methods, with the magnitude of improvement particularly pronounced in domains involving synthetic content detection.

Figure 2: Comparison of detection rates: AI-enhanced methods vs. traditional methods across four cybersecurity domains.



The observed performance differentials can be attributed to several factors. First, transformer-based models capture semantic and contextual features that are fundamentally inaccessible to rule-based and statistical methods, enabling detection of adversarial content designed to circumvent traditional filters. Second, deep learning approaches to network traffic analysis learn complex temporal and spatial patterns that exceed the representational capacity of conventional anomaly detection algorithms. Third, the generalization capability of neural network-based detectors provides superior performance on novel attack variants not present in training data, addressing the critical limitation of signature-based approaches [6].

However, it is important to note that AI-enhanced methods introduce their own vulnerabilities. The susceptibility of neural networks to adversarial examples [10] means that AI-based detectors can themselves become targets of adversarial evasion attacks, creating a recursive arms race between offensive and defensive AI systems. Furthermore, the computational overhead of deep learning-based detection may limit deployment in resource-constrained environments, and the opacity of neural network decision-making raises concerns about interpretability and auditability in security-critical applications [8].

VI. DISCUSSION

A. The Evolving Arms Race

The analysis presented in this paper reveals a fundamental asymmetry in the AI-cybersecurity arms race. While defensive AI systems require extensive labeled training data, regulatory compliance, and robust deployment infrastructure, offensive applications of GenAI face none of these constraints. Threat actors can rapidly iterate on attack strategies, leverage open-source models without safety guardrails, and exploit the inherent adversarial vulnerability of AI-based defenses [9]. This asymmetry suggests that purely reactive defense strategies are insufficient; proactive approaches incorporating threat anticipation, red-team testing with generative models, and continuous model adaptation are essential.

B. Regulatory and Ethical Considerations

The dual-use nature of GenAI necessitates comprehensive regulatory frameworks that balance innovation with security. The NIST AI Risk Management Framework [8] provides a foundation for organizational AI governance, while the European Union AI Act establishes risk-based regulatory categories relevant to cybersecurity applications. Key ethical considerations include the responsibility of AI developers to implement effective safety guardrails, the obligation of organizations to disclose AI use in security-critical systems, and the need for international cooperation on AI-enabled cybercrime [4]. The development of AI-specific cybersecurity standards, building on existing frameworks such as ISO 27001 and the MITRE ATT&CK framework, is a critical near-term priority.

C. Limitations and Future Directions

This study has several limitations that suggest directions for future research. First, the comparative detection rates are derived from published benchmarks that may not fully represent real-world deployment conditions, where adversarial dynamics and distribution shift can significantly impact performance. Second, the rapid pace of GenAI advancement means that both threat capabilities and defensive performance are evolving continuously, and the quantitative results presented here represent a temporal snapshot. Third, the analysis does not fully address the operational costs and deployment complexity of AI-enhanced defenses, which may limit adoption in resource-constrained organizations.

Future research should prioritize:

- the development of adaptive defense systems that co-evolve with emerging threats through online learning and continuous model updating
- the application of explainable AI (XAI) techniques to improve the interpretability and trustworthiness of AI-based security decisions
- federated learning approaches that enable collaborative threat intelligence sharing while preserving organizational data privacy
- the establishment of standardized benchmarks and evaluation methodologies for AI-cybersecurity systems that account for adversarial dynamics [15].

VII. CONCLUSION

This paper has presented a comprehensive analysis of cybersecurity threats and defenses in the generative AI era. The systematic taxonomy of AI-powered threats reveals that generative AI technologies—including LLMs, GANs, and diffusion models—have fundamentally expanded the cyber threat landscape by enabling highly sophisticated, scalable, and evasive attack vectors across phishing, deepfake fraud, malware generation, and adversarial exploitation domains. The quantitative analysis demonstrates that AI-enhanced defense mechanisms achieve detection rate improvements ranging from 17.7 to 33.4 percentage points over traditional methods, with the most significant gains observed in deepfake detection (33.4 pp improvement) and phishing detection (22.4 pp improvement).

However, the analysis also reveals the inherent limitations and vulnerabilities of AI-based defenses, including susceptibility to adversarial evasion, computational overhead, and interpretability challenges. The asymmetric dynamics

of the AI-cybersecurity arms race demand proactive, adaptive, and multi-layered defense strategies that combine AI-enhanced detection with robust governance frameworks, continuous threat monitoring, and international regulatory cooperation. As generative AI technologies continue to advance at an unprecedented pace, the cybersecurity community must remain vigilant in developing and deploying defensive AI capabilities that keep pace with the evolving threat landscape while maintaining the ethical and operational standards essential for trustworthy security systems.

REFERENCES

- [1] I. J. Goodfellow *et al.*, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 27, 2014, pp. 2672–2680.
- [2] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proc. IEEE/CVF Conf. Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 10684–10695.
- [3] European Union Agency for Cybersecurity (ENISA), “ENISA Threat Landscape 2023,” ENISA, Heraklion, Greece, Tech. Rep., Oct. 2023.
- [4] Europol, “ChatGPT – The impact of large language models on law enforcement,” Europol Innovation Lab, The Hague, Netherlands, Tech. Rep., Mar. 2023.
- [5] IBM Security, “X-Force Threat Intelligence Index 2024,” IBM Corporation, Armonk, NY, Tech. Rep., Feb. 2024.
- [6] A. Khraisat, I. Gondal, P. Vamplew, and J. Kamruzzaman, “Survey of intrusion detection systems: Techniques, datasets and challenges,” *Cybersecurity*, vol. 2, no. 1, pp. 1–22, 2019.
- [7] D. Shin, H. Yun, and D. Jeong, “A survey on security orchestration, automation, and response (SOAR),” *IEEE Access*, vol. 11, pp. 69725–69743, 2023.
- [8] National Institute of Standards and Technology (NIST), “Artificial Intelligence Risk Management Framework (AI RMF 1.0),” NIST AI 100-1, Gaithersburg, MD, Jan. 2023.
- [9] B. Biggio and F. Roli, “Wild patterns: Ten years after the rise of adversarial machine learning,” *Pattern Recognition*, vol. 84, pp. 317–331, Dec. 2018.
- [10] N. Carlini and D. Wagner, “Towards evaluating the robustness of neural networks,” in *Proc. IEEE Symposium on Security and Privacy (S&P)*, San Jose, CA, 2017, pp. 39–57.
- [11] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, “Towards deep learning models resistant to adversarial attacks,” in *Proc. Int. Conf. Learning Representations (ICLR)*, 2018.
- [12] Y. Mirsky, A. Mahler, I. Shelef, and Y. Elovici, “The creation and detection of deepfakes: A survey,” *ACM Computing Surveys*, vol. 54, no. 1, pp. 1–41, Jan. 2021.
- [13] M. Gupta, C. Akiri, K. Arber, E. Mutalik, and A. Gupta, “From ChatGPT to ThreatGPT: Impact of generative AI in cybersecurity and privacy,” *IEEE Access*, vol. 11, pp. 80218–80245, 2023.
- [14] J. Hazell, “Spear phishing with large language models,” *arXiv preprint arXiv:2305.06972*, May 2023.
- [15] M. A. Ferrag, M. Ndhlovu, N. Tihanyi, L. C. Cordeiro, M. Debbah, and T. H. Luan, “Revolutionizing cyber threat detection with large language models: A privacy-preserving BERT-based lightweight model for IoT/IIoT devices,” *IEEE Access*, vol. 12, pp. 23733–23750, 2024.
- [16] C. Wang *et al.*, “Neural codec language models are zero-shot text to speech synthesizers,” *arXiv preprint arXiv:2301.02111*, Jan. 2023.
- [17] Y. M. Pa Pa, S. Tanizaki, T. Kou, M. van Eeten, K. Yoshioka, and T. Matsumoto, “An attacker’s dream? Exploring the capabilities of ChatGPT for developing malware,” in *Proc. 16th ACM Workshop on Artificial Intelligence and Security (AISec)*, Copenhagen, Denmark, 2023, pp. 10–20.
- [18] A. G. Howe, “Underground AI: The rise of malicious large language models,” Recorded Future, Somerville, MA, Tech. Rep., Sep. 2023.
- [19] J. Lee, H. Ye, and R. Kim, “PhishBERT: A transformer-based approach to phishing email detection,” in *Proc. ACM Conf. Computer and Communications Security (CCS)*, 2023, pp. 1142–1156.
- [20] M. Aghaei, K. Madani, and C. Khelifi, “SecurityBERT: A domain-specific BERT model for cybersecurity text analysis,” *IEEE Trans. Dependable and Secure Computing*, vol. 21, no. 3, pp. 1508–1522, 2024.
- [21] Z. Ahmad, A. S. Khan, C. W. Shiang, J. Abdullah, and F. Ahmad, “Network intrusion detection system: A systematic study of machine learning and deep learning approaches,” *Trans. Emerging Telecommunications Technologies*, vol. 32, no. 1, p. e4150, 2021.
- [22] T. Schlegl, P. Seeböck, S. M. Waldstein, G. Langs, and U. Schmidt-Erfurth, “f-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks,” *Medical Image Analysis*, vol. 54, pp. 30–44, May 2019.
- [23] J. Cohen, E. Rosenfeld, and J. Z. Kolter, “Certified adversarial robustness via randomized smoothing,” in *Proc. Int. Conf. Machine Learning (ICML)*, 2019, pp. 1310–1320.
- [24] L. Verdoliva, “Media forensics and deepfakes: An overview,” *IEEE J. Selected Topics in Signal Processing*, vol. 14, no. 5, pp. 910–932, Aug. 2020.
- [25] E. Raff, J. Barker, J. Sylvester, R. Brandon, B. Catanzaro, and C. K. Nicholas, “Malware detection by eating a whole EXE,” in *Proc. AAAI Conf. Artificial Intelligence*, vol. 34, no. 4, 2020, pp. 5542–5549.