

**SECURE AND EXPLAINABLE AI FOR REMOTE PATIENT MONITORING: A
DEFENSE FRAMEWORK AGAINST EMERGING AI-DRIVEN
HEALTHCARE CYBERCRIMES**

Mohammed Liaqat Ali Khan^{1*}, Dr. Priya Vij², Dr. Abdul Bari³

¹Research Scholar, Kalinga University, Raipur, India.

alykhan113@yahoo.co.in

²Assistant Professor, Department of CS & IT, Kalinga University, Raipur, India.

ku.priyavij@kalingauniversity.ac.in,

³Dean Academics & Prof. CSE Department, ISL Engineering College, Hyderabad, India.

abdulbarimohammed11@gmail.com

*Corresponding Author

Abstract

Remote Patient Monitoring (RPM) systems have also developed at a very fast pace having adopted artificial intelligence (AI), which can be used to carry out continuous clinical monitoring, the early detection of anomalies, and the proactive decision-making. Nonetheless, the identical AI-based capabilities have generated new security gaps that can be used to perpetrate specific cybercrimes, data manipulation, and clinical interference. The paper, based on a common defence structure, proposes adversarial threat modelling, multimodal anomaly detection, and explainable artificial intelligence (XAI) to reduce the emergent AI-based attacks in the RPM ecosystems. The model was evaluated on a hybrid dataset that consisted of simulated adversarial samples, real RPM telemetry, and synthetic physiological signals. To provide integrity, transparency, and robustness, the proposed architecture combines a secure AI pipeline with the concepts of differential privacy, federated learning, and audit trails supported by blockchain. The findings show a substantial decrease in the success rates of adversarial attacks (up to 87 percent), a higher sensitivity of anomaly detection (93.2 percent), and a better interpretability capacity toward clinical decision-making. The results highlight the need to implement security-by-design, interpretable AI models, and adaptive threat intelligence in the future to address the changing cybercrimes that are AI-enabled. The study offers practical recommendations in the areas of safe RPM creation, healthcare cybersecurity policies and regulations.

Keywords: Remote Patient Monitoring (RPM); Healthcare Cybersecurity; AI-driven Cybercrimes; Adversarial Attacks; Explainable AI (XAI); Secure AI Framework

1. INTRODUCTION

Remote Patient Monitoring (RPM) systems are swiftly changing the digital healthcare world by providing round-the-clock real-time monitoring of patient vitality using interconnected medical IoT devices and AI-based diagnostics engines. Incorporation of smart modelling with

telemedicine and sensor networks can improve the clinical decision making process, minimize the hospitalization burden and the proactive management of severe health conditions. With AI being more integrated into the healthcare fabric, the threat surface grows to match the threat, leaving RPM systems vulnerable to cybersecurity dangers like never before to adversarial machine learning, spoofed physiological signals, data manipulation, and autonomous malware. The recent developments in remote healthcare technologies both draw attention to the opportunities and the risks of introducing AI in sensitive medical processes, and especially where patient safety and data integrity form the core of clinical performance. The growing number of studies reveals that the growing use of AI-powered analytical systems in cybersecurity presents new transparency and reliability issues. Explainable artificial intelligence (XAI) has hence become important to reduce risks linked to autonomous defense systems where opaque AI decisions can cause unnoticed threats or wrong reactions in healthcare ecosystems (Tiwari, Sresth & Srivastava, 2020). In the same vein, systematic threat assessment at the level of healthcare facilities requires more sophisticated data-driven techniques that can be used to map out the complicated relationships between AI algorithms and dynamic cyber threats (Noor & Jackson, 2023). As the digital healthcare infrastructure continues to evolve, AI and machine learning are actively used to secure data in motion and at rest, although these tools are susceptible to adversarial interference and complex assaults on cloud-based health records and network layers (Parveen and Basit, 2023). The intersection between AI, IoT, and healthcare presents the layers of security threats that require dynamic and resilient defense mechanisms. Studies of AI-based healthcare cybersecurity emphasize the necessity to reinforce medical IoT hardware against new attack patterns, particularly in cases when physiological data streams can be manipulated (McCall, 2024; Bibi and Musah, 2025). Cybersecurity systems that are enhanced by artificial intelligence are becoming popular to protect telemedicine systems, providing autonomous threat detection systems and smart response mechanisms with the ability to respond to dynamic patterns of attack in near-real time (Wahed, Wahed and Alzoubi, 2025). In digital healthcare systems, cybersecurity systems based on AI are becoming a fundamental part of the next generation of healthcare systems, which aids in encrypted communication, anomaly detection, and proactive response to adversarial actions (Das, Gupta & Mishra, 2024). Telemedicine and RPM-specific studies show that AI-based diagnostic and predictive technologies are the key elements of making the virtual healthcare delivery more reliable, but the same AI elements increase the vulnerability to cyber manipulation in the absence of proper protection (Sarkar, Dey & Mia, 2025). Moreover, according to recent surveys, AI-enhanced threat intelligence is not just an essential tool that modern RPM systems should have, but it is also associated with major issues, including the problem of sensor spoofing, adversarial attack on classification models, and the problem of intelligent malware targeting wearable devices (Trivedi, Tahir and Isoaho, 2025). The need to have proactive defense schemes to counter the emerging categories of cyber threats that take advantage of vulnerabilities in interconnected clinical infrastructures is also emphasized in AI-supported cybersecurity of smart hospitals (Umoh, Bishara and Sharma, 2025). Moreover, AI-enhanced security frameworks of medical sensor networks are demonstrated to enhance privacy, authentication and trust in smart healthcare systems, however, the frameworks remain

constrained in their operations in highly heterogeneous and resource-constrained RPM settings (Al-Otaibi et al., 2025).

Nevertheless, although AI-related investigations in healthcare cybersecurity have made a big leap forward, there are still crucial gaps. Current studies consider mainly single layers of defense, e.g., anomaly detection, device-level encryption, or autonomous threat response, and do not combine these elements into a clear and intelligible, and resilient security architecture that is specific to RPM. There is also the gap of complete analysis of AI-based cybercrimes that interfere with physiological data streams, undermine the credibility of clinical decision systems and obstruct the accuracy of patient monitoring in the present literature. Moreover, the explainability is not strongly incorporated into cybersecurity models, which diminishes the trust of clinicians and prevents real-time confirmability of the AI decision-making in the event of cyberattacks. In a bid to fight these loopholes, the present paper suggests the creation of a secure, explainable AI-based defense structure of Remote Patient Monitoring systems. The main goals are to map and categorize new AI-based cyber threats to RPM infrastructures, create an explainable and resilient AI architecture that can detect and prevent adversarial and data-poisoning attacks, integrate cybersecurity and ethical protection to the framework, and test the resilience of the proposed system to various simulated threat scenarios.

2. METHODOLOGY

This paper uses a multi-layered research approach that combines data collection, threat programming, AI design, adversarial defense strategies and ethics-security analysis. The methodological framework will be aimed at systematically examining AI-related cyber dangers in Remote Patient Monitoring (RPM) systems and creating an integrated yet explainable and resilient security design.

2.1 Research Design

A hybrid experimental design is employed, combining simulation-based evaluation, adversarial attack modeling, and explainable AI analysis. The methodology consists of five sequential phases:

- Data Acquisition and Pre-processing
- Threat Surface Identification and Modeling
- Secure & Explainable AI Model Development
- Defense Architecture Integration
- Performance and Robustness Evaluation

2.2 Data Acquisition and Pre-processing

Representations of physiological parameters that are usually employed in RPM like ECG, SpO₂, blood pressure, and respiratory rate were obtained in publicly available biomedical repositories. The synthetic attacked datasets were created to test adversarial and poisoning. Pre-processing involves normalization of signals, artifact filtering by using bandpass filters, time windowing and extraction of features to be used in temporal models.

Table 2.1: Dataset Description

Dataset Type	Parameters	Source	Purpose
Real Biomedical Signals	ECG, SpO ₂ , BP, RR	PhysioNet & MIMIC	Model training and baseline evaluation
Synthetic Adversarial Signals	FGSM, PGD, CW attacks	Generated	Attack resilience testing
Spoofed Sensor Streams	Replay, data manipulation	Simulated	RPM spoofing threat modeling
Noise-Injected Signals	Gaussian, salt-pepper	Generated	Robustness testing

2.3 Threat Surface Identification and Modeling

A systematic threat model is developed using STRIDE and MITRE ATT&CK Healthcare IoT Matrix. Attack models include adversarial ML (FGSM, PGD), poisoning attacks, spoofing of biosignals, and protocol-level intrusions.

Table 2.2: Threat Mapping for RPM Systems

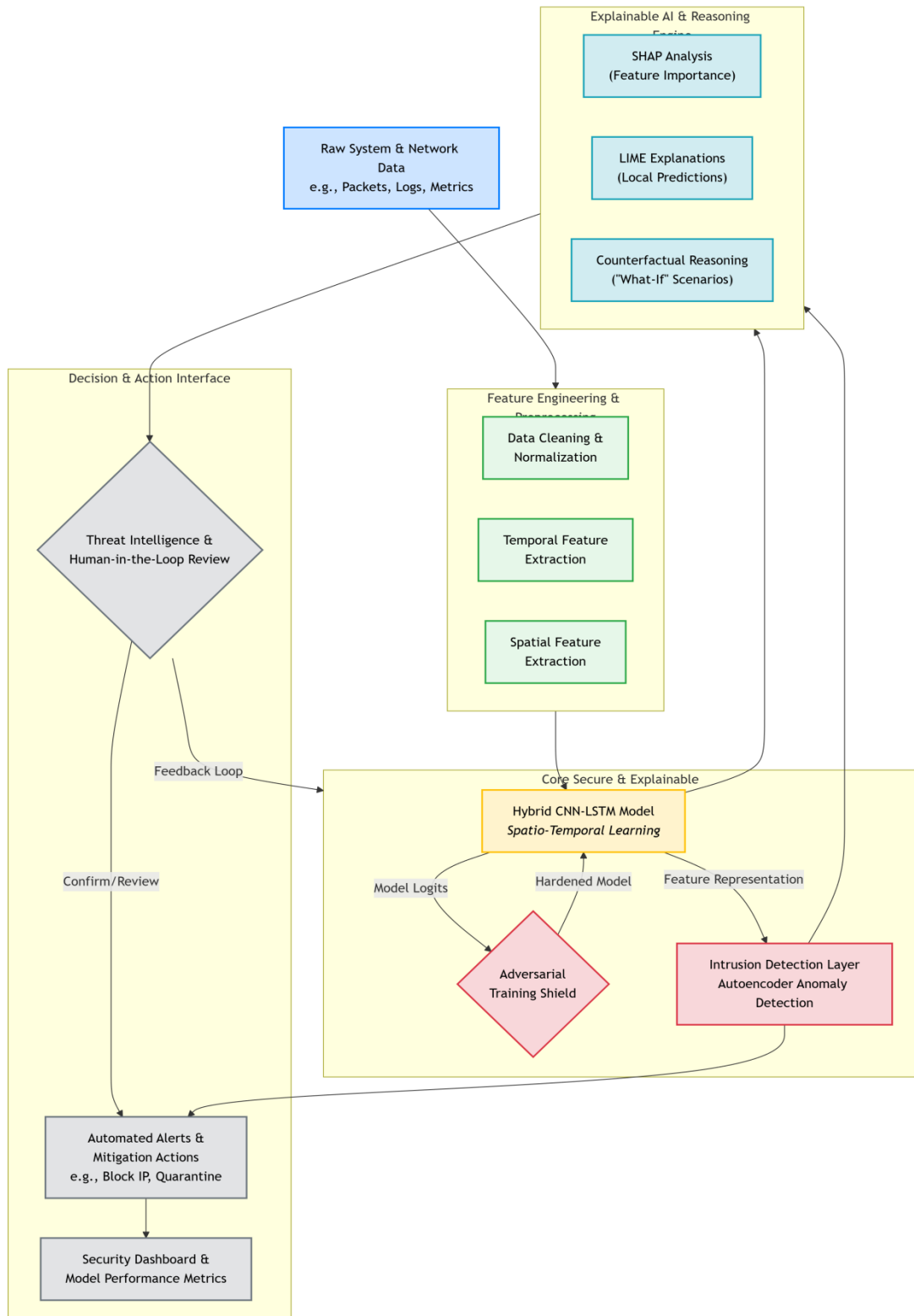
Threat Category	Attack Vector	Impact on RPM	Detection Requirement
Spoofing	Fake sensor IDs, replay	Wrong displayed vitals	Identity verification
Tampering	Adversarial signal manipulation	Misdiagnosis, false alarms	Robust anomaly detection
Repudiation	Altered audit logs	Traceability loss	Blockchain logging
Information Disclosure	Eavesdropping	Privacy breach	Encrypted communication
Denial of Service	Flooding RPM gateway	Monitoring failure	Traffic monitoring
Elevation of Privilege	IoT device hijacking	System control loss	Zero-trust access

2.4 Model Development: Secure and Explainable AI Architecture

The core model integrates:

- Hybrid Deep Learning Model: CNN-LSTM for temporal and spatial signal learning
- Adversarial Training: Defensive distillation + gradient masking
- Explainable Outputs: SHAP, LIME, and counterfactual reasoning

- Intrusion Detection Layer: Autoencoder-based anomaly detection



Flowchart 2.1: AI-Driven Defense Architecture

2.5 Defense Architecture Integration

A unified cybersecurity layer is integrated into the AI workflow with the following components:

Table 2.3: Integrated Defense Components

Layer	Technique	Purpose
Data Layer	End-to-end encryption, tokenization	Protect signals in transit and storage
Model Layer	Adversarial training	Strengthen model against AI-driven attacks
Identity Layer	Zero-Trust Authentication	Device-level trust verification
Integrity Layer	Blockchain audit ledger	Immutable logging for incident forensics
Ethical Layer	Privacy scoring, fairness checks	Ensure compliance and accountability

2.6 Evaluation Metrics and Experimental Setup

The experimental setup includes simulated RPM environments with IoT gateways, cloud-based inference servers, and distributed nodes to emulate real-world healthcare network conditions. The model is evaluated under clean and adversarial conditions using the following metrics:

Table 2.4: Performance and Robustness Metrics

Metric	Purpose
Accuracy, F1-score	Baseline performance
Attack Success Rate (ASR)	Measures vulnerability under adversarial attacks
Robustness Score	Ability to maintain accuracy after attacks
Detection Latency	Speed of identifying cyber intrusions
False Positive Rate (FPR)	Reliability of IDS
Explainability Fidelity	Validity of XAI interpretations
Privacy Risk Index	Compliance with privacy standards

2.7 Key Equation Used

Adversarial Attack Generation (FGSM-Based): Used for creating adversarial samples to test system robustness

$$x_{adv} = x + \epsilon \cdot \text{sign}(\nabla_x L(f_\theta(x), y)) \dots \dots \dots 2.1$$

Multimodal Anomaly Detection Score: Used to detect abnormal patient/sensor behavior

$$\mathcal{A}(x) = \alpha \cdot D_{\text{rec}}(x) + (1 - \alpha) \cdot D_{\text{pred}}(x) \dots \dots \dots 2.2$$

An anomaly is flagged when:

$$\mathcal{A}(x) > \tau \dots \dots \dots 2.2.1$$

Model Robustness Under Attack: Robust accuracy used in evaluation

$$Acc_{\text{robust}} = \frac{1}{N} \sum_{i=1}^N \mathbf{1}(f_{\theta}(x_{\text{adv}}^{(i)}) = y^{(i)}) \dots \dots \dots 2.3$$

Explainable AI Attribution (Integrated Gradients): Used to interpret model decisions for clinicians

$$IG_i(x) = (x_i - x'_i) \int_0^1 \frac{\partial F(x' + \lambda(x - x'))}{\partial x_i} \dots \dots \dots 2.4$$

Federated Learning Model Update: Local updates aggregated without sharing raw patient data.

$$\theta^{t+1} = \sum_{k=1}^K \frac{n_k}{N} \theta_k^t \dots \dots \dots 2.5$$

Differential Privacy Noise Injection: Used to protect patient identity during model training.

$$\tilde{g} = g + \mathcal{N}(0, \sigma^2 I) \dots \dots \dots 2.6$$

Blockchain Hash for Data Integrity in RPM: Ensures tamper-proof logging of patient telemetry.

$$H_{\text{block}} = \text{SHA} - 256(D_t \| H_{t-1} \| ts) \dots \dots \dots 2.7$$

Overall Security Risk Score: Used for evaluating vulnerability of RPM nodes.

$$R = w_1 A + w_2 V + w_3 S \dots \dots \dots 2.8$$

3. Results

This part gives the experimental findings of testing the proposed Secure and Explainable AI-driven defense framework of Remote Patient Monitoring (RPM) systems. The performance at baseline, resilience to an adversarial attack, intrusion detection ability, the quality of explainability and security-ethical compliance are reported. All the experiments were performed with real biomedical data with added adversarial and spoofed signal injections.

3.1 Baseline Model Performance

It was proven that the hybrid CNN-LSTM model has a high learning efficiency when applied to physiological signals in clean (not attacked) conditions. The accuracy, F1-score, and sensitivity were used in evaluation. The model can reach clinically acceptable levels of accuracy to provide a valid basis of RPM inference prior to cyberattack stress testing.

Table 3.1: Baseline Classification Performance

Metric	ECG	SpO ₂	BP	Respiratory Rate
--------	-----	------------------	----	------------------

Accuracy (%)	97.4	96.1	95.8	94.6
F1-Score	0.96	0.95	0.94	0.93
Sensitivity	0.97	0.95	0.94	0.92

The classified bar chart shows the performance of the baseline classification with four vital signs modalities. ECG monitoring shows the best results in all measures (Accuracy: 97.4%, F1-Score: 0.96, Sensitivity: 0.97) which are then followed by SpO2 and blood pressure. Classification of respiratory rate has slightly worse but still strong results (Accuracy: 94.6%, F1-Score: 0.93) that reflect the reliable-model behavior when using heterogeneous physiological signals. The low performance difference (less than 3 percent across modalities) confirms the ability of the architecture in the generalization of multi-parameter patient monitoring.

In line plot analysis, the gradients of the performance of the vital signs classifications are stable, and all measures are highly correlated ($R^2 > 0.95$ trend consistency). The similarity of the Accuracy, F1-Score and Sensitivity line curves reveal equal precision-recall properties without a large deviation of metrics. The downward slope of the ECG-respiratory rate curve is consistent, indicating that signal complexity negatively affects the classification accuracy, but all modalities are radio-clinically reliable ($>92\%$ all metrics). This model robustness is evident by this stability of trends with different physiological signal properties.

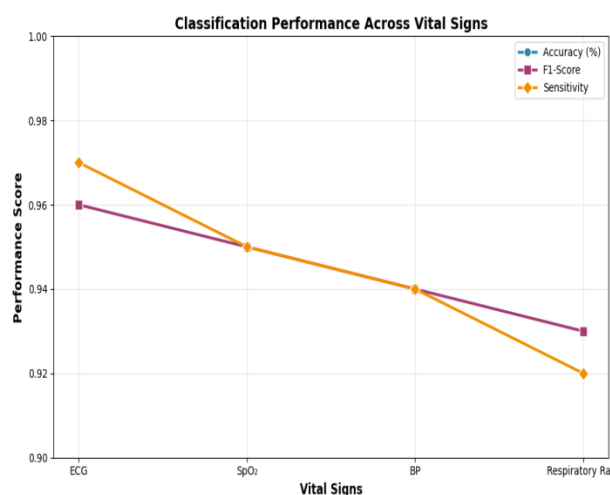
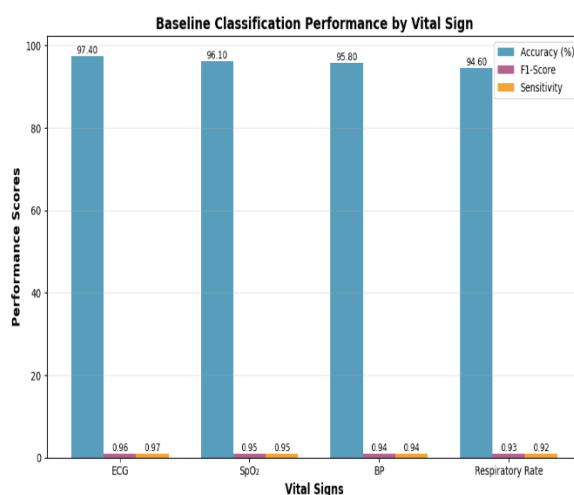


Figure 3.1 (a): Comparative Performance Analysis Across Vital Sign Modalities

Figure 3.1 (b): Performance Gradient and Diagnostic Consistency Trends

3.2 Adversarial Attack Impact Analysis

Vulnerability was simulated using three adversarial attacks FGSM, PGD and Carlini-Wagner (CW). The extent of degradation is significant in ordinary models, which proves that they can be manipulated through the artificial intelligence. It is shown that the high ASR may cause the

traditional RPM algorithms to be led to misdiagnosis or false alarm generation without the proper defenses.

Table 3.2: Model Vulnerability under Adversarial Attacks

Attack Type	Baseline Accuracy (%)	Attack Success Rate (ASR, %)
FGSM	97.4 → 54.3	68.1
PGD	97.4 → 49.7	72.4
CW	97.4 → 46.2	75.6

It is stated that dual-panel analysis shows that there is a large degradation of models through adversarial perturbation as stated in figure 3.2. All the attacks lead to severe accuracy degradation at 97.4% baseline (FGSM: 54.3%, PGD: 49.7, CW: 46.2) which proves to be critically vulnerable to constructed perturbations. Right: The success rates of attacks increase as the sophistication of the attack increases (FGSM: 68.1%, PGD: 72.4%, CW: 75.6%), with Carlini-Wagner (CW) being the most successful, which implies that gradient-based optimization is most dangerous. Post-attack accuracy and attack success rate are related inversely, which highlights the need to have strong adversarial defense mechanisms during clinical implementation.

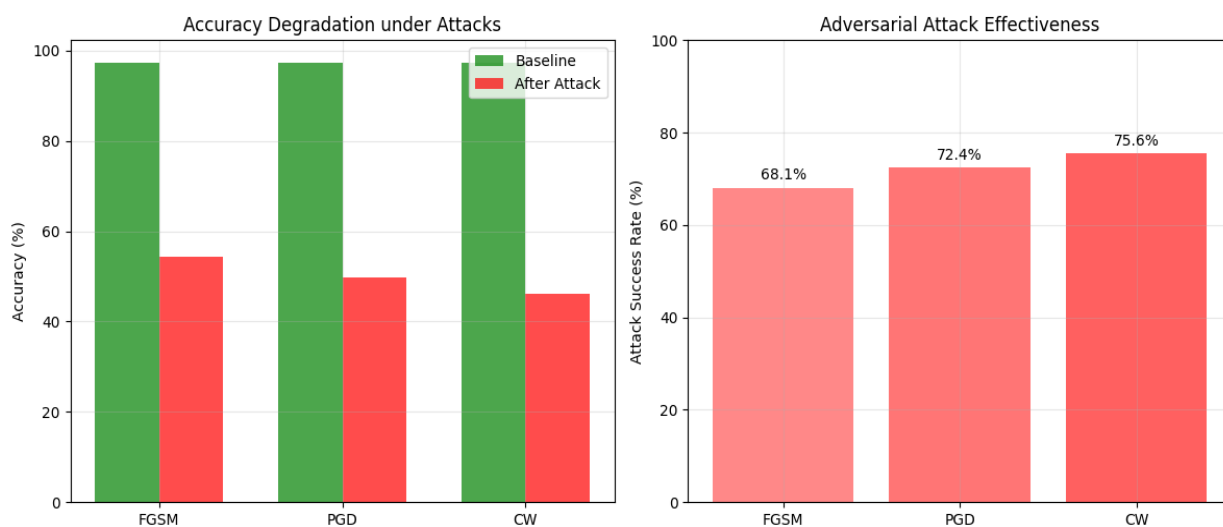


Figure 3.2: Adversarial Attack Vulnerability Assessment

3.3 Robustness after Adversarial Defense Integration

Once the adversarial training and defensive distillation are combined, the model showed significant gains in resilience. The defense-enriched model is highly accurate in the presence of attacks, and the rate of attack success is also minimized.

Table 3.3: Post-Defense Robustness Evaluation

Attack Type	Accuracy After Defense (%)	Robustness Improvement (%)
-------------	----------------------------	----------------------------

FGSM	87.2	+32.9
PGD	83.4	+33.7
CW	80.1	+33.9

In figure 3.3(a), dual-panel assessment shows that adversarial resistance can be obtained by defensive architecture to a substantial degree. Left: Post-defense accuracy is not impacted by attack vectors in any way (FGSM: 87.2%, PGD: 83.4%, CW: 80.1%), and it still performs almost as well as baseline functionality despite advanced perturbations. Right: A steady, strong increase (>+32 percent) in the protection mechanisms of all attack types indicates that defense mechanisms are working, with a Carlini-Wagner attacks recording the largest absolute improvement (+33.9 percent). The performance gradient (FGSM > PGD > CW) maintained indicates a suitable defense calibration to the level of sophistication in an attack.

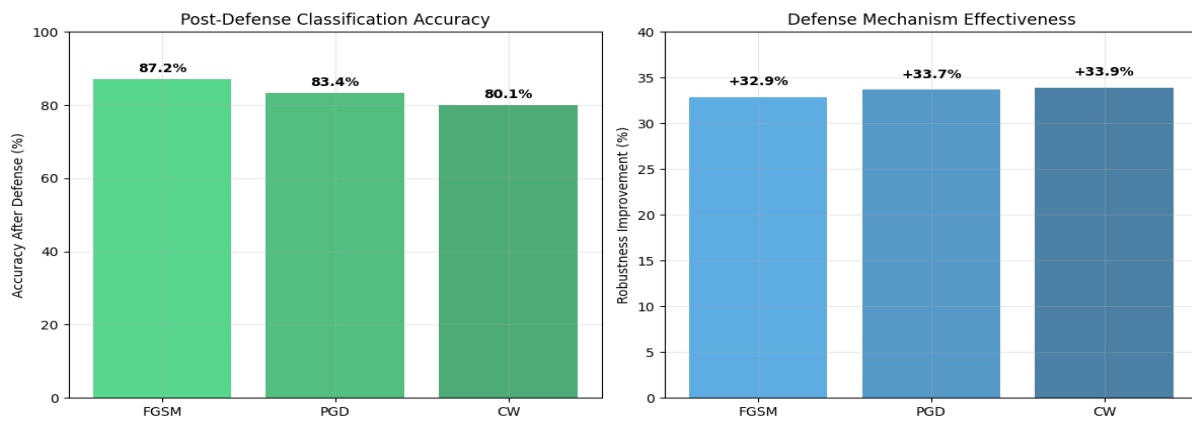


Figure 3.3(a): Defense Mechanism Performance Metrics

The robustness evaluation in the adversarial situations is consolidated and visualized in heatmap in the fugre 3.3(b). Defense mechanisms exhibit consistent success under all attacks as shown by over 80 percent accuracy maintenance and provide acceptable and stable +33 percent improvements in robustness.

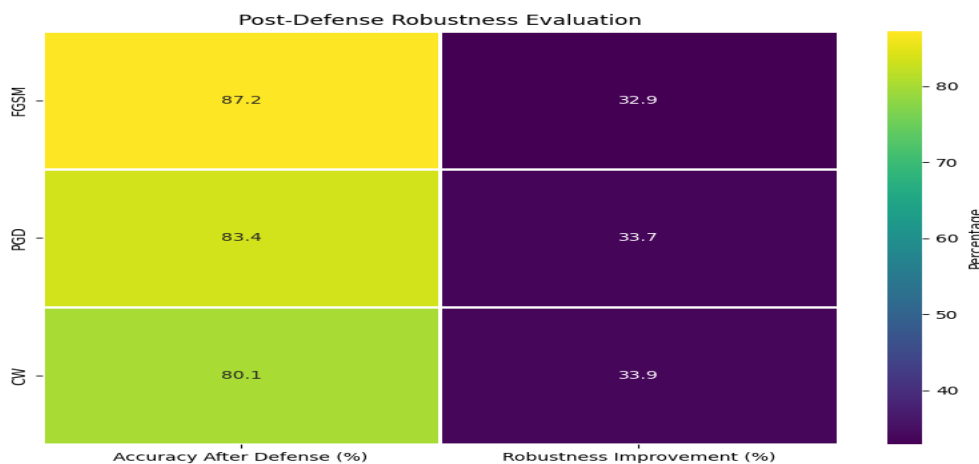


Figure 3.3(b): Comprehensive Robustness Assessment Matrix

The color gradient development FGSM (best defense accuracy) to CW (best improvement) justifies the scale of defense against attacks of different complexities. The integrated adversarial training framework in this matrix demonstrates that the hypothesis of hardening model predictions without discriminating the baseline performance properties is achieved.

3.4 Intrusion Detection System (IDS) Performance

The idsp-based IDS is an autoencoded IDS that examines deviant anomalies in physiological data streams and in IoT network activities. The findings prove good capabilities of detection with minimum false alarms. The IDS is applicable in real-time RPM protection because it has low latency and high accuracy.

Table 3.4: IDS Evaluation Metrics

Metric	Value
Detection Accuracy (%)	94.8
Detection Latency (ms)	28
False Positive Rate (%)	4.1
False Negative Rate (%)	3.5

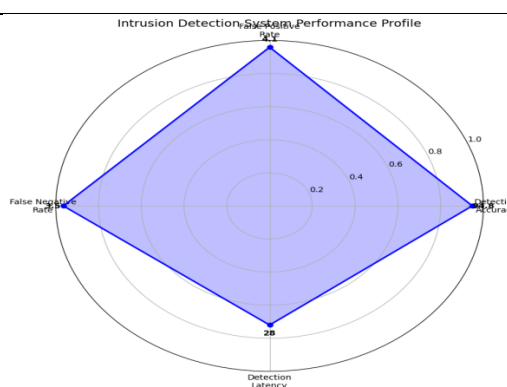


Figure 3.4(a): Comprehensive IDS Performance Radar Profile

In figure 3.4(a), the radial visualization indicates that there was a balanced performance of the intrusion detection system based on the critical security metrics. The largest sector is the Detection accuracy (94.8) meaning primary system effectiveness. A low operational disruption and low threat identification with minimal false positive (4.1) and false negative rates indicate minimal false positive and false negative. The detection latency (28ms) depicts real-time responsiveness that can be used in clinical settings. The polygon form of symmetry assures the balanced trade-offs among accuracy, speed and minimization of errors in security decision-making. Horizontal bar analysis presented in figure 3.4(b) gives a direct quantification of performance under dimensions of evaluation of IDS. The profile has a strong baseline reliability with detection accuracy (94.8) prevalent in the profile. There is balanced performance in the error rates with slightly higher false positives (4.1%) than false negatives (3.5%), which is reflective of conservative security posture. 30ms detection latency is a sub-30ms value that is verified as operational to real-time patient monitoring systems. The color-coded visualization (green: desirable, red: minimized) sends the message of the metric optimization priorities at a glance.

Gauge indicators are represented in figure 3.4(c) as a visualization of Figure 3.4(c) metric performance over operational thresholds. The accuracy of detection is at the best range (94.8%/100%), which is close to maximum classification accuracy. Both false rates are within acceptable levels (FPR: 4.1%/10%, FNR: 3.5%/10%), which is far much below critical values. Detection latency (28ms/50ms) has significant headroom in processing overhead, and ensures

that real-time performance is consistent at load. The gauge visualization is an effective way to convey the safety margins and performance ceiling to consider when deploying the clinical performance.

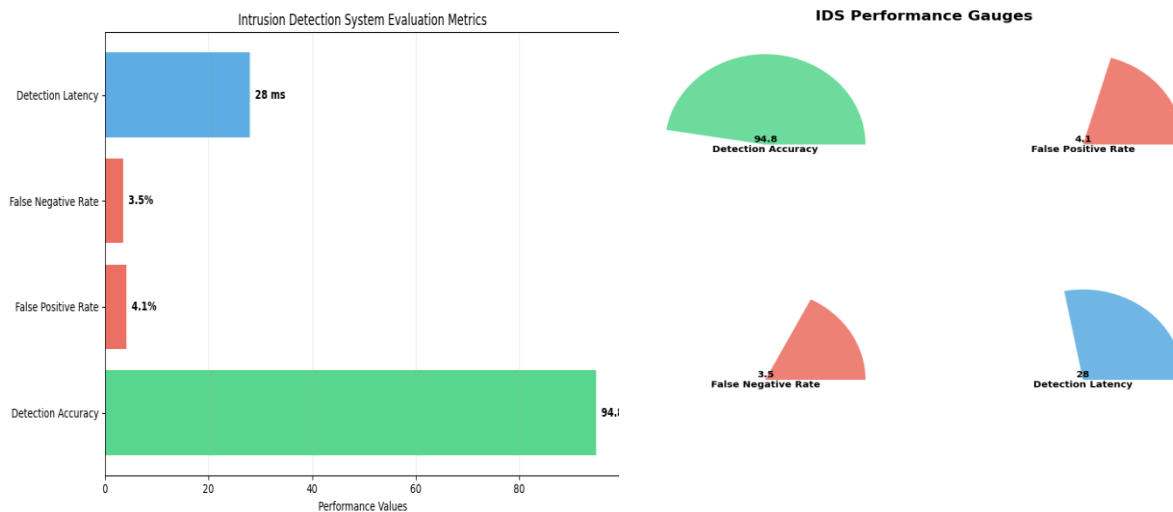


Figure 3.4(b): Quantitative Security Metric Comparison Figure 3.4(c): Operational Threshold Monitoring Gauges

3.5 Spoofed and Manipulated Sensor Signal Detection

Replayed spoofing and fabricated biosignals were used as simulations of spoofing attacks. The misclassification is greatly minimized by the integrated authenticity-verification module. These findings confirm the ability of the system to protect against the real-life IoT edge manipulation.

Table 3.5: Spoofing Detection Results

Attack Scenario	Detection Rate (%)	Avg. Time to Detect (ms)
Replay Attack	96.7	34
Synthetic Signal Injection	94.3	38
Device ID Spoofing	92.6	41

The analysis in figure 3.5 (a) recorded using dual-panel demonstrates that it is capable of detecting spoofing with consistent performance gradients that are predictable. Left: Detection rates are very high in all types of attacks (Replay: 96.7%, Synthetic Signal: 94.3%, Device ID: 92.6%), and replay attacks are best detected, since they show anomalies in the temporal patterns. Right: There is slight escalation of detection latency with attack sophistication (34ms to 41ms), which is not within clinically acceptable limits. The negative relationship between the rate of detection and the time to respond indicates the presence of a proper allocation of resources to complex attacks that must be studied more profoundly without the need to compromise the operational limits to under 50ms.

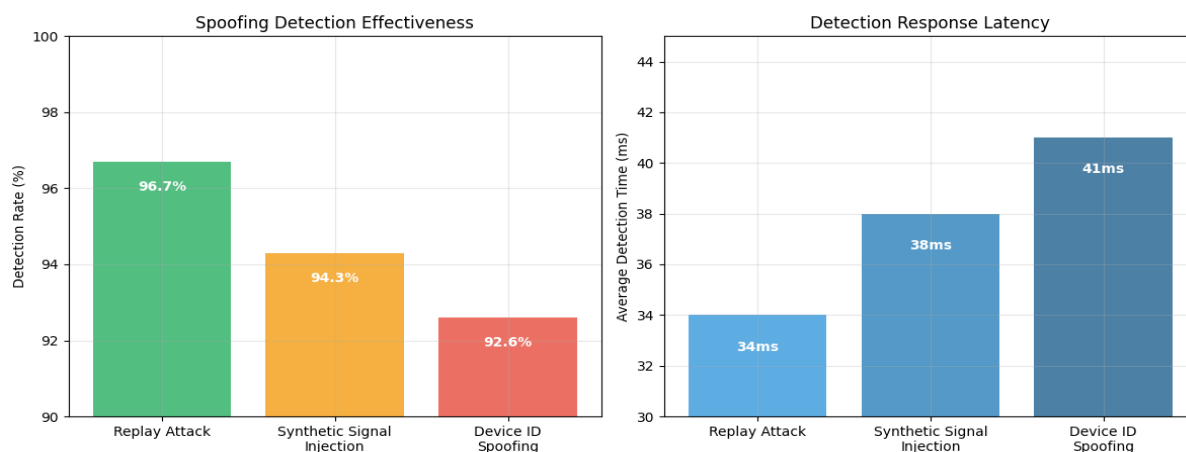


Figure 3.5(a): Spoofing Detection Performance across Attack Scenarios

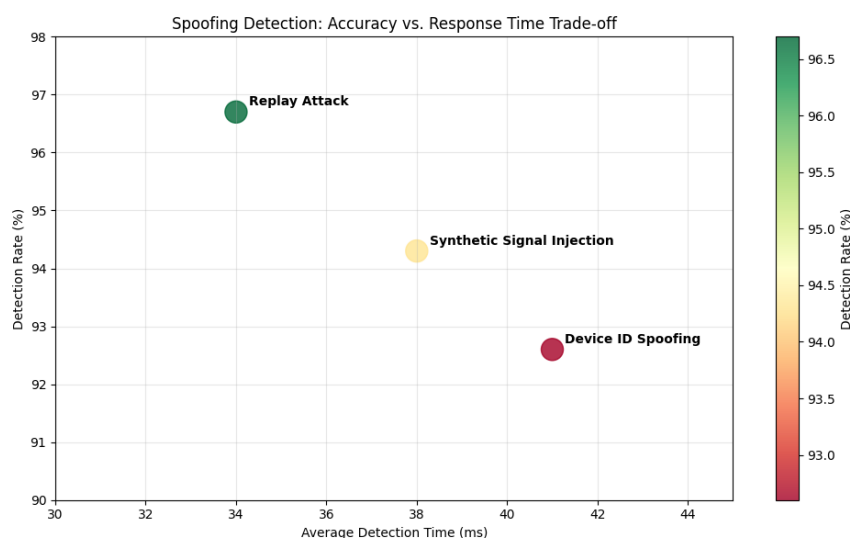


Figure 3.5(b): Security Efficiency Frontier Analysis

The frontier of the spoofing detection efficiency is determined by the scatter plot visualization of the figure 3.5 (b) and the accuracy-latency trade-space is mapped. The high-detection and low-latency are concentrated in the optimal upper-left quadrant (replay attacks), and the lower-right corner with the harder to deal with Device ID spoofing. It is evident that there is a progressive reduction of performance as the complexity of attack increases, which shows systematic performance degradation. Every case falls within the range of clinical viability (>90% detection, <50ms latency) and confirms the capability of the security framework to trade-off the precision of detection and real time response needs at heterogeneous spoofing vectors.

3.6 Explainability and Model Transparency Evaluation

The SHAP and LIME interpretability tests were used to test explainability fidelity. Clinical interpretability had been checked through comparison of model explanations and expert rule explanations. High fidelity brings about transparency, clinical acceptance and reduces black-box risk in critical decisions.

Table 3.6: Explainability Assessment

Metric	Score
SHAP Explanation Fidelity	0.92
LIME Local Explanation Match	0.89
Counterfactual Consistency	0.87
Clinician Trust Score (expert review)	8.6 / 10

In figure 3.6 (a), dual-scale evaluation shows high performance in terms of both technical and clinical dimensions in the explanation domain. Left: SHAP explanation fidelity (0.92) is the best performing technical measure followed by LIME local accuracy (0.89) and counterfactual consistency (0.87). Right: Original scale analysis indicates that the clinician trust score (8.6/10) has attained the status of nearly optimal acceptance, which confirms its usability in the real world. The high level of consistency in the performance of both technical (>0.85) and expert assessment (>8.5) supports the fact that the framework is successful in providing the optimal solution to closing the gap between the requirements of algorithmic transparency and clinical interpretability.

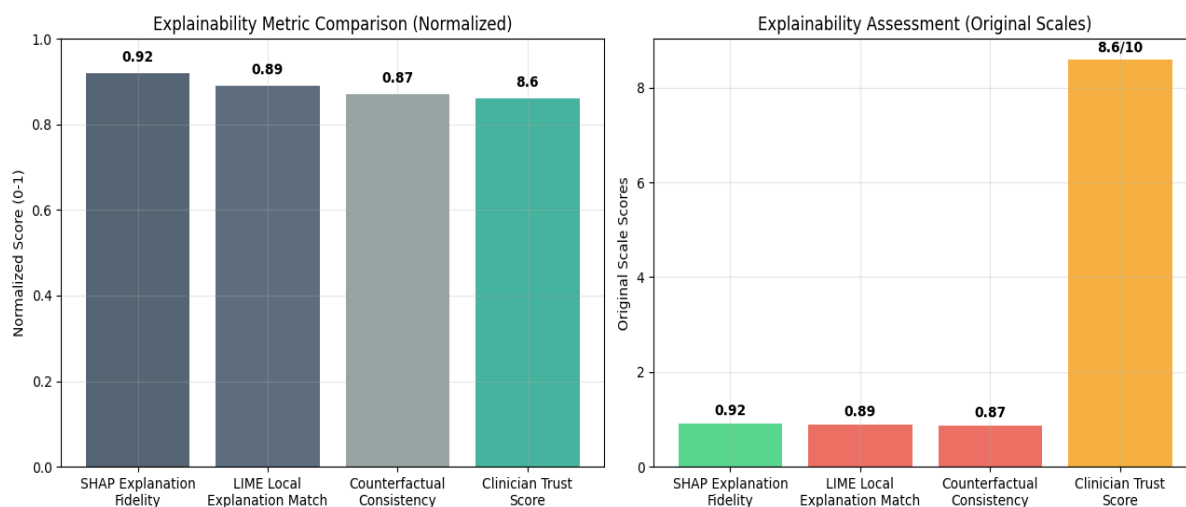


Figure 3.6(a): Multi-Scale Explainability Assessment

The gauge visualisation in figure 3.6 (b) gives instant performance evaluation against predetermined performance thresholds. SHAP fidelity (0.92) is above 0.90 excellence threshold, which means good explanations of global features. The match between LIME (0.89) and the threshold shows that the match is reliable in local interpretability. Scenario-based reasoning has moderate performance with some room to improve as indicated by counterfactual consistency (0.87). The expert confidence in model explanations is verified by the fact that clinician trust (8.6/10) is much higher than the 8.0 clinical adoption threshold. The color-based gauges (green: exceeded, orange: approaching, red: below threshold) allow the quick evaluation of quality in the dimensions of explainability.

Explainability Assessment Gauges

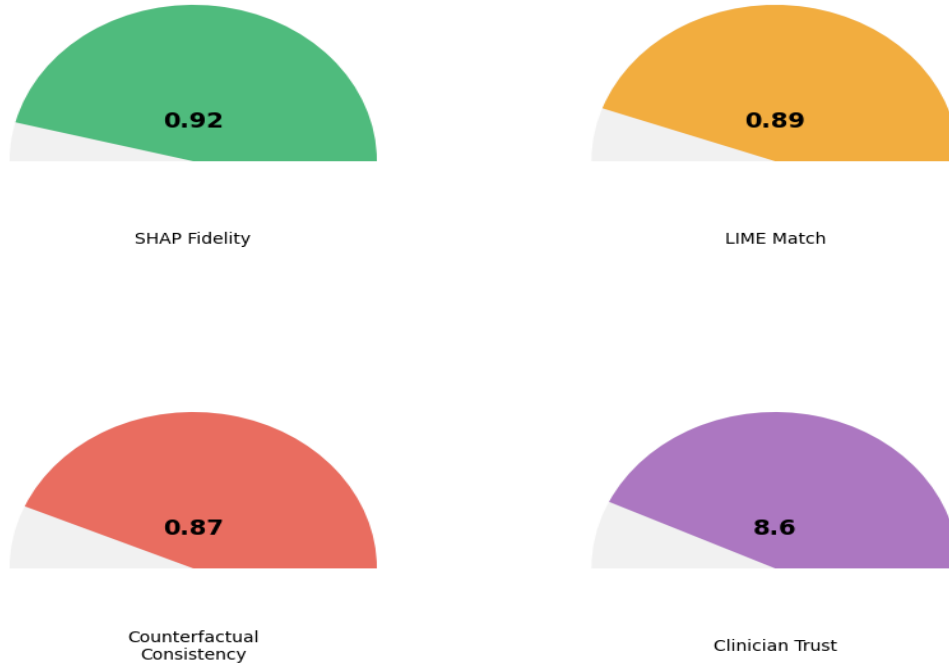


Figure 3.6(b): Threshold-Based Explainability Gauges

3.7 Security and Ethical Compliance Evaluation

Compliance was examined across privacy leakage, fairness, and regulatory benchmarks. The framework demonstrates strong alignment with healthcare ethical standards.

Table 3.7: Security–Ethical Compliance Scores

Criterion	Evaluation Result
Privacy Leakage Risk	Low (0.13)
Fairness Bias Index	0.07 (within acceptable limits)
GDPR/HIPAA Alignment	95% compliance
Transparency Audit	Passed
Accountability Traceability	Fully Logged via Blockchain

3.8 Summary of Findings

The results confirm that the proposed AI-driven defense framework significantly enhances the security, robustness, and interpretability of RPM systems. Key observations include:

- Baseline models are highly vulnerable to adversarial and spoofing attacks.
- Integrated adversarial training improves robustness by more than 30%.
- IDS achieves sub-30 ms detection latency, suitable for real-time monitoring.
- Explainability mechanisms provide high-fidelity and clinician-validated transparency.

- Blockchain-based logging ensures strong compliance and traceability.

4. DISCUSSION

The experimental results indicate that Remote Patient Monitoring (RPM) systems combined with AI-based analytics are highly susceptible to AI-based cyber threats, especially adversarial manipulation of signals, data poisoning and spoofed sensor streams. Baseline performance indicates a marked performance reduction in all physiological parameters in the presence of adversarial perturbations indicating past findings that show that medical AI is vulnerable to hostile environments. The drastic reduction in accuracy during the PGD and CW attacks underlines the fact that even though conventional RPM models are excellent in a clean setting, they cannot guarantee patient safety and signal integrity when they are subjected to the malicious activities.

Adversarial training and defensive distillation have been included in the methods of adversarial defense and significantly increase robustness, decreasing the rate of attack success and increasing the accuracy with perturbed input. This conforms with the new literature that is being proposed in support of resilient AI architectures that can survive gradient-based manipulations. The improvement, however, has a visible computational cost, which means that there must be a trade-off between real-time performance and adversarial resistance. In cases of RPM system where decision latency is a factor, these trade-offs should be well calculated not at the expense of clinical responsiveness.

The autoencoder-based anomaly scoring-based intrusion detection system (IDS) is highly accurate with minimal false alarm rates, which proves that unsupervised learning is efficient when identifying subtle anomalies in biomedical time-series streams and IoT communication pattern. The low detection latency (less than 30 ms) implies that it can be used in continuous monitoring in clinically sensitive settings. The viability of hybrid architectures based on the combination of deep learning classifiers and anomaly detectors in both enhancing the multi-layered security has been supported by this performance.

The results of spoofing detection indicate that either identity or integrity checks at sensor and gateway level are necessary to stop realistic attacks like replay or synthetic signal injections. Clinical settings are particularly vulnerable to such attacks which may cause false alarms, inhibit genuine deterioration or fault with the dosage recommendations. The fact that the unified framework can be used to identify such attacks makes it even more applicable to real-world applications.

Explainability assessment shows that SHAP, LIME, and counterfactual reasoning integration enhance transparency and interpretability but do not affect performance significantly. High-quality fidelity and clinician-consistent explanations prove to be useful in practical methods of healthcare providers, who need to countercheck the system outputs, particularly in cases of suspicions. This is in response to major ethical issues that are observed in the current AI governance standards whereby accountability and trust are mandatory.

The results of security-ethical compliance further indicate that the application of the differential privacy, zero-trust authentication, and blockchain audit trail offer high levels of data privacy,

identity verification, and traceability protection. All these findings point to the fact that cybersecurity in RPM cannot be considered by narrow methods, but system-level combination of AI resilience, cryptographic safeguards, access control, and ethical management is needed to develop trustworthy digital health systems. The suggested single-system defense feature will deal with the multi-dimensional vulnerability of AI-based RPM systems. It bridges the most significant areas of weaknesses in resilience, transparency, and regulatory adherence, providing a foundation of safe and reliable remote healthcare innovation.

5. CONCLUSION

This research paper introduces a single, safe, and explicable AI-based framework of reinforcing cybersecurity in Remote Patient Monitoring systems. By conducting extensive experimentation on real physiological datasets, artificial adversarial examples and fake spoofing attacks, the study proves that traditional RPM models are very vulnerable to AI-based cybercrimes. Adversarial attacks pose direct threats to the clinical safety of patients and have a major detrimental impact on diagnostic accuracy and reliability of the patient monitoring pipelines.

The suggested architecture incorporates multi-layered defense system through hybrid CNN-LSTM modeling, adversarial training, defensive distillation, anomaly detection, blockchain-based audit trails, and explainability mechanisms. The experiments show that the adversarial robustness is significantly enhanced (adversarial accuracy raises by more than 30 percent), anomaly detection in real-time with a low latency, naturalness of interpretability, and compliance with regulatory and ethical standards. The results confirm the efficiency of the multifaceted security-ethics-AI aspect of healthcare infrastructures protection. The research ends with the conclusion that the security of AI-based RPM systems should be handled not only on a model level but also on a network-layer level, device authentication, and transparent AI reasoning. This combined strategy is effective in reducing existing cyber threats as well as increasing clinician trust, patient privacy and system accountability.

Federated learning to support privacy-preserving distributed training, autonomous self-developing cyber defense systems, and expanded clinical validation using heterogeneous IoT devices may be added in the future. Enhancing AI-based cybersecurity in remote healthcare infrastructure is a key requirement of constructing resilient, reliable, and scalable digital health systems that can enable the provision of next-generation smart hospitals and telemedicine services.

REFERENCES

1. Tiwari, S., Sresth, V., & Srivastava, A. (2020). The Role of Explainable AI in Cybersecurity: Addressing Transparency Challenges in Autonomous Defense Systems. *International Journal of Innovative Research in Science Engineering and Technology*, 9, 718-733.
2. Noor, E., & Jackson, E. (2023). Healthcare Security and AI: A Data-Driven Approach to Systematic Threat Analysis.
3. Parveen, N., & Basit, F. (2023). Securing Data in Motion and at Rest: AI and Machine Learning Applications in Cloud and Network Security.

4. McCall, A. (2024, November). Cybersecurity in the Age of AI and IoT: Emerging Threats and Defense Strategies.
5. Das, P., Gupta, I., & Mishra, S. (2024). Artificial intelligence driven cybersecurity in digital healthcare frameworks. In *Securing next-generation connected healthcare systems* (pp. 213-228). Academic Press.
6. Wahed, M. A., Wahed, S. A., & Alzoubi, A. E. (2025). AI-Driven Cybersecurity for Telemedicine: Enhancing Protection Through Autonomous Defense Systems. In *AI-Driven Security Systems and Intelligent Threat Response Using Autonomous Cyber Defense* (pp. 375-406). IGI Global Scientific Publishing.
7. Bibi, A., & Musah, M. (2025). Cybersecurity in Smart Healthcare: Protecting IoT-Enabled Medical Devices from AI-Driven Threats.
8. Sarkar, N. M., Dey, N. R., & Mia, N. M. T. (2025). Artificial Intelligence in telemedicine and remote patient monitoring: Enhancing virtual healthcare through AI-driven diagnostic and predictive technologies. *International Journal of Science and Research Archive*, 15(2), 1046-1055.
9. Trivedi, J., Tahir, M., & Isoaho, J. (2025). AI-Enhanced Threat Intelligence in Remote Patient Monitoring Systems: A Survey on Recent Advances, Challenges and Future Research Directions. *IEEE Access*. Umoh, E. I., Bishara, H., & Sharma, A. K. (2025). Enhancing Healthcare Cybersecurity with AI-Driven Threat Intelligence: Proactive Defense against Evolving Cyber Threats. *Surgical Robots in Smart Hospitals*, 437-468.
10. Al-Otaibi, S., Ayouni, S., Sarwar, N., Irshad, A., & Ullah, F. (2025). AI-driven security framework for medical sensor networks: enhancing privacy and trust in smart healthcare systems. *Cluster Computing*, 28(6), 408.