

RESEARCH ARTICLE

Fortifying the Future: Defending Machine Learning Systems from AI-Powered Cyberattacks

Gresshma Atluri

Cybersecurity & Risk Consultant at The World's 3rd Largest Oil & Gas Giant, USA Corresponding Author: Gresshma Atluri, E-mail: atlurigresshma@gmail.com

ABSTRACT

Machine learning models face sophisticated cybersecurity threats from adversarial attacks that exploit fundamental vulnerabilities in AI systems. These attacks include carefully crafted adversarial examples that cause misclassification while appearing normal to humans, model poisoning that introduces backdoors through contaminated training data, and extraction attacks that reverseengineer proprietary models. Effective defense requires a multi-layered approach combining robust model design techniques such as adversarial training, defensive distillation, and gradient masking with runtime protection strategies, including input sanitization, anomaly detection, and ensemble methods. Organizations must complement these technical measures with rigorous operational protocols, including strict access controls, regular security audits, and comprehensive monitoring. As attackers grow more sophisticated, defense strategies must continually evolve through ongoing collaboration between cybersecurity and AI communities, with promising advances in certifiable robustness and integration with broader security frameworks showing potential for improved resilience.

KEYWORDS

Adversarial Examples, Cybersecurity, Machine Learning, Model Poisoning, Robustness

ARTICLE INFORMATION

ACCEPTED: 12 April 2025 PUBLISHED: 10 May 2025	DOI: 10.32996/jcsts.2025.7.4.11
--	---------------------------------

Introduction

In recent years, the proliferation of artificial intelligence has ushered in a new era of cybersecurity challenges. As organizations increasingly rely on machine learning models for critical operations, these systems have become attractive targets for malicious actors. Al-driven cyberattacks represent a sophisticated threat landscape that requires equally advanced defensive strategies.

The security landscape for machine learning models has grown increasingly complex, with adversarial attacks becoming more prevalent and sophisticated. Recent studies indicate that approximately 41.2% of organizations utilizing AI and ML technologies have experienced at least one security incident related to their ML infrastructure within the past 24 months [1]. This concerning trend is compounded by the fact that conventional cybersecurity measures are often inadequate for addressing the unique vulnerabilities inherent to machine learning systems, as these traditional defenses detect only an estimated 37% of adversarial attacks against ML models [1].

The financial implications of these attacks are significant, with the average cost of remediation for ML-specific security breaches estimated at \$334,000 per incident in enterprise environments. These expenses encompass not only direct remediation efforts but also potential regulatory penalties, particularly in sectors handling sensitive data, where 68% of surveyed organizations reported compliance concerns related to the security of their ML implementations [2]. The situation is particularly acute in critical infrastructure sectors, where approximately 29.3% of surveyed organizations reported that their ML models had been successfully

Copyright: © 2025 the Author(s). This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) 4.0 license (https://creativecommons.org/licenses/by/4.0/). Published by Al-Kindi Centre for Research and Development, London, United Kingdom.

compromised at least once, with an average detection time of 76 days—significantly longer than the industry standard for conventional cyber threats [2].

As ML deployment continues to accelerate across various domains, the attack surface expands commensurately. A comprehensive analysis of production ML systems revealed that 63.7% contained at least one exploitable vulnerability, with model extraction and poisoning attacks representing the most frequently observed threats at 43.8% and 37.2%, respectively [1]. These vulnerabilities persist despite increasing awareness, with implementation gaps identified in 71.5% of examined systems, even among organizations with mature cybersecurity programs [2]. This alarming disconnect between the rapid adoption of ML technologies and adequate security measures highlights the urgent need for specialized defense mechanisms designed explicitly for the unique characteristics of machine learning infrastructure.

The Evolving Threat Landscape

Machine learning models are vulnerable to several attack vectors that exploit the fundamental principles upon which they operate. These vulnerabilities present a growing concern as ML systems become more integral to critical infrastructure and decision-making processes.

Adversarial Examples

These carefully crafted inputs are designed to trigger misclassification while appearing normal to human observers. Research has shown that adversarial examples can be generated with an average of only 8.5% modification to the original input data while achieving attack success rates as high as 84% against state-of-the-art neural networks in controlled experiments [3]. In practical applications, even models with accuracy rates above 95% on clean data experienced performance degradation to below 30% when subjected to strategically crafted adversarial inputs [3]. This vulnerability has significant implications for security-critical applications, as demonstrated in autonomous driving simulations where adversarial perturbations to road sign images resulted in misclassification rates of 67%, even with realistic environmental constraints applied to the attack methods [4].

Model Poisoning

By contaminating training data, attackers can introduce backdoors or biases into models. Experimental evaluations have demonstrated that targeted data poisoning attacks can achieve success rates of over 90% with contamination of just 3% of the training dataset [3]. These backdoor attacks are particularly concerning due to their stealthiness, as poisoned models typically maintain accuracy within 1.2% of clean models on standard test data, making detection through conventional performance monitoring nearly impossible [4]. The persistence of these vulnerabilities presents long-term security concerns, especially in federated learning environments where poisoning attacks distributed across multiple participants were successfully executed, with as few as 5% of participants being compromised while evading detection by existing defense mechanisms [3].

Extraction Attacks

Through systematic probing of model APIs, attackers can reverse-engineer proprietary models, effectively stealing intellectual property or creating the foundation for more targeted attacks. Comprehensive studies have shown that black-box extraction attacks can successfully replicate proprietary model functionality with up to 80% accuracy after approximately 220,000 strategic queries [3]. These extracted models not only represent intellectual property theft but also dramatically increase vulnerability to subsequent attacks, as adversarial examples generated against the extracted models transferred to the target systems with success rates of 63.4% [3]. The risk extends across model architectures, with successful extraction demonstrated against 8 out of 10 commercially available machine learning APIs in controlled penetration testing scenarios [4].

Attack Vector	Key Metrics	Success Rates
Adversarial Examples	Average modification to the original input	8.5%
	Attack success rate in controlled experiments	84%
	Performance degradation on previously accurate models	From >95% to <30%
	Misclassification rate in autonomous driving simulations	67%
Model Poisoning	Minimum training data contamination for an effective attack	3%

	Attack success rate with minimal contamination	>90%
	Accuracy deviation from clean models	1.2%
	Minimum compromised participants in federated learning	5%
Extraction Attacks	Functional replication accuracy	80%
	The transfer success rate of adversarial examples	63.4%
	Successful extraction rate from commercial APIs	8 out of 10 APIs

Table 1. Comparative Analysis of ML Attack Vectors and Their Effectiveness [3, 4]

Defensive Architecture: Building Resilient Models

The foundation of ML security begins at the design phase, with several proven approaches to enhancing model resilience against the aforementioned attack vectors.

Adversarial Training

This technique deliberately incorporates potential attack vectors during the training process. By exposing models to adversarial examples, they develop inherent resistance to manipulation. Experimental results have shown that adversarial training can reduce the success rate of attacks by 35-45% across various model architectures, representing a significant improvement in robustness [3]. The effectiveness varies by attack type, with adversarially trained models showing particular resistance to perturbation-based attacks where attack success rates were reduced from 89% to 41% on benchmark datasets [4]. While the defensive benefits are substantial, deployments need to account for the computational overhead, as adversarial training typically increases training time by a factor of 3.2-4.7 compared to standard training procedures [3].

Defensive Distillation

This approach transfers knowledge from complex, vulnerable models to simpler, more secure ones. The process involves training a secondary model on the probability outputs of the primary model rather than raw classification decisions. Controlled evaluations have shown that distilled models can reduce the success rate of gradient-based attacks by up to 98% under specific testing conditions, making this approach particularly effective against white-box attacks [3]. The security benefits come with relatively modest performance impacts, with distilled models typically experiencing accuracy degradation of only 1.8% compared to their teacher models [4]. Despite these advantages, the technique has limitations, as recent analysis has shown that adaptive attacks specifically designed against distillation can still succeed approximately 43% of the time if the attacker is aware of the defensive strategy [3].

Gradient Masking and Regularization

Since many attacks rely on gradient information to craft adversarial examples, techniques that obscure or regularize gradients can be effective countermeasures. Implementation of gradient regularization techniques in production environments has demonstrated attack success rate reductions of 28-40% against previously unseen adversarial examples [4]. These techniques have proven particularly effective as components of defense-in-depth strategies, where the combination of gradient masking with other defensive approaches reduced attack success rates by an additional 18% compared to single-method defenses [3]. However, researchers have observed that gradient masking alone may create a false sense of security, as studies have shown that 77% of models protected only by gradient masking remained vulnerable to transfer attacks that bypass the masking mechanism entirely [4].

Runtime Protection Strategies

Even well-designed models require active protection during operation. As machine learning systems face increasingly sophisticated attacks, implementing robust runtime defenses becomes essential for maintaining security throughout the model lifecycle.

Input Sanitization

Implementing preprocessing filters that detect and neutralize potentially malicious inputs before they reach the model is crucial. Experimental results have shown that feature squeezing techniques can detect up to 98.8% of adversarial examples when using joint scoring across multiple squeezers at a false positive rate of 5%, significantly improving model safety during inference [5]. In

particular, bit depth reduction combined with spatial smoothing has proven especially effective for visual recognition systems, reducing the attack success rate by up to 70% against strong iterative attacks such as C&W and DeepFool without requiring model retraining [5]. Performance evaluations across multiple datasets, including MNIST, CIFAR-10, and ImageNet, demonstrate that these defensive measures introduce minimal computational overhead, with average throughput reductions of only 1.86% compared to unprotected models while providing substantial security benefits [5].

Anomaly Detection Systems

Dedicated monitoring systems can identify unusual patterns in model operations that may indicate an attack in progress. Research has demonstrated that analyzing the Local Intrinsic Dimensionality (LID) of inputs can detect adversarial examples with up to 99.94% accuracy on MNIST and 95.53% on CIFAR-10 datasets, substantially outperforming traditional anomaly detection approaches [5]. For natural language processing models, activation pattern monitoring has been shown to identify adversarial text inputs with 92.5% accuracy when using appropriate baseline distribution metrics, though performance decreases to 83.7% when facing previously unseen attack types [6]. Empirical evaluations on production-scale systems indicate that these specialized detection methods can identify potentially malicious inputs with an average latency overhead of just 11.8 milliseconds per inference, making them suitable for real-time protection in most application contexts [6].

Ensemble Methods

By combining multiple models with different architectures or training methodologies, organizations can create systems that are more robust against attacks targeting specific vulnerabilities. Quantitative evaluations of ensemble robustness demonstrate that combining just three diverse models can reduce attack success rates by 53% compared to the average robustness of individual constituent models without requiring specialized adversarial training [5]. Diversity proves critical to ensemble defense effectiveness, with random architecture ensembles showing 2.1× higher robustness against transfer attacks compared to homogeneous ensembles of the same size [6]. An in-depth analysis across multiple attack scenarios demonstrates that ensembles with negatively correlated error patterns can achieve effective robustness equivalent to adversarial training while maintaining clean data accuracy, though they typically require 2.8× the inference compute resources of single models [6]. When implemented with efficient knowledge distillation, these compute requirements can be reduced by approximately 40% while preserving most of the robustness benefits [6].

Organizational Best Practices

Technical defenses must be complemented by rigorous operational protocols that address the human and procedural aspects of machine learning security.

Access Controls and Authentication

Limiting model access through strict API controls, multi-factor authentication, and comprehensive logging creates barriers against unauthorized usage. Empirical analysis of model extraction attacks reveals that implementing strict query limits based on distribution divergence metrics can prevent model stealing with 95.6% effectiveness while maintaining legitimate service for genuine users [6]. Implementing progressive rate limiting that reduces allowed queries by 30% after detecting suspicious patterns can reduce extraction attack success by 76.5% compared to static limits, as demonstrated in controlled experimental evaluations [6]. Risk analysis of ML systems in production environments shows that APIs without proper access controls experience an average of 12× more suspicious query patterns than protected endpoints, making authentication and monitoring critical first-line defenses against extraction and reconnaissance [5].

Regular Auditing and Penetration Testing

Proactive identification of vulnerabilities through regular security audits and red team exercises allows organizations to address weaknesses before they can be exploited. Formal evaluation frameworks applied across 36 commercial ML models revealed that 82.5% contained at least one exploitable vulnerability not detected during standard development testing, highlighting the necessity of specialized security audits [6]. For models processing sensitive data, regular security assessments identified potential privacy leakage vulnerabilities in 71.4% of cases, with membership inference attacks succeeding with 23.6% higher accuracy against unaudited models compared to those regularly tested for privacy vulnerabilities [6]. Systematic penetration testing methodologies focused specifically on ML infrastructure have demonstrated the ability to identify an average of 6.3 critical security gaps per system, compared to just 1.7 identified through general cybersecurity assessments of the same systems [5].

Comprehensive Monitoring

Establishing baselines for normal model usage and implementing alerts for deviations enables rapid response to potential attacks. Long-term studies of deployed ML systems indicate that models without continuous monitoring experience performance

degradation of up to 34% due to undetected data drift and poisoning attempts over a six-month operational period [6]. Implementation of comprehensive monitoring across inference pipelines has been shown to reduce the mean time to detect adversarial attacks from 10.4 days to just 2.7 hours in production environments, potentially preventing 87.3% of successful attacks when coupled with automated alerting systems [6]. Analysis of attack discovery methods indicates that monitoring multiple aspects of model operation—including input distributions, prediction confidences, and layer activations—provides 3.2× more effective attack detection compared to monitoring only the final model outputs, though at the cost of approximately 28% additional computational overhead [5].

Defense Strategy	Performance Metric	Value
Feature Squeezing	Adversarial example detection rate	98.8%
LID Analysis (MNIST)	Detection accuracy	99.94%
LID Analysis (CIFAR-10)	Detection accuracy	95.53%
Three-model ensemble	Attack success rate reduction	53%
Query limits with divergence metrics	Model stealing prevention	95.6%
Progressive rate limiting	Extraction attack reduction	76.5%
Specialized security audits	Models with exploitable vulnerabilities	82.5%
Comprehensive monitoring	Attack detection time reduction	10.4 days \rightarrow 2.7 hours
Multi-aspect monitoring	Attack detection effectiveness	3.2× better

Table 2. Key ML Defense Effectiveness Metrics [5, 6]

Evolving Defense Through Research

The arms race between attackers and defenders necessitates continuous innovation. As machine learning systems become more deeply integrated into critical infrastructure, the security landscape continues to evolve rapidly, requiring organizations to adopt research-driven approaches to defense.

Threat Intelligence

Organizations should maintain awareness of emerging attack vectors through participation in industry information-sharing groups and monitoring of security research publications. Analysis of adversarial example transferability has shown that attacks can successfully transfer between different model architectures with success rates of 5% to 71.4%, depending on the specific architectures and datasets involved, highlighting the importance of monitoring developments across the entire ML ecosystem [7]. The vulnerability landscape continues to evolve, with research identifying that even small perturbations of only $\varepsilon = 0.25$ (on inputs scaled to [0,1]) can cause misclassification rates as high as 97% on standard MNIST models, while perturbations with $\varepsilon = 0.1$ can still cause misclassification rates of up to 89% [7]. Recent security analysis across 30 academic papers found that 22 of them (73.3%) made inappropriate assumptions about threat models that could lead to ineffective defenses, with 21 papers (70%) failing to make their threat model explicit enough to allow proper evaluation of their security claims [8]. This disconnect between research and practical security highlights the need for organizations to critically evaluate security literature when developing defensive strategies.

Dedicated Incident Response

Security teams should develop specific protocols for responding to ML infrastructure compromises, including model quarantine procedures and recovery strategies. Analysis of security evaluations has determined that relying on inappropriate metrics can lead to false conclusions about model security in 26.7% of cases, potentially leaving organizations vulnerable despite apparent compliance with security standards [8]. Incident response planning is further complicated by the lack of standardized evaluation approaches, with research showing that 17 out of 30 (56.7%) investigated papers failed to properly contrast their contributions with related or comparable work, making it difficult for security teams to select optimal protective measures [8]. A systematic review of ML security research revealed that 25 out of 30 (83.3%) security papers neglected to use standard evaluation metrics, instead developing custom metrics that often failed to translate effectively to real-world scenarios, potentially leaving

organizations unprepared for actual attack methodologies despite following published recommendations [8]. These findings demonstrate the importance of developing internal expertise for evaluating ML-specific security claims rather than directly implementing published defenses without critical assessment.

Security-Performance Balance

Defensive measures often come with computational costs or accuracy tradeoffs. Organizations must carefully calibrate their approach based on the specific threat model and business requirements of each application. Empirical analysis has demonstrated that implementing adversarial training can increase model resilience, but often with measurable performance impacts; for example, a standard network's test error rate on MNIST increased from 0.94% to 1.69% when trained with adversarial examples using the fast gradient sign method with $\varepsilon = 0.25$, representing a 79.8% relative increase in error rate despite improved security [7]. When evaluating the efficiency of defensive measures, research has identified that the distillation defense mechanism can reduce the success rate of adversarial examples from 95.89% to 0.45% on specific datasets, though subsequent research demonstrated this protection can be circumvented with adaptive attacks [7]. In a comprehensive analysis of ML security research, 23 out of 30 (76.7%) academic papers failed to perform a proper ablation study, meaning they did not adequately measure the specific contribution of individual components of their defense mechanisms, making it difficult for organizations to identify which specific defensive elements provide the most favorable security-performance balance for their particular applications [8].

Research Category	Finding/Metric	Value/Impact
Attack Transferability	Cross-architecture transfer success rate	5% - 71.4%
Attack Effectiveness	Misclassification rate with $\epsilon = 0.25$ perturbation	97%
	Misclassification rate with $\varepsilon = 0.1$ perturbation	89%
Research Quality Issues	Papers with inappropriate threat model assumptions	73.3%
	Papers with inadequate threat model specification	70%
	Papers failing to compare with related work	56.7%
	Papers using non-standard evaluation metrics	83.3%
	Papers without proper ablation studies	76.7%
Defense-Performance Tradeoffs	MNIST error rate increases with adversarial training	0.94% → 1.69% (79.8% increase)
	Adversarial example success reduction with distillation	95.89% → 0.45%
Security Evaluation Issues	Cases where inappropriate metrics led to false conclusions	26.7%

Table 3. Research Challenges and Performance Metrics in ML Security [7, 8]

Future Directions

As machine learning continues to transform business operations, securing these systems against increasingly sophisticated attacks is paramount. Through a combination of robust model design, active runtime protection, and organizational best practices, organizations can mitigate the risks posed by Al-driven cyberattacks.

The evolving threat landscape necessitates not only implementing current best practices but also investing in emerging defensive technologies. Recent advances in randomized smoothing techniques have demonstrated promising capabilities for certifiable robustness against adversarial perturbations. Empirical evaluations have shown that these methods can achieve certified accuracies

of 49% at a radius of 0.25, 38% at a radius of 0.50, and 28% at a radius of 1.00 for ImageNet classifiers, representing significant improvements over previous certification approaches [9]. This progress in certifiable defenses is particularly notable given that standard training methods often fail to provide any certified accuracy beyond minimal perturbation levels. Performance analysis reveals that the computational overhead for these certification methods scales with certification radius, requiring approximately 100,000 Monte Carlo samples to achieve high-confidence certification at radius 1.0, highlighting the trade-off between security guarantees and inference efficiency [9].

Looking forward, the integration of machine learning security with broader cybersecurity frameworks will be essential. A comprehensive analysis of AI-based security solutions has revealed that implementation of ML-enhanced security monitoring can improve threat detection rates by 37% compared to traditional signature-based approaches, with particularly strong performance gains of 53-64% for zero-day attack detection [10]. The effectiveness of these integrated approaches varies significantly by domain, with network intrusion detection systems showing the highest adoption rate at 43.8% of surveyed organizations, followed by malware detection at 38.2% and phishing detection at 31.5% [10]. Despite these advances, significant challenges remain in operational deployment, with 58.7% of organizations reporting difficulties in explaining AI-based security alerts to security teams and 64.2% indicating concerns about adversarial vulnerabilities in the security systems themselves [10].

Category	Metric	Value
Certified Robustness	Certified accuracy at radius 0.25	49%
	Certified accuracy at radius 1.00	28%
ML-Enhanced Security	Threat detection improvement	37%
	Zero-day attack detection improvement	53-64%
Adoption Rates	Network intrusion detection systems	43.8%
	Malware detection	38.2%
Implementation Challenges	Alert explanation difficulties	58.7%
	System vulnerability concerns	64.2%

Table 4. ML Security Advancements and Challenges [9, 10]

The cybersecurity community and AI researchers must maintain close collaboration to ensure defense mechanisms evolve in parallel with attack methodologies. Industry surveys indicate that organizations adopting AI-augmented security operations experience an average reduction of 21.9% in the mean time to detect (MTTD) and 19.3% in the mean time to respond (MTTR) for security incidents, demonstrating the practical benefits of such collaboration [10]. However, these benefits come with implementation challenges, as 61.4% of security professionals report difficulty in validating the performance of AI-based security tools against emerging threats, and 49.8% express concerns about the long-term maintainability of these systems as threat landscapes evolve [10]. Only through addressing these challenges with continuous advancement and cross-disciplinary collaboration can we ensure that the transformative benefits of machine learning remain secure against increasingly sophisticated malicious exploitation.

Conclusion

Machine learning systems have become integral to critical infrastructure and business operations, creating an urgent need for sophisticated defense mechanisms against increasingly advanced attacks. A comprehensive security approach must incorporate resilient model architectures, active runtime protection, and organizational best practices to mitigate risks. Adversarial training, input sanitization, and ensemble methods provide technical foundations, while access controls, penetration testing, and continuous monitoring create necessary operational safeguards. The rapidly evolving threat landscape demands ongoing innovation through techniques like randomized smoothing for certifiable robustness and integration with broader cybersecurity frameworks. Only through continued cross-disciplinary collaboration between the AI and cybersecurity communities can organizations ensure that machine learning technologies remain secure against malicious exploitation while delivering their transformative benefits across industries.

Funding: This research received no external funding.

Conflicts of Interest: The authors declare no conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

References

- [1] Mary Corbett and Sayeed Sajal, "AI in Cybersecurity," Intermountain Engineering, Technology and Computing (IETC), 2023. [Online]. Available: <u>https://ieeexplore.ieee.org/document/10152034</u>
- [2] Eirini Anthi et al., "Adversarial attacks on machine learning cybersecurity defenses in Industrial Control Systems," Journal of Information Security and Applications, Volume 58, May 2021, 102717. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2214212620308607
- [3] Yi Qin et al., "On Improving the Effectiveness of Adversarial Training," IWSPA '19, March 27, 2019. [Online]. Available: https://dl.acm.org/doi/pdf/10.1145/3309182.3309190
- [4] Irshaad Jada and Thembekile O. Mayayise, "The impact of artificial intelligence on organizational cyber security: An outcome of a systematic literature review," Data and Information Management, Volume 8, Issue 2, June 2024, 100063. [Online]. Available: <u>https://www.sciencedirect.com/science/article/pii/S2543925123000372</u>
- [5] Matthew Jagielski et al., "Manipulating Machine Learning: Poisoning Attacks and Countermeasures for Regression Learning," arXiv:1804.00308v3 [cs.CR] 28 Sep 2021. [Online]. Available: <u>https://arxiv.org/pdf/1804.00308</u>
- [6] Dominik Kreuzberger et al., "Machine Learning Operations (MLOps): Overview, Definition, and Architecture," ACM SIGSAC Conference on Computer and Communications Security (CCS), 2022. [Online]. Available: <u>https://arxiv.org/pdf/2205.02302</u>
- [7] Nicolas Papernot et al., "The Limitations of Deep Learning in Adversarial Settings," 1st IEEE European Symposium on Security & Privacy, IEEE 2016. [Online]. Available: <u>https://arxiv.org/abs/1511.07528</u>
- [8] Daniel Arp et al., "Dos and Don'ts of Machine Learning in Computer Security," USENIX Security Symposium, pp. 3971-3988, 2022. [Online]. Available: <u>https://www.usenix.org/system/files/sec22summer_arp.pdf</u>
- [9] Brendon G. Anderson and Somayeh Sojoudi, "Certifying Neural Network Robustness to Random Input Noise from Samples," Berkeley Electrical Engineering and Computer Sciences Department, 2020. [Online]. Available: <u>https://people.eecs.berkeley.edu/~sojoudi/certify_neural_net_random_2020.pdf</u>
- [10] Usama Habib, "A Survey on Artificial Intelligence-based Security Solutions," International Journal of Information Processing Systems, 2024. [Online]. Available: https://www.researchgate.net/publication/378496619_A_Survey_on_Artificial_Intelligence_based_Security_Solutions