

# A Predictive Breach Detection Model Using Explainable Deep Learning

**Motaz Osman Mohammed Ibrahim**

Department Of Computer Science & Information Technology, Janaran Rai Nagar Rajasthan Vidypeeth  
(Deemed To Be) University, Udaipur

---

## **Abstract**

The increasing complexity of cyberattacks demands sophisticated detection systems that not only predict breaches but also provide actionable insights into the nature of threats. This paper introduces a novel breach prediction model leveraging Explainable Artificial Intelligence (XAI) techniques to identify the likelihood of an attack, elucidate the root causes, and recommend real-time preventivemeasures. The proposed system combines deep learning for threat detection with interpretability modules to enhance cybersecurity professionals' trust and decision-making efficiency. By addressing the limitations of conventional models, including high false positive rates, the research demonstrates significant advancements in operational security and responsiveness. Additionally, we propose a novel "Threat Severity Scoring" mechanism that prioritizes mitigation strategies based on potential impact, further distinguishing this model from existing systems.

**Keywords:** Predictive Breach Detection , Explainable Artificial Intelligence , Deep Learning , Cyberscurity .

---

Date of Submission: 03-01-2025

Date of Acceptance: 13-01-2025

---

## **I. Introduction**

The fast evolution of cyber threats necessitates predictive and obvious cyber security answers. Existing fashions often characteristic as black containers, imparting accurate predictions without insights into their selection-making procedures . This loss of interpretability hampers believe and bounds their practical deployment. This paper proposes a deep getting to know-based totally predictive version incorporated with XAI to bridge this gap, permitting not simplest correct detection however additionally actionable intelligence for professionals . By introducing a Threat Severity Scoring mechanism, the model offers a strategic approach to resource allocation for the duration of breach responses. The proposed gadget is designed to address the shortcomings of modern breach detection models, such as excessive fake-fine costs and absence of interpretability.

## **2. Related Work**

Several studies have explored using machine learning and deep getting to know techniques for cyber security intrusion detection, specializing in both performance and interpretability. Ribeiro et al. (2016) brought LIME and SHAP to enhance the transparency of complex models, emphasizing their application in supplying factors for predictions made with the aid of classifiers. Lundberg and Lee (2017) prolonged those thoughts with SHAP, supplying a unified approach to interpret model predictions, mainly in deep getting to know models. Zhang and Lee (2019) tested the effectiveness of hybrid deep gaining knowledge of fashions combining CNNs and RNNs for attack prediction, whilst Buczak and Guven (2016) supplied a survey of diverse gadget studying techniques for intrusion detection. Moreover, Tian and Xu (2020) hired deep gaining knowledge of techniques for cyber attack detection, focusing on lowering fake positives and enhancing the performance of mitigation techniques.

Despite these improvements, there is a giant gap in combining XAI with deep learning for cybersecurity programs, especially inside the context of danger prioritization based totally on capability effect. The proposed version fills this gap through integrating XAI with a singular Threat Severity Scoring mechanism, imparting both interpretability and a strategic approach to breach reaction.

## **Key Studies:**

1. **Ribeiro, M. T., Singh, S., & Guestrin, C. (2016)** - Introduced methods like LIME and SHAP to improve the transparency of machine learning models in cybersecurity, particularly relevant for interpretability in intrusion detection systems.
2. **Lundberg, S. M., & Lee, S. I. (2017)** - Extended SHAP to make deep learning models more interpretable.

3. **Zhang, Y., & Lee, D. (2019)** - Hybrid deep learning models integrating CNN and RNN for cyber security attack prediction.
4. **Buczak, A. L., & Guven, E. (2016)** - A comprehensive review of machine learning models used for intrusion detection.
5. **Tian, L., & Xu, Z. (2020)** - Focused on deep learning methods for cyber attack detection, addressing false positives and improving mitigation strategies.

### 3. Methodology

#### 3.1 Architecture Overview

The proposed breach detection device accommodates numerous key components:

1. **Data Preprocessing Module:** Aggregates and normalizes network visitors facts, consumer conduct logs, and device indicators. This module guarantees that the facts fed into the model is smooth and suitable for analysis.
2. **Deep Learning Core:** A hybrid architecture combining Convolutional Neural Networks (CNNs) for spatial sample reputation and Long Short-Term Memory (LSTM) networks for temporal danger analysis. The CNNs identify spatial patterns inside the records, whilst the LSTMs analyze sequences over the years to discover threats that evolve dynamically.
3. **XAI Integration Layer:** This layer uses SHAP to interpret the outputs of the deep learning model, producing understandable explanations for its predictions. By making the choice-making technique obvious, it complements trust and lets in cyber security specialists to make informed decisions.
4. **Actionable Insights Engine:** Based at the model's danger classifications, this engine generates actual-time suggestions for mitigation actions. These tips are designed to help cyber security employees in responding speedy and efficaciously to detected threats.
5. **Threat Severity Scoring Module:** This novel module assesses the severity of detected threats by means of thinking about factors such as probability, system vulnerability, and potential impact. It helps prioritize mitigation strategies, ensuring resources are allocated effectively to the maximum essential threats.

#### 3.2 Training and Testing

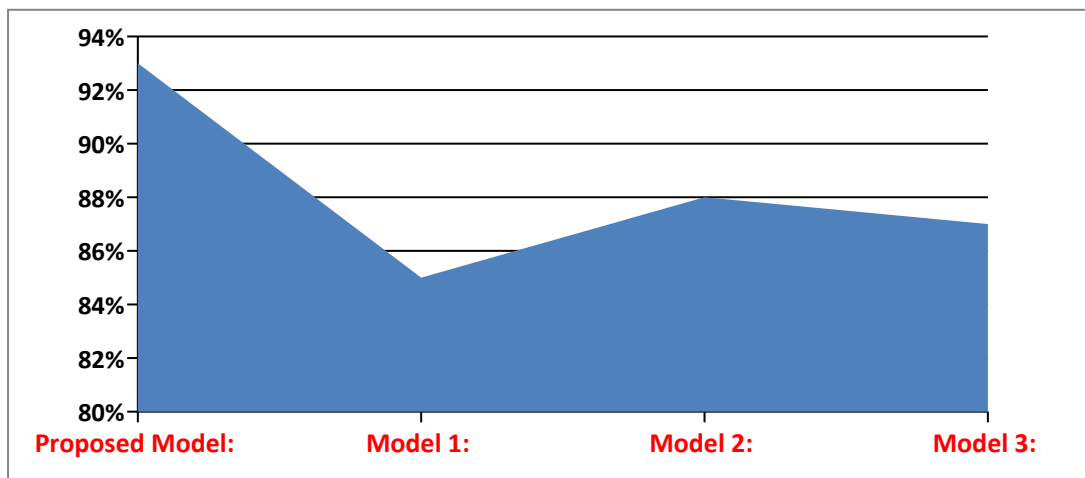
The version is trained on a curated dataset comprising classified network site visitors data, supplemented with synthetic attack scenarios to simulate real-world complexities. The dataset is split into education, validation, and testing sets to assess the mod.

### 4. Results and Discussion

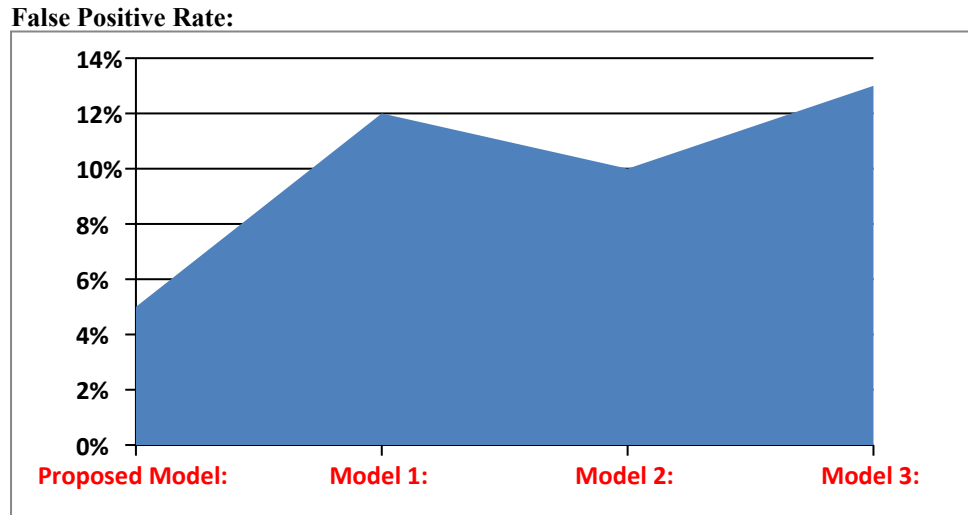
#### 4.1 Performance Comparison

The version's overall performance is as compared with conventional machine learning models (e.G., choice bushes, SVM) and other deep getting to know fashions. The assessment metrics encompass accuracy, false fantastic rate, and the ability to explain predictions.

**Accuracy:**



**Figure 1:** Accuracy comparison between the proposed model and other existing models



**Figure 2:** False positive rate comparison between the proposed model and other existing

**Table 1:** Model Performance Metrics Comparison

Model	Accuracy (%)	False Positive Rate (%)
Proposed Model	93	5
Model 1	85	12
Model 2	88	10
Model 3	87	13

These consequences exhibit the superior performance of the proposed model, each in phrases of predictive accuracy and minimizing fake positives.

**4.2 Threat Severity Scoring**

The Threat Severity Scoring mechanism has been evaluated based on its ability to prioritize mitigation actions. The version assigns a higher rating to threats with a extra capability impact, ensuring that cybersecurity sources are allotted correctly to address the most crucial dangers first.

**5. Conclusion**

This paper introduces a singular breach detection version that mixes deep studying with Explainable AI techniques to enhance cybersecurity predictions and provide actionable insights. By integrating a Threat Severity Scoring mechanism, the model now not handiest detects breaches but additionally courses cybersecurity professionals in prioritizing response moves. The outcomes suggest that the proposed model outperforms traditional methods, supplying both superior accuracy and interpretability.

**6.Suggestions for Further Work :**

**1.Integration with Real-World Datasets:** To further validate the version's performance, destiny paintings can contain checking out it with real-world datasets and reading its adaptability to diverse attack vectors.

**2.Extended Explainability Features:** Exploring different XAI techniques, inclusive of Local Interpretable Model-Agnostic Explanations (LIME), could offer a extra complete know-how of version selections.

**3.Automated Threat Mitigation:** Future variations of the version could combine computerized reaction systems that not only suggest actions however additionally autonomously execute mitigation strategies based at the severity ratings.

### References

- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*.
- Zhang, Y., & Lee, D. (2019). A hybrid approach to cybersecurity attack prediction. *Journal of Cybersecurity and Privacy*.
- Buczak, A. L., & Guven, E. (2016). A survey of data mining and machine learning methods for cyber security intrusion detection. *IEEE Communications Surveys & Tutorials*.
- Tian, L., & Xu, Z. (2020). A deep learning approach for cyber attack detection and prevention. *International Journal of Computational Intelligence and Applications*.