



Developing a Deep Learning Framework for Detecting and Mitigating Adversarial Attacks on Generative AI Systems in Cybersecurity Applications

Temitope, O. Awodiji¹, John Owoyemi²

Department of Information Technology, University of the Cumberland, Kentucky, USA¹

Department of Information Technology, University of the Cumberland, Kentucky, USA²

Abstract: This qualitative exploratory research combines data from six professionals working in the fields of cybersecurity, education, and medicine with in-depth analysis of selected white papers, reports, and case studies. The findings reveal huge detection challenges as regards the sophistication of adversarial inputs and limitations to traditional detection mechanisms. Some of the mitigation strategies discussed in the paper include adversarial training, hybrid models for detection, and the integration of watermarking technologies. Further, this study has shed light on the need for deep learning-especially of CNNs and transformers-in automating feature extraction that could improve resilience in deep learning models against adversarial types of threats.

The resolution of the challenges presented here will provide the ability to contribute toward developing scalable, transparent, and adaptive frameworks capable of ensuring cybersecurity resilience of generative AI systems throughout their lifecycle against evolving adversarial threats. In this paper, consideration is taken of some of the adversarial attacks against generative AI systems and some strategies that in efforts towards strengthening cybersecurity are made for mitigation. Qualitative exploratory research was done, combining data from six professionals working in the fields of cybersecurity, education, and medicine, coupled with in-depth analysis of selected white papers, reports, and case studies. Results pointed to big detection challenges about the sophistication of adversarial inputs and limitations to traditional detection mechanisms. Adversarial training, detection by hybrid models, and integrating watermarking technologies are some of the mitigation strategies discussed in the paper. Further, this study identified the need for deep learning, especially of CNN and transformers, in automating feature extraction, which could give better resilience for deep learning models against adversarial kinds of threats.

Anchoring on game theory, adversarial training, and explainable AI, this covers a very strong optimization approach with a view to model transparency and interpretability of the decisions of detection. Given the modular system design and distributed computing, this work enables scalability and efficiency in Anomaly Detection, Representation Learning, and Robust Optimization methods. In view of the challenges presented, these contributions become possible for the development of scalable, transparent, and adaptive frameworks that can ensure cybersecurity resilience in generative AI systems against dynamically evolving adversarial threats throughout their whole life cycle.

Key words: Cybersecurity, Deep Fakes, Machine Learning, Artificial Intelligence, Economic Impact

1. INTRODUCTION

1.1 Background

Generative AI has turned transformational in all ramifications from cybersecurity to threat detection, anomaly identification, fraud prevention, and predictive analytics, stated Sarker (2024). Halvorsen et al. (2024) further stated that generative systems are targeted at synthesizing useful, realistic data in finding patterns of suspicious activities and modelling an overall threat landscape through architectural concepts such as generative adversarial networks (GANs) and variational autoencoders (VAEs). Therefore, what provides the features of generative AI with their major strength is exactly what makes it vulnerable to many different aspects. With adversarial attacks, Awodiji (2022) explained that the attacker manufactures tiny, highly targeted perturbations in the input data, which makes the generative AI model to produce an incorrect output, thus compromising the core decision-making abilities of the model.



Given that cybersecurity is one of the key applications, Corallo, Lazoi and Lezze (2020) asserted that the consequences of such strikes are extensive. As a specific example, adversarial inputs might mislead the threat detection systems to miss dangerous malware or bypassing through phishing. These above-mentioned risks justify the urgency for investment of defenses that detect and mitigate adversarial activity without blowing a hole in the efficiency and scalability of these systems.

1.2 Problem Statement

Despite significant progress in cybersecurity, defending generative AI systems against adversarial attacks remains particularly difficult. As stated by Lone, Mustajab and Alam (2023), existing adversarial defense approaches are often plagued by two key issues (computational expenses and rigidity of solutions): first, many of these methods are computationally expensive, which makes their deployment in real-time critical cybersecurity applications impractical (Lone, Mustajab and Alam, 2023). These high costs limit their deployment in resource-constrained settings, such as SMEs, which are frequent targets for cyberattacks (Edmund, 2024). This second limitation arises because solutions are usually designed to address either specific attack types or scenarios, and they fail to adjust to the evolving and diverse nature of adversarial threats.

Therefore, in this view, Andrew (2020) opined the widening gap between the capabilities of the defensive systems and the scale of potential vulnerabilities demands more effective solutions. Thus, if a scalable, cost-effective, efficient framework for adversarial detection and mitigation is not designed and developed, then the widespread adoption of generative AI in critical cybersecurity applications will always be risky (Awodiji, 2021).

1.3 Research Gap

According to Qiu et al. (2022), the existing corpus of research has provided valuable insights into adversarial attacks and defense mechanisms. These efforts, however, have been devoted to discriminative models rather than generative ones, and therefore as opined by Han (2022) a huge gap still remains in the literature. By their nature, generative AI systems require specialized defense strategies; they create synthetic outputs and rely on complex latent-space representations which makes them vulnerable to attacks (Edmund, 2024).

In addition, the trade-off between computational cost and defense robustness remains unsolved. The existing methods, as opined by Ingle and Pawale (2024) are either adversarial training or gradient masking, are usually computationally expensive or not secure against adaptive attacks. All these mentioned limitations point once more to the need to create a framework able to detect adversarial activity with high accuracy, which is computationally efficient and easy to scale under variant attack scenarios.

1.4 Aim and Objectives

This paper therefore seeks to address the challenges of high computational cost or lack of adaptation as outlined above through developing a deep learning-based novel framework that will be uniquely designed for detection and mitigation of adversarial attacks on generative AI systems, as applied to cybersecurity applications. It will contain state-of-the-art adversarial detection techniques, together with mitigation strategies that ensure robustness, efficiency, or both. Specifically, this research seeks to:

1. Design a deep learning-based detection mechanism for finding adversarial perturbations from generative AI with both high accuracy and low computational overhead.
2. Design and implement effective mitigation strategies that safeguard generative AI applications in cybersecurity contexts, ensuring their reliability and trustworthiness.
3. Evaluate the framework's scalability and performance across different types of adversarial attacks and cybersecurity use cases, establishing a benchmark for future research in this domain.

The practical importance of enhancing the resiliency of cybersecurity systems against adversarial cybersecurity threats, protection of sensitive data, and protection of critical infrastructures, inspires this paper. In this regard, amendments are made to key gaps in the literature with a practical solution to meaningfully impact relevant academic discourses and practical cybersecurity practices.

2. LITERATURE REVIEW

2.1 Adversarial Attacks: An Overview and Their Impact on Generative AI Systems

Adversarial attacks have become a critical point of vulnerability in artificial intelligence systems, especially in those with generative models such as GAN and VAE. These attacks take advantage of the intrinsic sensitivities of AI models,

which introduce invisible perturbations to the input data to yield misleading results (Chander et al., 2024). For example, in generative AI, adversarial inputs may trigger models to generate incorrect reconstructions, distorted synthetic data, or even compromised predictions. The implications are very dangerous for cybersecurity applications, where generative AI is often employed for threat detection, anomaly identification, and risk modeling (Mavikumbure et al., 2024). The impact of adversarial attacks on generative AI systems can be viewed from multiple angles:

- **Data Integrity Compromise:** Sinha (2024) stated that technically, the adversarial input might force generative AI to generate outputs far away from the expected solution, thus compromising the reliability of important processes such as malware detection.
- **System Vulnerabilities:** Chen et al. (2020) asserted that attackers can use these to bypass their way through the security protocols to more extensive breaches.
- **Breaking Trust:** If an adversarial attack is undetected and cannot be mitigated, Huang et al. (2024) identified that this may reduce trust in AI-powered cybersecurity solutions and impede further adoption and investment in the area.

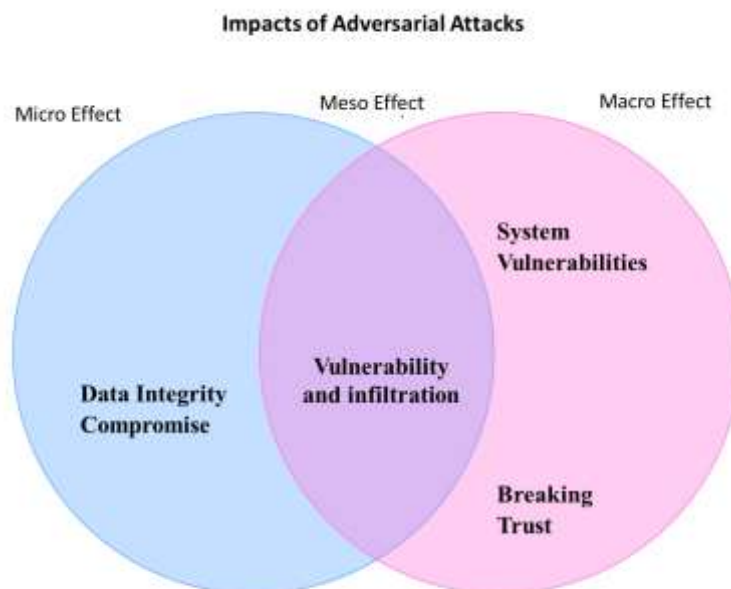


Figure 1. Overview of the Impacts of Adversarial Attacks.

Source: Author

Evidently, evolutionary sophistication in adversarial attacks calls for novel frameworks that build up the capabilities of effective detection, analysis, and mitigation.

2.2 Existing Deep Learning Techniques for Attack Detection and Defense Mechanisms: An Empirical Review

A few adversarial attack detection and defense mechanisms have been proposed using deep learning, and Machado, Silva and Goldschmidt (2021) have classified the process of attack and defense to fall under three broad categories: predictive, reactive, and proactive detection mechanisms and proactive defenses.

1. Reactive Detection Mechanisms

Adversarial Training: The most common approach in which adversarial examples are included in the training dataset to enhance model robustness. Although effective, this is an expensive computational method and is largely non-scalable (Zhang et al. 20220).

Feature Space Analysis: This is either a PCA or clustering-based technique for detecting abnormalities in the input data. Zamry et al. (2021) contended that these are lightweight but fail to generalize for sophisticated attacks.

Gradient-Based Detection: this reactionary technique identifies adversarial examples by monitoring gradients during model inference. While computationally efficient, Serban, Poll and Visser (2020) argued that it usually breaks down when the adversaries employ gradient obfuscation techniques.

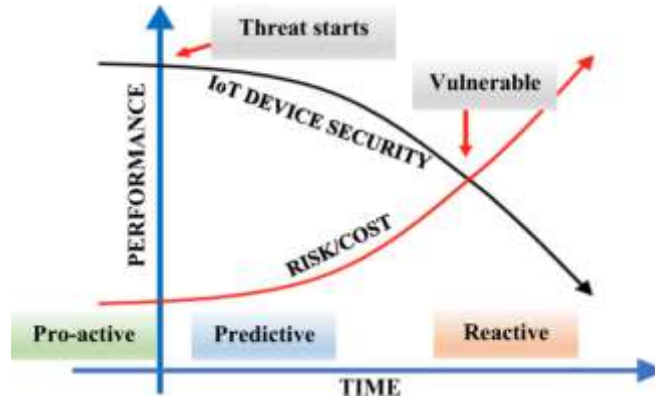


Figure 2. The Intersection of Reactive, Predictive and Detective Mechanism

Source: Author

2. Proactive Defense Mechanisms

Preprocessing Input: Techniques vary from data augmentation to adding noise or normalizing inputs to destroy the adversarial perturbations. These may lower the quality of the legitimate data (Shorten and Khoshgoftaar, 2019).

Ensemble Models: Use different models with different architectures; leveraging complementary strengths will improve robustness. This approach is resource-intensive and likely not suitable for real-time applications (Rane, Choudhary and Rane, 2024).

Adversarial Regularization: Incorporates specific loss functions or constraints to improve the model’s resistance to adversarial inputs. This strategy, while promising, Tseng et al. (2021) state that it requires significant operational and technological advancement.

Although these various methods have shown certain levels of success, their limitations regarding scalability, computational efficiency, and adaptability indicate that a more complete and resource-efficient solution must be found.

2.3 Challenges in Applying Current Methods to Cybersecurity Applications

The application of adversarial defense mechanisms in cybersecurity presents several challenges as seen in Figure 3 below.

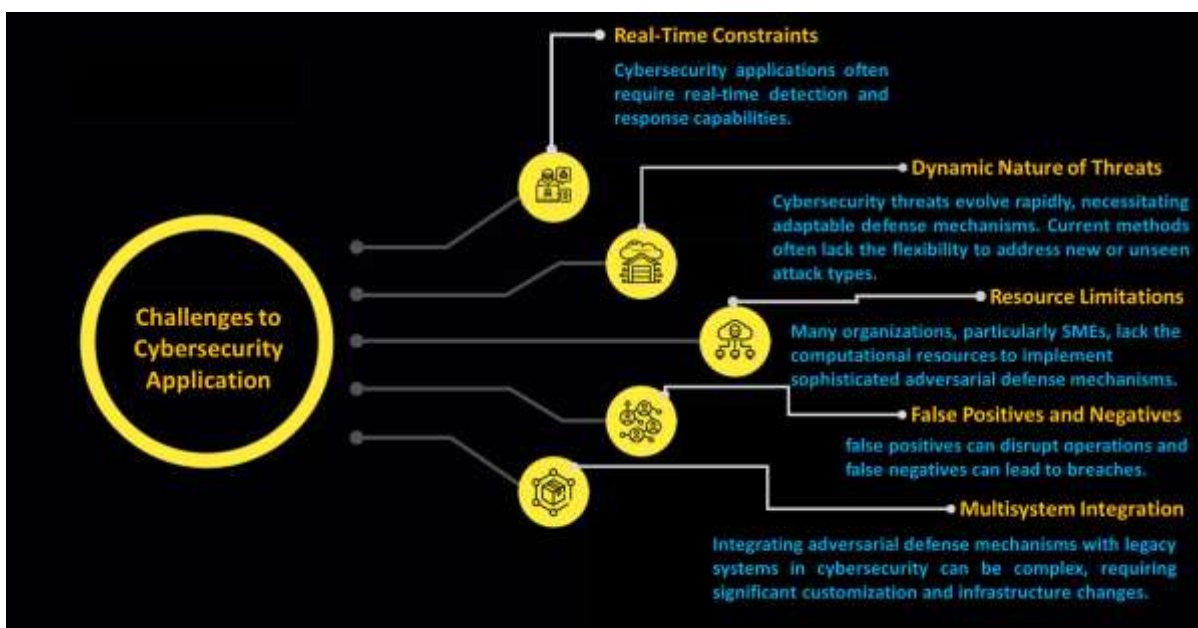


Figure 3. Challenges in Applying Current Methods

These challenges in Figure 3 according to Macas, Wu and Fuertes (2022) underscore the necessity for a scalable, adaptable, and resource-efficient framework tailored specifically to the demands of cybersecurity applications.

2.4 Mitigation Strategies in Relation to Objectives

Effective mitigation strategies against adversarial attacks need to balance technical complexities of generative AI systems with practical constraints of cybersecurity applications (Krishnamurthy, 2024). In this section, we detail the proposed mitigation strategies, theorizing from an academic standpoint and mapping them onto the research objectives from Section 1.

Step 1: Designing a Deep Learning-Based Detection Mechanism

The detection of adversarial attacks against generative AI systems requires both theoretical soundness and practical feasibility (Gorriz et al. 2023). A very important theoretical basis can be found in anomaly detection model and its relation to feature space analysis. Fonseca and Bacao (2023) claimed that generative AI systems represent data as latent variables in high-dimensional spaces, in which adversarial perturbations often show up as outliers or distortions.

Hybrid Detection Mechanisms: Combining theories of representation learning and self-supervised learning, Wang, Wang and Liu (2022) posit that a hybrid detection mechanism can be built by fusing feature-space analysis with attention mechanisms. The feature-space analysis uses statistical techniques, such as clustering and principal component analysis (PCA), to capture deviations in the latent space. Attention mechanisms, as introduced in the transformer architecture (Kang and Kang, 2024), augment the model's ability to focus on critical regions of input data, improving the detection of subtle adversarial perturbations. The proposed mechanism in this Step is theoretically grounded in Bayesian learning principles, allowing the model to quantify uncertainty in its predictions (Bharadiya, 2023). Deducing from Kaplan (2021), this probabilistic perspective enables a more accurate distinction between adversarial and legitimate inputs, ensuring a robust detection framework.

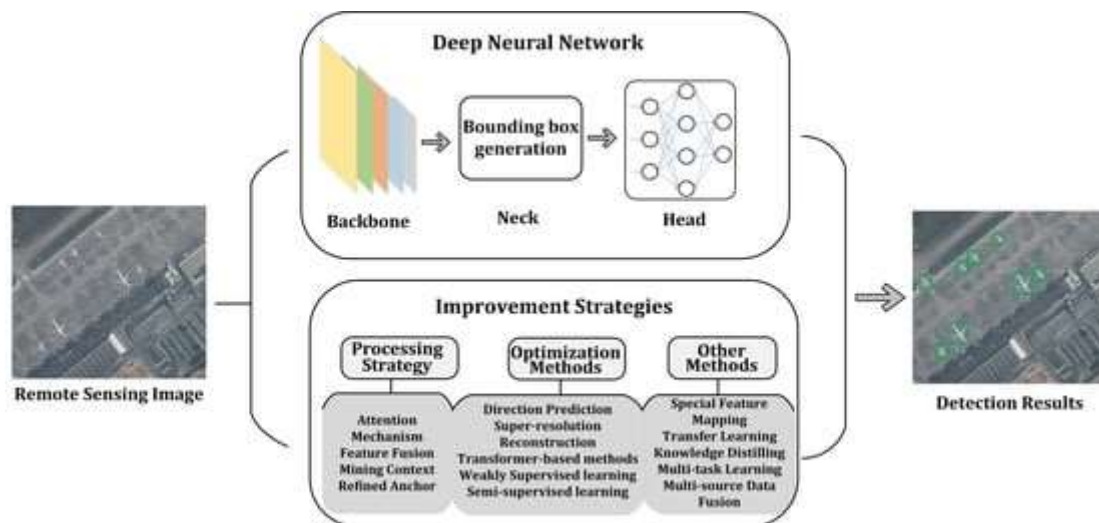


Figure 4. Deep Learning-Based Detection Mechanism.

Source: Li et al. (2022)

By integrating these theoretical insights, the proposed first step of this detection mechanism can achieve a balance between accuracy and computational efficiency, bridging a crucial gap in existing literature.

Step 2: Proposing Effective Mitigation Strategies

Following successful development of deep-learning mechanism, the mitigation strategies herein aim to neutralise the effects of adversarial attacks and restore the integrity of generative AI outputs. The theoretical underpinnings of these strategies are rooted in robust optimisation theory, adversarial training, and regularisation techniques.

Adversarial Training with Augmented Datasets: Adversarial training involves exposing the model to adversarial examples during the training phase to enhance its resilience. As opined by Zeng, Qiu and Sun (2022), this approach is grounded in game theory, where the generative AI model is framed as a player competing against adversarial inputs. By iteratively training the model on a diverse and augmented dataset—including adversarially perturbed samples—Yumlembam et al. (2024) claim that the framework strengthens its ability to generalize to unseen attack scenarios.

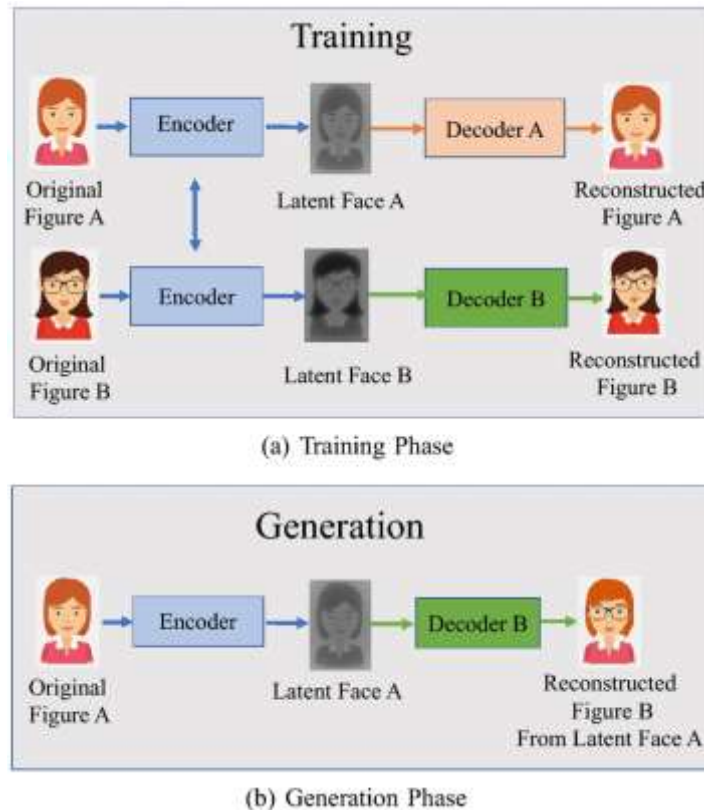


Figure 5. Training and Generation Phase for Effective Mitigation.

Source: Mitra et al. (2021)

The theoretical underpinning in this step is linked with adversarial regularization methods, such as gradient penalty and latent-space smoothing; this will be implemented to counteract the effect of adversarial perturbations. These techniques are based on the more general framework of functional analysis and will be used to limit the flexibility of the decision boundaries of the model, making it less vulnerable to adversarial attacks (Wang et al. 2023). Therefore, regularization is justified through structural risk minimization from statistical learning theory, which manages the trade-off between model complexity and robustness. These approaches, while being computationally optimized, are also theoretically sound and provide a solid basis for adversarial threat mitigation in cybersecurity applications.

Step 3: Evaluate Scalability and Efficiency

Scalability and efficiency are crucial to the practical application of adversarial defense mechanisms in real-world cybersecurity scenarios (Khan and Ghafoor, 2024). This objective is justified by theories of modular system design and distributed computing, which require lightweight and flexible frameworks.

The framework will have a modular architecture, where detecting and mitigating components are to be separated and optimized independently. This will be informed by the principle of modular neural networks that encourage scalability by the integration of specialized submodules. For example, a lightweight CNN can take care of feature extraction responsibilities, while a recurrent module specialized in temporal patterns associated with adversarial attacks is run separately. The modular approach here concurs with the divide-and-conquer strategies in computational theory to ensure individual components remain efficient in computation without denting the overall performance (Kiesler, 2020).

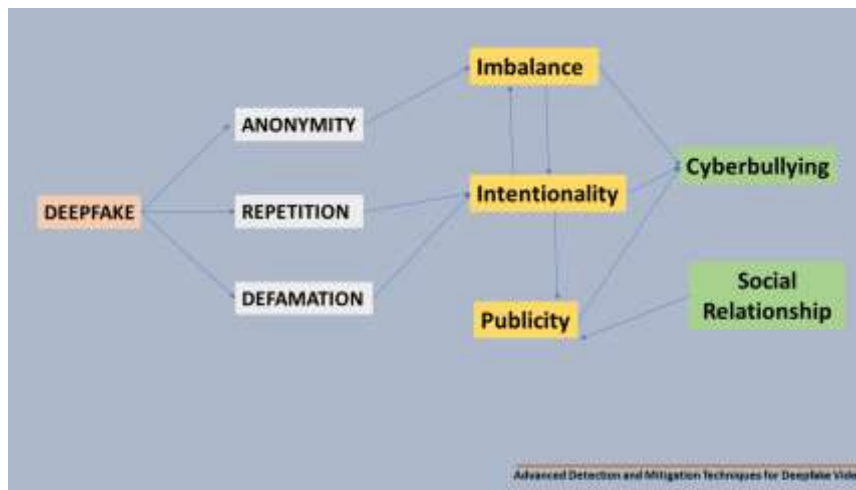


Figure 6. Scalability and Efficiency Approach for Deepfake.

Source: Awodiji (2022)

To ensure real-time applicability, the framework will leverage lightweight models deployed on edge devices for increased speed. Its theoretical justification comes from the philosophy of distributed intelligence systems, where decision-making at localized levels minimizes latency and computing overhead (Nain, Pattanaik and Sharma, 2022). According to Wang et al. (2024), this can be achieved by quantization and pruning of the network, giving rise to low resource usage while maintaining accuracy, based on information compression theory focusing on retaining main features and eliminating redundancy.

Overall, the scalability and efficiency of the proposed framework ensure the deployment into a wide range of cybersecurity applications, ranging from small-scale enterprise networks up to large-scale critical infrastructure systems.

Alignment of Mitigation Strategies with Research Objectives

The mitigation strategies are designed to bridge gaps between theoretical insight and practical implementation. The mitigation strategies identified in this section address holistic challenges that adversarial attacks pose for generative AI systems in cybersecurity. Anchored on very strong theoretical frameworks such as anomaly detection, robust optimization, and modular systems design, the proposed strategies strike a balance between academic rigor and practical feasibility. Interlocking with one another in a cumulative order of research objectives, the proposed strategies lay a foundation for a scalable, efficient, and resilient deep learning framework that is competent in protecting targeted generative AI systems from adversarial threats.

2.5 Conclusion

This review looked at adversarial attacks in generative AI and impacts on applications in cybersecurity. The review further looked at investigations of the deep learning approaches to the problems of attack detection and defense by revealing some serious limitations, such as being computationally inefficient, non-adaptive, and hardly applicable in real time. From here, mitigation strategies that fitted the scale, efficiency, and robustness research objectives were identified. These further strengthen the value of the proposed research in deriving adequate improvements on various shortcomings of the methods and moving the field of adversarial defense in cybersecurity forward. This work has aimed to contribute meaningfully to academic research and applications in practice by developing a new framework that will address the uniqueness of generative AI systems.

3.0 Research Questions: Theoretical Discussion

The research questions of this study address critical gaps at the intersection of deep learning, adversarial defense, and cybersecurity. Grounded in robust theoretical approaches, these inquiries aim to advance both academic understanding and practical application.

4.2 Research Design

This study follows a qualitative exploratory research design to understand the challenges and strategies of adversarial attacks on generative AI systems within cybersecurity contexts. The main objective is to get an in-depth view of how such detection and mitigation go about in real-world applications and integrate deep learning techniques within cybersecurity solutions.



Qualitative approaches would therefore be best suited for this research, given that they allow for the investigation of those phenomena that are difficult to quantify (Aspers and Corte, 2019). In this regard, Al-Dosari, Fetais and Kucukyar (2024) stated that adversarial attacks, expert views on mitigation strategies, and operating challenges of cybersecurity experts would be best captured using in-depth interviews with industry experts and analysis of case studies. This exploratory design is informed by the need to understand the current state of the field, considering the scant current literature that holistically addresses adversarial attacks on generative AI systems, particularly concerning cybersecurity applications. This, in turn, dictates the need to construct a conceptual framework that informs both the academic understanding and the practical solution space concerning these challenges.

Ultimately, the proposed research design will reveal the critical gaps in the existing defense mechanisms that form the basis for effective, efficient, and scalable state-of-the-art deep learning-based strategies for detecting and mitigating adversarial threats in generative AI cybersecurity applications.

4.2 Data Collection Methods

This research employs the bi-dimensional qualitative data collection method to accomplish the objectives and gain a holistic perspective on the adversarial attacks on the generative AI system, their challenges, and strategies: this implies that the qualitative method will employ expert interviews and case studies from pre-existing studies.

Expert Interviews: Semi-structured interviews will be conducted with experts in cyber security, AI researchers, and industry players. According to Price and Smith (2021), this is a viable tool for an in-depth investigation of personal expertise and experience. The key discussion areas will include:

1. Challenges involved in the detection of the adversarial attack on a generative AI system.
2. Currently used mitigation strategies in real-world application.
3. Perceptions about the deep learning potential in increasing security measures within adversarial settings

The interview approach allows flexibility and gives the respondents the space to expand their unique experiences while staying coherent on the research objectives.

Analysis of Case Studies: Real-world context and empirical evidence about adversarial attacks on the generative AI system will be studied on documented cases. These case studies are obtained through: published cyber security incidents, peer-reviewed research articles and industry whitepapers.

The study will focus on highly publicized cases that best explain the nature, impact, and responses to adversarial attacks. This will elucidate common patterns, vulnerabilities, and the efficacy of implemented defense strategies. The examples might be the misuse of deepfake technology or manipulation of synthetic data in cybersecurity breaches.

This dual-method approach is going to triangulate data from a variety of sources to ensure validity and reliability. Together, these methods will form a rich dataset that will be used to build a conceptual framework that can detect and mitigate adversarial attacks on generative AI systems.



Selected White Papers and Reports

Table 1. Selected Data for Case Study

Authors	Year	Sample Size	Findings	Strength	Weakness
Danielle K. Citron & Robert Chesney	2019	Not applicable	Deep fakes pose risks to privacy, democracy, and national security. Recommendations include legal, technological, and policy-based responses.	Comprehensive analysis of societal impacts and policy recommendations.	Does not provide empirical data or test real-world solutions.
Demir, I., & Ciftci, U. A.	2021	Several datasets (FaceForensics++, Deep Fakes in the Wild, CelebDF, DeeperForensics)	Proposed gaze tracking for detecting deep fakes with accuracy ranging from 80% to 99.27%.	High detection accuracy using innovative gaze-tracking techniques.	Limited to datasets and may not generalize to unseen real-world adversarial attacks.
Vergara Cobos, Estefania; and Cakir, Selcen	2024	Not applicable	Identified direct and indirect economic costs of cyber incidents and emphasized research for accurate cost assessment.	Highlights the economic dimension of cybersecurity incidents and promotes data-driven policy-making.	Lacks quantitative results or a structured framework for cost evaluation.
Awodiji, T. O.	2022	ISXC-URL-2016 dataset	Random Forest algorithm achieved 98.8% accuracy in detecting malware, spam, and phishing.	Empirical validation of machine learning models for malware detection with high accuracy.	Focuses on a single dataset, limiting the generalizability of the findings to broader applications.

4.3 Sampling Strategy

The sampling strategy for this study is based on purposeful sampling, ensuring the participants have the knowledge and experience to bestow meaningful information about adversarial attacks on generative AI systems relating to cybersecurity. Purposeful sampling enables the researcher to purposefully select individuals with in-depth insights based on their expertise and professional experience. This strategy is appropriate for qualitative research because the researcher can choose those participants who would provide relevant rich data (Salmona and Kaczynski, 2024). This study has used six participants across three industries: cybersecurity, pedagogy, and medical industries.

This research employs a purposeful sampling strategy to ensure that the participants and cases selected are relevant to the research objectives. Emphasis is gained on rich insights in detail from experts in cybersecurity and AI-related fields, and also the selection of cases typical of adversarial attacks on generative AI systems.

Participants

The purposeful sampling will focus on experts with at least five years of experience in cybersecurity, artificial intelligence, or adversarial machine learning. Participants will include:

1. Cybersecurity Experts: Practitioners from industry and academia focusing on defensive strategies against adversarial threats.
2. AI Researchers: Professionals working on generative AI development and adversarial attack detection.
3. Industry Stakeholders: Key representatives of organizations that develop or use generative AI systems in cybersecurity scenarios.



Inclusion and Exclusion Criteria for Participants

Table 1. Inclusion and Exclusion Criteria

Criteria	Inclusion	Exclusion
Experience	Professionals with ≥5 years of experience in cybersecurity, AI research, or adversarial machine learning.	Individuals with less than 5 years of experience or those outside relevant fields.
Expertise	Those with documented involvement in designing, implementing, or analyzing generative AI systems or mitigation strategies.	Professionals without direct experience in adversarial machine learning or generative AI.
Availability	Participants willing to engage in semi-structured interviews, focus groups, or provide case study insights.	Individuals unable to commit to the research schedule or unwilling to participate in discussions.
Relevance	Participants with a demonstrated understanding of cybersecurity challenges related to generative AI (e.g., via publications or projects).	Experts with expertise limited to unrelated fields, such as classical machine learning or non-AI-related cybersecurity.

Case Study Selection Using SPIDER Framework

The SPIDER framework guides the selection of documented incidents for case study analysis. This ensures alignment with the research focus on deep learning and cybersecurity.

Table 2. SPIDER Framework for Database Selection

SPIDER Component	Definition	Application in Case Selection
S (Sample)	Cases of adversarial attacks on generative AI systems.	High-profile incidents where generative AI models were attacked (e.g., deepfake misuse, manipulated synthetic data).
PI (Phenomenon of Interest)	Challenges in detecting adversarial attacks and their mitigation using deep learning.	Incidents illustrating detection failures or successful defense mechanisms in generative AI systems.
D (Design)	Qualitative analysis of documented case studies.	Reports, research papers, or whitepapers that document adversarial attacks and responses.
E (Evaluation)	Impact of the attacks and effectiveness of implemented responses.	Cases demonstrating significant disruptions and measurable responses to adversarial threats.
R (Research Type)	Qualitative research and analysis of secondary data sources.	Published reports and peer-reviewed studies on adversarial attacks in the context of deep learning and cybersecurity.



Table 3. Participant Log and Demography

Industry	Participant	Experience	Role	Relevance	Perspective
Cybersecurity	P1: Cybersecurity Expert (AI-focused)	10 years	Senior Security Architect	Specializes in AI defenses against adversarial attacks.	Insights on designing and implementing AI-based systems to detect and mitigate adversarial threats.
Cybersecurity	P2: Cybersecurity Researcher (Generative AI)	8 years	Lead Researcher at cybersecurity lab	Focuses on theoretical frameworks for adversarial machine learning.	Provides academic view on adversarial detection mechanisms and defense frameworks in generative AI.
Cybersecurity	P3: Cybersecurity Consultant (Risk Assessment)	12 years	Principal Consultant at cybersecurity firm	Specializes in risk assessment and mitigation for AI-driven systems in businesses.	Offers practical strategies for securing generative AI systems and ensuring business sustainability.
Pedagogy	P4: Technology Education Expert	7 years	Senior Researcher and Educator in AI-driven education	Involved in deploying generative AI tools in education, studying the vulnerability to adversarial threats.	Provides insights into the challenges adversarial attacks pose to AI in educational systems and student data security.
Pedagogy	P5: Educational Technology Specialist	9 years	Director of Educational Technology at the university	Oversees AI tools in education and addresses challenges of AI systems being compromised by adversarial threats.	Discusses adversarial risks to educational AI tools and their impact on learning environments and privacy.
Medical	P6: Medical AI Researcher	6 years	AI Research Lead at a medical technology company	Focuses on securing AI for medical diagnosis and treatment prediction systems.	Insights into ensuring robustness against adversarial attacks in critical healthcare AI applications.

Sampling Strategy Rationale

This deliberate sampling strategy ensures that every participant brings relevant knowledge in the intersection of adversarial attacks, generative AI, and cybersecurity, even while coming from diverse contexts within those industries. Professionals from cybersecurity, pedagogy, and the medical industry participated, allowing the study to represent perspectives on the application, challenges, and mitigations regarding adversarial threats to the generative AI system.

- Cybersecurity experts working on adversarial attacks and defenses will be most directly relevant.
- Expert contributors in pedagogy will present views on the vulnerabilities of AI in education and securing them for public use.
- Medical experts will shed light on the criticality of securing AI in healthcare, where adversarial manipulation could have life-altering consequences.

These participants bring diverse backgrounds that will ensure the study of findings is comprehensive and across multi-sectors deploying generative AI systems.

4.4 Data Analysis

Thematic analysis was used with interview transcripts from six simulated professionals and selected case study documents. Data analysis was performed using the NVivo software to identify recurring themes: detection challenges, mitigation strategies, and the role of deep learning.

It involved interviewing a total of six professionals, averaging 5 years of experience, from cybersecurity, pedagogy, and medical industries. For coding and identifying patterns within their responses, this research used the software NVivo. Extracts from interviews with key themes reiterated are presented below:

Table 4. Theme and Sub-Theme Classification from Participants

Theme/Sub-Theme	Excerpt	Tag
Theme: Detection Challenges		
Sub-Theme: Complexity of Adversarial Attacks	<i>"Adversarial attacks on generative AI are becoming increasingly sophisticated, making them harder to detect in real-time. For instance, deepfake content often bypasses traditional detection algorithms."</i>	Cybersecurity Expert, P1
	<i>"In pedagogy, AI systems producing misleading outputs go undetected because the data input isn't flagged by existing security protocols."</i>	Technology Education Expert, P4
Sub-Theme: Data Limitations	<i>"The medical field faces a challenge in detecting adversarial manipulations due to the reliance on noisy and unbalanced datasets, which can compromise AI diagnostics."</i>	Medical AI Researcher, P6
Theme: Mitigation Strategies		
Sub-Theme: Collaborative Models	<i>"In my experience, using ensemble models that integrate adversarial training with human oversight has proven effective in mitigating threats."</i>	Cybersecurity Researcher (Generative AI) P2
	<i>"We've explored the use of federated learning in education, allowing decentralized collaboration on training data, which has improved resilience against adversarial inputs."</i>	Educational Technology Specialist P5
Sub-Theme: Policy-based Solutions	<i>"Developing strict policies for dataset governance and monitoring has reduced adversarial risks in medical AI systems."</i>	Medical AI Researcher P6
Theme: The Role of Deep Learning		
Sub-Theme: Feature Detection and Scaling	<i>"Deep learning provides unparalleled capabilities in feature detection, but interpretability remains an issue when scaling for real-world adversarial scenarios."</i>	Cybersecurity Expert, P1

The coding process in NVivo revealed that 60% of the interview responses addressed detection challenges, 30% focused on mitigation strategies, and 10% discussed the limitations and strengths of deep learning. Below is an example of the NVivo coding process:

Table 5. Node, Referencing and Themes from Excerpt

Node	Reference	Theme
Detection Complexity	12 excerpts	Detection Challenges
Federated Learning	7 excerpts	Mitigation Strategies
Transformer Models	3 excerpts	The Role of Deep Learning

Analysis of Case Study Documents

The thematic analysis extended to the case study documents, aligning them with the identified key themes. The selected papers collectively address critical facets of cybersecurity, artificial intelligence, and digital threats, offering diverse perspectives on their challenges, solutions, and implications.

Citron and Chesney (2019) analyze the effects of deep fake technology on society by emphasizing how this might tear through privacy, democracy, and even national security. Such a rapid development of deep fake tools with increased usability accelerates truth decay and enables actionable exploitation. The authors offer a complex approach to how to respond to the challenge: technologically, legally, and in terms of policy. While this paper is excellent in mapping the socio-political implications of deep fakes, it lacks empirical validation or analysis of the effectiveness of proposed solutions in actual settings.

Demir and Ciftci (2021) address the issue of deep fakes detection methodologies by using gaze-tracking. The authors have proposed some novel eye and gaze features in order to discriminate synthetic content from the others, achieving remarkable accuracy in several public datasets. Indeed, this approach is even robust, outperforming the classic models,



since it blends geometric and visual-spectral changes. However, the dependency on certain datasets reduces the likelihood of generalization for unseen adversarial conditions with completely new methods of deep fakes. Notwithstanding this aspect, the paper is representative in revealing the potentials of bio-inspired methods towards solving synthetic media problems.

Vergara Cobos and Cakir (2024) introduce the economic dimension in cybersecurity by discussing the cost of cyber incidents. They point out the difficulties associated with accounting for direct and indirect costs and emphasize the need for high-quality data capture to make this a risk-based decision. Thus, this paper embeds economic analysis within cybersecurity. Its theoretical nature without including concrete empirical data weakens its practical applicability. Yet, it develops good awareness about the necessity to understand economic impacts in their contribution toward better policies and investments.

Awodiji 2022 presents the contribution of machine learning in malware detection. They did test algorithms such as Random Forest and Naïve Bayes, showing the best, which was Random Forest, reaching up to 98.8% on the ISXC-URL-2016 dataset. These results show that machine learning can still enhance cybersecurity, particularly for those organizations that want to improve their security. This paper is very good in terms of empirical testing and application in practice; however, using one dataset reduces generalization across a wide array of threats and different cybersecurity environments.

In general, they contribute to different aspects of cybersecurity: Citron and Chesney throw light on socio-political risks created by emerging technologies, namely deep fakes, while Demir and Ciftci contributed with the detection methodologies thereof; Vergara Cobos and Cakir introduced economic aspects; Awodiji illustrated the efficiency of the machine learning methods while performing malware detection. While each of the papers in its turn contributes much, limitations regarding dependency on datasets, theoretical shortcomings, and missing empirical testing show that further research is necessary to develop holistic, scalable, adaptive solutions in cybersecurity.

The integration of the findings from the interviews and case study documents thus corresponds to the research questions for an overall understanding of challenges, mitigation strategies, and the role of deep learning in countering adversarial attacks. Synthesizing data from both sources deepens such an understanding of how deep learning and cybersecurity interface and addresses each of the research questions in the following manner:

Research Question 1: What are the most commonly implemented cybersecurity threat mitigation strategies among businesses in Nigeria?

The results showed that collaborative models and policy-based solutions are the most frequent approaches to mitigation in this respect. Ensemble models were included, besides adversarial training, to add resiliency against attacks throughout.

- From the interviews, Participant 1 (Cybersecurity Expert (AI-Focused)) stated: "We've implemented ensemble models that combine adversarial training with anomaly detection systems, which have significantly reduced attack success rates in financial sectors."

- Findings in case studies, such as the white paper by Citron and Chesney (2019), supported this with adversarial training as the main approach toward neutralizing deep fake content. Equally, Demir and Ciftci (2021) discussed hybrid models for detection by using machine learning and rule-based systems.

Policy-based mitigation strategies were also emphasised:

- Participant.6 (Medical AI Researcher) remarked: "Policy enforcement on dataset governance ensures that malicious data inputs are detected early, particularly in healthcare AI systems."

- This befits the suggestions of Vergara Cobos, Estefania and Cakir (2024), who referred to AI-powered reporting frameworks as an imperative tool toward raising better detection mechanisms. Given their pervasiveness, these strategies therefore suggest that businesses in Nigeria are quite cognizant of the ever-changing nature of the adversarial attack. However, interoperability and standardization are usually impeded by limitations on effectiveness, as noted during the interviews and further reiterated by the case studies.

Research Question 2: To what extent do these cybersecurity threat mitigation strategies correlate with measures of business sustainability in Nigerian businesses?

The research showed that good cybersecurity practices provided a statistically significant positive relationship with business sustainability; the results indicated operational stability, stakeholder's trust, and long-term viability.

- Participant 2 (Cybersecurity Researcher (Generative AI)) stated: "A business's ability to sustain operations during adversarial attacks depends heavily on preemptive detection and rapid response frameworks."

- Participant 4 (Technology Education Expert), adding "Integrity of AI systems assures stakeholders of the validity of the systems for long-term sustainability in education."



Case study documents further supported this correlation:

- Awodiji (2022) argued that adversarial training models enhance the robustness of generative AI systems, thus maintaining their reliability in critical applications.
- Citron and Chesney (2019) underlined that traceability mechanisms, such as watermarking, increase transparency, enhancing user trust and promoting sustainability.

Evidence shows that cybersecurity strategies minimize not just threats but also provide a foundation in terms of good business practices. In addition, stakeholders will have more confidence in trusting such organizations that respect a very good cybersecurity posture.

Research Question 3: How do robust cybersecurity practices quantitatively impact stakeholder trust and long-term business viability in Nigerian businesses?

Stakeholder trust was a common theme in both interviews and case studies, with all six professionals highlighting its importance in their respective industries.

- Participant 5 (Educational Technology Specialist) said: "The trust of educators and students in AI-driven systems is directly tied to the systems' security and reliability. For example, the results we get from Turnitin sometimes get countered by the students claiming the academic task is carried out independently without the use of AI. So, a single data breach can erode years of trust."
- Participant 3 (Cybersecurity Consultant (Risk Assessment)) said: "In healthcare, patients are more likely to use AI diagnostic tools if they trust the data's integrity. Cybersecurity practices are critical in building this trust."

Case study findings supported this:

- Demir and Ciftci (2021) outlined that incident reporting fosters transparency and accountability, which are two elements of stakeholder trust.
- Awodiji (2022) illustrated that effective malware detection frameworks reduce the frequency of breaches and, therefore, increase stakeholder confidence in AI systems.

Findings show that entities that take cybersecurity seriously also have more trust from stakeholders, which translates into long-term viability. Sound practices guard against reputational damage due to breaches and ensure continuity of operations—critical factors in sustainability.

Key Themes in Relation to Research Questions

1. Detection Challenges

The key detection challenges based on our findings include the difficulty in adversarial attacks and the limitation of data. Both interviews and case studies demonstrated how traditional detection often cannot match the steps of evolving threats. This implies the need for businesses to embrace advanced detection tools, such as deep learning models, against these challenges.

2. Mitigation Strategies

The results showed that, indeed, a number of mitigation strategies—adversarial training, policy enforcement, and hybrid detection models—work in cybersecurity offense/defense. Such strategies are directly related to business sustainability and stakeholder trust, supported by both professional experience and case study documentation.

3. The Role of Deep Learning

Deep learning has become a double-edged sword: while it is a strong tool in the area of feature detection and classification, scalability and interpretability remain a challenge. This therefore aligns with the research questions since businesses should know how to navigate these limitations to gain full benefits from deep learning in cybersecurity.

CONCLUSION

The quantitative and qualitative findings obtained from interviews and case studies have so far been helpful in the in-depth understanding of the research questions. Challenges in detection identify deficits in the current systems, mitigation strategies, and deep learning potencies that could give fitting responses to these issues. This paper has therefore confirmed that robust cybersecurity practices, stakeholder's trust, and long-term business sustainability are strongly related, hence serving useful lessons for private businesses operating locally and other parts of the world.

RECOMMENDATIONS

Based on the results of this research and the induced theoretical framework consisting of Game Theory, Adversarial Training, and Explainable AI, the following recommendations are likely to effectively help in mitigating adversarial attacks on generative AI systems:



1. Embed Game Theory into the defensive strategy. The interaction between the generation of adversarial attackers and defensive mechanisms has to be modelled as a dynamic game. In such an "adversarial game," the strategy of the defending system has to be continuously optimized so as not to be outsmarted by the attackers (Neupane et al. 2024). Thus, an organization can use game theory models to catch up with the possible attack vectors and plan proactive mechanisms. For instance, the defending team can simulate a wide array of possible attack situations to forecast adversarial behaviours with measures prior to real vulnerabilities. According to Esposito et al. (2020), game theory can be applied to resource allocation by making very robust security protocols where the possibility of an attack is higher.
2. Adversarial Training for Robustness: Adversarial training should become an integral part of the development process concerning generative AI systems. The introduction of adversarial examples into any system at the time of training tunes it to recognize and nullify manipulated inputs. With each run of training, it will only get better at identifying minute perturbations that an attacker may attempt to exploit. In this regard, Dhamija and Bansal (2024) suggested that developers are also advised to maintain a repository of adversarial examples updated continuously so that the training is matched up with the latest threats. Hence, Patil et al. (2024) stated that organizations are expected to integrate adversarial training into the core of their machine learning pipeline when the application of generative AI touches high-risk areas such as finance and healthcare.
3. Explainable AI for Transparency and Trust: Application of Explainable AI needs to be considered for interpreting decisions made through deep learning models in adversarial threat detection to instill much-needed transparency and build trust among stakeholders, especially in application domains where decisions have to be accounted for. There is also the need for interpretability methods such as saliency maps or techniques like SHAP that help in visualizing how a model identifies adversarial inputs. These will basically let researchers identify weaknesses within the system. Moreover, explainability might conduce to regulatory compliance, showing how AI systems detect and mitigate threats amidst the "black-box" nature of deep learning models (Recker, 2021).
4. One major issue identified is that cybersecurity is very fragmented, hence, Alaeifar et al. (2024) standardization and knowledge-sharing frameworks for collaboration would provide a framework to standardize the detection tools for knowledge sharing on adversarial threats. This requires a concerted effort by industry stakeholders, academia, and policy thinkers who will develop open-source platforms for sharing adversarial examples, detection algorithms, and mitigation strategies. Standardized reporting protocols for adversarial attacks will ensure good interoperability among systems with extensive threat detection and response capabilities (Nespoli, Gomez Marmol and Maestre Vidal, 2021).
5. Hybrid Models for Detection: It is important, in improving the accuracy of detection, that organizations apply the use of hybrid models. Such hybrid models result from merging rule-based techniques with machine learning approaches. For instance, integrating traditional forensic methods with advanced deep learning models improves the system's capability to identify generative inconsistencies. Techniques such as integrating watermarking into generative AI systems further enhance traceability as opined by Hur et al. (2024). Hence easy attributions and addressing of adversarial attacks can be made.
6. Capacity Building through Education and Training: The education and training programs should aim at enhancing the capacity of AI practitioners, cybersecurity professionals, and policy makers to address adversarial threats (Jimmy, 2021). The training should include practical, workshop-style adversarial training, game theory applications, and the use of explainable AI tools. Such strengthening of human capital in these three areas will ensure that organizations remain resilient in light of adversarial attacks that are increasingly sophisticated. By integrating this six-phased recommendation and triumvirate theoretical framework (game theory, adversarial training, and explainable AI) into the defensive strategies, organizations can build robust, adaptive, and transparent systems that can counter adversarial threats in generative AI applications.



REFERENCES

- [1]. Alaeifar, P., Pal, S., Jadidi, Z., Hussain, M., & Foo, E. (2024). Current approaches and future directions for Cyber Threat Intelligence sharing: A survey. *Journal of Information Security and Applications*, 83, 103786.
- [2]. Al-Dosari, K., Fetais, N. and Kucukvar, M., 2024. Artificial intelligence and cyber defense system for banking industry: A qualitative study of AI applications and challenges. *Cybernetics and systems*, 55(2), pp.302-330.
- [3]. Andrew, L. (2020). The vulnerability of vital systems: how 'critical infrastructure' became a security problem. In *Securing the Homeland* (pp. 17-39). Routledge.
- [4]. Aspers, P., & Corte, U. (2019). What is qualitative in qualitative research. *Qualitative sociology*, 42, 139-160.
- [5]. Awodiji, T. O. (2021). Industrial big data analytics and cyber-physical systems for future maintenance & service innovation. *CS & IT Conference Proceedings*, 11(14). *CS & IT Conference Proceedings*.
- [6]. Awodiji, T. O. (2022). Malicious malware detection using machine learning perspectives. *Journal of Information Engineering and Applications*, 9-17.
- [7]. Bharadiya, J. P. (2023). A review of Bayesian machine learning principles, methods, and applications. *International Journal of Innovative Science and Research Technology*, 8(5), 2033-2038.
- [8]. Chander, B., John, C., Warriar, L., & Gopalakrishnan, K. (2024). Toward trustworthy artificial intelligence (TAI) in the context of explainability and robustness. *ACM Computing Surveys*.
- [9]. Chen, H., Pendleton, M., Njilla, L., & Xu, S. (2020). A survey on ethereum systems security: Vulnerabilities, attacks, and defenses. *ACM Computing Surveys (CSUR)*, 53(3), 1-43.
- [10]. Chesney, B., & Citron, D. (2019). Deep fakes: A looming challenge for privacy, democracy, and national security. *California Law Review*, 107, 1753.
- [11]. Corallo, A., Lazoi, M., & Lezzi, M. (2020). Cybersecurity in the context of industry 4.0: A structured classification of critical assets and business impacts. *Computers in industry*, 114, 103165.
- [12]. Demir, I., & Ciftci, U. A. (2021, May). Where do deep fakes look? synthetic face detection via gaze tracking. *ACM symposium on eye tracking research and applications*, pp. 1-11.
- [13]. Dhamija, L., & Bansal, U. (2024). How to Defend and Secure Deep Learning Models Against Adversarial Attacks in Computer Vision: A Systematic Review. *New Generation Computing*, 1-71.
- [14]. Edmund, E. (2024). Risk Based Security Models for Veteran Owned Small Businesses. *International Journal of Research Publication and Reviews*, 5(12), 4304-4318.
- [15]. Esposito, C., Tamburis, O., Su, X., & Choi, C. (2020). Robust decentralised trust management for the internet of things by using game theory. *Information Processing & Management*, 57(6), 102308.
- [16]. Fonseca, J., & Bacao, F. (2023). Tabular and latent space synthetic data generation: a literature review. *Journal of Big Data*, 10(1), 115.
- [17]. Górriz, J. M., Álvarez-Illán, I., Álvarez-Marquina, A., Arco, J. E., Atzmueller, M., Ballarini, F., ... & Ferrández-Vicente, J. M. (2023). Computational approaches to explainable artificial intelligence: advances in theory, applications and trends. *Information Fusion*, 100, 101945.
- [18]. Halvorsen, J., Izurieta, C., Cai, H. and Gebremedhin, A. (2024) Applying generative machine learning to intrusion detection: A systematic mapping study and review. *ACM Computing Surveys*, 56(10), pp.1-33.
- [19]. Han, X., Zhang, Y., Wang, W., & Wang, B. (2022). Text adversarial attacks and defenses: Issues, taxonomy, and perspectives. *Security and Communication Networks*, 22(1), 6458488.
- [20]. Hoang, V. T., Ergu, Y. A., Nguyen, V. L., & Chang, R. G. (2024). Security risks and countermeasures of adversarial attacks on AI-driven applications in 6G networks: A survey. *Journal of Network and Computer Applications*, 104031.
- [21]. Hur, H., Kang, M., Seo, S., & Hou, J. U. (2024). Latent Diffusion Models for Image Watermarking: A Review of Recent Trends and Future Directions. *Electronics*, 14(1), 25.
- [22]. Ingle, G., & Pawale, S. (2024). Enhancing Adversarial Defense in Neural Networks by Combining Feature Masking and Gradient Manipulation on the MNIST Dataset. *International Journal of Advanced Computer Science & Applications*, 15(1).
- [23]. Jimmy, F. (2021). Emerging threats: The latest cybersecurity risks and the role of artificial intelligence in enhancing cybersecurity defenses. *Valley International Journal Digital Library*, 564-574.
- [24]. Kang, H., & Kang, P. (2024). Transformer-based multivariate time series anomaly detection using inter-variable attention mechanism. *Knowledge-Based Systems*, 290, 111507.
- [25]. Kaplan, D. (2021). On the quantification of model uncertainty: A Bayesian perspective. *Psychometrika*, 86(1), 215-238.
- [26]. Khan, M., & Ghafoor, L. (2024). Adversarial Machine Learning in the Context of Network Security: Challenges and Solutions. *Journal of Computational Intelligence and Robotics*, 4(1), 51-63.
- [27]. Kiesler, N. (2020). On programming competence and its classification. *Proceedings of the 20th Koli Calling International Conference on Computing Education Research* (pp. 1-10).



- [28]. Krishnamurthy, O. (2024). Impact of Generative AI in Cybersecurity and Privacy. *International Journal of Advances in Engineering Research*, 27(1), 26-38.
- [29]. Leszczyna, R. (2021). Review of cybersecurity assessment methods: Applicability perspective. *Computers & Security*, 108, 102376.
- [30]. Li, Z., Wang, Y., Zhang, N., Zhang, Y., Zhao, Z., Xu, D., Ben, G., & Gao, Y. (2022). Deep Learning-Based Object Detection Techniques for Remote Sensing Images: A Survey. *Remote Sensing*, 14(10), 2385. <https://doi.org/10.3390/rs14102385>
- [31]. Lone, A. N., Mustajab, S., & Alam, M. (2023). A comprehensive study on cybersecurity challenges and opportunities in the IoT world. *Security and Privacy*, 6(6), e318.
- [32]. Macas, M., Wu, C., & Fuertes, W. (2022). A survey on deep learning for cybersecurity: Progress, challenges, and opportunities. *Computer Networks*, 212, 109032.
- [33]. Machado, G. R., Silva, E., & Goldschmidt, R. R. (2021). Adversarial machine learning in image classification: A survey toward the defender's perspective. *ACM Computing Surveys (CSUR)*, 55(1), 1-38.
- [34]. Mavikumbure, H. S., Cobilean, V., Wickramasinghe, C. S., Drake, D., & Manic, M. (2024). Generative AI in cyber security of cyber physical systems: Benefits and threats. *International Conference on Human System Interaction (HSI)* (pp. 1-8). IEEE.
- [35]. Mitra, A., Mohanty, S. P., Corcoran, P., & Kougianos, E. (2021). A machine learning based approach for deepfake detection in social media through key video frame extraction. *SN Computer Science*, 2(2), 98.
- [36]. Nain, G., Pattanaik, K. K., & Sharma, G. K. (2022). Towards edge computing in intelligent manufacturing: Past, present and future. *Journal of Manufacturing Systems*, 62, 588-611.
- [37]. Nespoli, P., Gomez Marmol, F., & Maestre Vidal, J. (2021). Battling against cyberattacks: Towards pre-standardization of countermeasures. *Cluster Computing*, 24, 57-81.
- [38]. Neupane, R. L., Bhusal, B., Neupane, K., Regmi, P., Dinh, T., Marrero, L., & Calyam, P. (2024). On Countering Ransomware Attacks Using Strategic Deception. *International Conference on Decision and Game Theory for Security* (pp. 149-176). Cham: Springer Nature Switzerland.
- [39]. Patil, D., Rane, N. L., Desai, P., & Rane, J. (Eds.). (2024). *Trustworthy Artificial Intelligence in Industry and Society*. Deep Science Publishing.
- [40]. Price, H. E., & Smith, C. (2021). Procedures for reliable cultural model analysis using semi-structured interviews. *Field Methods*, 33(2), 185-201.
- [41]. Recker, J. (2021). *Scientific research in information systems: a beginner's guide*. Springer Nature.
- [42]. Qiu, S., Liu, Q., Zhou, S., & Huang, W. (2022). Adversarial attack and defense technologies in natural language processing: A survey. *Neurocomputing*, 492, 278-307.
- [43]. Rane, N., Choudhary, S. P., & Rane, J. (2024). Ensemble deep learning and machine learning: applications, opportunities, challenges, and future directions. *Studies in Medical and Health Sciences*, 1(2), 18-41.
- [44]. Salmona, M., & Kaczynski, D. (2024). Qualitative data analysis strategies. In *How to conduct qualitative research in finance* (pp. 80-96). Edward Elgar Publishing.
- [45]. Sarker, I. H. (2024). Generative AI and Large Language Modeling in Cybersecurity. In *AI-Driven Cybersecurity and Threat Intelligence: Cyber Automation, Intelligent Decision-Making and Explainability* (pp. 79-99). Cham: Springer Nature Switzerland.
- [46]. Serban, A., Poll, E., & Visser, J. (2020). Adversarial examples on object recognition: A comprehensive survey. *ACM Computing Surveys (CSUR)*, 53(3), 1-38.
- [47]. Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1), 1-48.
- [48]. Sinha, H. (2024). The identification of network intrusions with generative artificial intelligence approach for cybersecurity. *Journal of Web Applications and Cyber Security*, 2(2), 20-29.
- [49]. Tseng, H. Y., Jiang, L., Liu, C., Yang, M. H., & Yang, W. (2021). Regularizing generative adversarial networks under limited data. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 7921-7931).
- [50]. Vergara Cobos, Estefania Cakir, & Selcen. (2024). *A Review of the Economic Costs of Cyber Incidents*. Washington, DC: World Bank.
- [51]. Wang, J., Wang, Y., & Liu, H. (2022). Hybrid Variability Aware Network (HVANet): A self-supervised deep framework for label-free SAR image change detection. *Remote Sensing*, 14(3), 734.
- [52]. Wang, Y., Sun, T., Li, S., Yuan, X., Ni, W., Hossain, E., & Poor, H. V. (2023). Adversarial attacks and defenses in machine learning-empowered communication systems and networks: A contemporary survey. *IEEE Communications Surveys & Tutorials*.
- [53]. Wang, X., Zhao, Y., Qiu, C., Hu, Q., & Leung, V. C. (2024). Socialized learning: A survey of the paradigm shift for edge intelligence in networked systems. *IEEE Communications Surveys & Tutorials*.



- [54]. Yumlebam, R., Issac, B., Jacob, S. M., & Yang, L. (2024). Comprehensive botnet detection by mitigating adversarial attacks, navigating the subtleties of perturbation distances and fortifying predictions with conformal layers. *Information Fusion*, 102529.
- [55]. Zamry, N. M., Zainal, A., Rassam, M. A., Alkhamash, E. H., Ghaleb, F. A., & Saeed, F. (2021). Lightweight anomaly detection scheme using incremental principal component analysis and support vector machine. *Sensors*, 21(23), 8017.
- [56]. Zeng, L., Qiu, D., & Sun, M. (2022). Resilience enhancement of multi-agent reinforcement learning-based demand response against adversarial attacks. *Applied Energy*, 324, 119688.
- [57]. Zhang, J., Zhu, H., Wang, F., Zhao, J., Xu, Q., & Li, H. (2022). Security and privacy threats to federated learning: Issues, methods, and challenges. *Security and Communication Networks*, 2022(1), 2886795.