

# Intrusion Detection in Industrial Control Systems Using ML-Based Log Analysis

Doaa Mohammad Majed<sup>1</sup>, Aya Falah<sup>2</sup> and Vishal Jain<sup>3</sup>

<sup>1</sup>*Al-Turath University, Baghdad, Iraq, 10013 Baghdad, Iraq*

<sup>2</sup>*Medical Technical College, Al-Farahidi University, 10065 Baghdad, Iraq*

<sup>3</sup>*Sharda School of Computing Science & Engineering, Sharda University, 201310 Greater Noida, India*  
*doaa.mohammad@uoturath.edu.iq, aya.falah@life-rdh.org, vishal.jain@sharda.ac.in*

**Keywords:** Industrial Control Systems (ICS), Intrusion Detection, Machine Learning, Log Analysis, Anomaly Detection, Autoencoder, Cybersecurity, Feature Importance, ROC-AUC, Critical Infrastructure.

**Abstract:** The foundation of operational critical infrastructure is known as Industrial Control Systems (ICS) and due to their growing interconnection they have become targets of advanced cyber-attacks. The conventional intrusion detection systems (IDS) are limited because of being rule-based and are unable to identify zero-day exploits and subtle anomalies. This paper offers an intrusion detection system, which is a machine learning (ML) system that will examine ICS log data to detect the anomalies efficiently and accurately. The process includes systematized data gathering, preprocessing, feature discovery and mixed ML modeling by autoencoders and classifiers. The results of the experiment prove that the proposed system is more precise, has higher recall, and AUC values compared to traditional approaches. The distribution of anomaly scores indicates the existence of a distinct boundary between normal behavior and attack behavior, whereas the analysis of the feature importance offers operational information on the important log parameters. The framework has an AUC of 0.984 indicative of its strength and capability to be used in real-time. Also, the architecture of the model is modular which facilitates future scalability and explainability. The superiority of the system to other existing log based and network based IDS models is proved by comparative benchmarking. The research points out the opportunities of ML to improve the ICS cybersecurity using data-driven, adaptive, and explainable ways to do it.

## 1 INTRODUCTION

The Industrial Control Systems (ICS) are the important means of controlling and automation of processes of energy, water treatment, manufacturing, and transportation sectors. These systems have become much more vulnerable to cyber threats as they are exposed to the corporate IT networks and the internet as well. The emergence of advanced cyber-physical events, such as Stuxnet, Triton, and Industroyer, has highlighted the pressing need to consider the development of strong security measures that are applicable in the context of ICS. They are particularly susceptible to attacks as traditional security measures have a tendency to be insufficient to identify slow-moving, subtle yet domain-specific threats that characterize ICS attacks [1], [2].

Offering valuable insights on detecting anomalies, log information gathered by ICS, including event logs, system alerts, and communication traces, can be

used as valuable information. Log-based intrusion detection is passive, scalable, and less intrusive unlike active monitoring tools. Nevertheless, due to the huge volume, heterogeneity, real-time character of log data, it becomes impossible to access the information manually, and intelligent approaches to automated detection are required (Jeffrey et al., 2023) [3]. The traditional rule-based and signature-based intrusion detection systems (IDS) cannot keep up with the pattern of new or changing attacks, and usually, they result in a high level of false-positive data or even fail to detect emerging threats at all (Ahmed et al., 2016) [4].

To address these shortcomings, scholars have shifted their attention more towards machine learning (ML) and deep learning (DL) techniques. These models are capable of directly learning complicated patterns of normal behavior and abnormal behavior using past data and are thus very suitable in the dynamic threat environment of ICS. Supervised

algorithms like the Random Forests and Support Vector Machines and unsupervised algorithms like autoencoders and isolation forests have been promising in this respect. More complex deep learning algorithms such as Convolutional Neural Networks (CNNs) have been considered to classify spatial and temporal patterns of process data (Kravchik et al., 2018) [1]. Nevertheless, these models have shortcomings that include balancing of classes, unlabeled attack data, and variability by domain that do not help them to generalize in other settings of ICS (Pan et al., 2009) [5].

Recent reports highlight the concept of modification of the traditional ICS security mechanisms to intelligent, adaptive mechanisms. At the smart grid level, in particular, ML-based methods are actively considered to identify both known and unknown threats without disrupting the stability of the operations (Mirzaee et al., 2022) [6]. Accurate, real-time, and scalable detection systems are also becoming a condition that arises in next-generation infrastructures, including Positive Energy Districts (PEDs) and smart buildings, where data-driven security is also becoming a requirement (Han et al., 2024) [7].

In this paper, we seek to transform and test a machine learning-based intrusion detection system which uses ICS log data to identify the threats. We make the following contributions: (i) an application-specific preprocessing pipeline of multi-source ICS logs, (ii) a comparative study of ML models to detect known and zero-day attacks, and (iii) information on the problem of deployment and mitigation to real-world ICS deployments.

## 2 LITERATURE REVIEW

The development of cybersecurity threats to the Industrial Control Systems (ICS) has rendered much attention to machine learning (ML) and deep learning (DL) techniques of anomaly detection. Although the conventional rule-based systems are still used, research has shown them to be incapable of identifying more advanced, zero-day or slow evolving threats. Specifically, the challenges that are peculiar to ICS environments include few labelled datasets, real-time detection necessity, and essential safety restrictions. The literature review presents a summary of the major recent works that have been created to resolve these issues based on the advanced approaches to ML and DL, as stated in Table 1.

Aslam et al. (2025) [8] undertook an extensive survey that benchmarked deep learning system architectures including convolutional neural networks (CNNs), recurrent neural networks (RNNs), and autoencoders to detect intrusion in the ICS system. Their review emphasized the increasing dependence on data-oriented methods, and their effectiveness as well as the limitations in practice in the industrial environment. On the same note, Umer et al. (2022) [9] have reviewed the use of ML applications in many other ICS platforms such as SCADA systems and power grid elements. They highlighted the difficulties in managing unbalanced data sets, dynamic network settings and the applicability of trained models to real world applications.

Resorting to hybrid solutions that involve a combination of several methods has also become popular. Aslam et al. (2025, August) [10] came up with a multi-feature anomaly detection framework, which combines the use of ML, DL, and statistical approaches to enhance detection accuracy. Their approach delivered improved resilience with a combination of heterogeneous features of ICS logs, network traffic and system status information. This was furthered by Gulzar and Mustafa (2025) [11] who proposed an interdisciplinary deep learning framework of anomaly detection in diverse ICS domains. They showed that hybrid models are more amenable to unobservable patterns of data, but the costs of training and computation are of concern.

Explainable AI (XAI) is now an essential variable in the ICS anomaly detection, as the human operators should be able to interpret the system alerts promptly. Birihanu and Lendak (2025) [12] solved this issue by presenting a detection technique that uses a correlation with the interpretability attribute. Their model offered transparency in the process of identifying anomalies and giving domain experts the opportunity to test system behavior. On the same note, Jadidi et al. (2023) [13], [14] have also used statistical correlation of log sequences to reveal abnormal dependencies among ICS components, with more focus placed on time-series structure and sequential relationships.

Although the majority of studies use tabular or network data, Zhang et al. (2025) [15] is one of the recent studies to offer a wider field of research to image-based ICS diagnostics. Their high-resolution industrial image anomaly detection model showed that visual inspection, combined with log data, is able to improve the localization of faults and detection accuracy in sensor-rich environments.

Table 1: Summary of reviewed literature on ML/DL-based intrusion detection in ICS.

Ref. No.	Authors	Focus Area	ML/DL Models Used	ICS Component	Key Contribution	Limitations Identified
[8]	Aslam et al. (2025)	Survey of DL methods for ICS security	CNN, RNN, Autoencoders	General ICS	Comparative analysis of DL models	Limited discussion on interpretability
[10]	Aslam et al. (2025 Aug)	Hybrid anomaly detection framework	ML + DL + Statistical	ICS Process Logs	Feature fusion improves detection accuracy	Computational overhead
[9]	Umer et al. (2022)	Survey on ML in ICS	SVM, RF, KNN	SCADA & Network Logs	Applications, datasets, evaluation metrics	No deployment-oriented discussion
[12]	Birihanu & Lendák (2025)	Explainable anomaly detection	Correlation-based + XAI	ICS Network	Transparent models for operator trust	May miss deep latent features
[15]	Zhang et al. (2025)	Image-based anomaly detection	High-resolution DL models	Visual ICS Data	Shift toward visual-log fusion in ICS	Needs sensor-rich environments
[11]	Gulzar & Mustafa (2025)	Interdisciplinary DL-based detection	Deep Neural Networks (DNN)	General ICS	Flexible model integration across domains	Training time and scaling issues
[13]	Jadidi et al. (2023)	Correlation anomaly detection	Correlation + Statistical	ICS Network Logs	Time-series correlation between log streams	Evaluation limited to a single dataset

Table 2: Summary of ICS log dataset used in the study.

Dataset Name	Log Type	Size (MB)	Records	Attack Types	Duration	Labels
ICS-WaterSim	System & Network	145	25,000	DoS, Replay, MITM	7 days	Binary

To conclude, recent studies indicate that there is a significant advancement in implementing ML and DL towards ICS anomaly detection. Nonetheless, interpretability, domain adaptability and effective real-time operation challenges continue to exist. Table 1 offers a comparative perspective of these contributions in terms of their focus, techniques applied and limitations established.

### 3 METHODOLOGY

This section describes the proposed machine learning-based intrusion detection system for Industrial Control Systems (ICS) using structured log data. The system pipeline covers all key stages of the workflow, including log collection, preprocessing, feature engineering, model training, anomaly detection, and alert generation.

#### 3.1 System Architecture Overview

The suggested architecture involves the first stage of gathering raw logs in ICS settings (i.e. SCADA systems or Programmable Logic Controllers (PLCs))

and subsequent structured preprocessing, feature extraction and application of supervised/unsupervised models. The final output will consist of real-time score of anomalies and intrusion alerts [16]-[18]. All the elements of the framework are combined to provide low-latency and high-reliability in industrial environments.

#### 3.2 Dataset Description and Log Collection

The data that is used in this study is the real time ICS logs (normal and attack) logs. System events and network activities logs were obtained in a simulation environment of a water treatment. The data set is tagged, making it possible to train supervisedly and to assess it. The most important statistics of the dataset are summarized in Table 2.

#### 3.3 Preprocessing and Feature Engineering

Raw logs were read and converted into structured format with regular expression templates and time-stamped to make them synchronous. The null entries

were deleted, and categorical fields (e.g., protocol type) were coded on numbers. Continuous features were normalized using a z-score scaling:  
 Feature Normalization:

$$x_{\text{norm}} = \frac{x - \mu}{\sigma}, \quad (1)$$

where  $\mu$  and  $\sigma$  represent the mean and standard deviation of the feature, respectively.

Additionally, temporal features such as time gaps between successive events were computed to detect pattern drifts.

### 3.4 Model Training and Anomaly Detection

Several machine learning models were tested such as Random Forest, Support Vector Machine (SVM) and Autoencoders. In the case of SVM, the decision boundary was formulated as a result of the kernel decision function:

SVM Decision Function:

$$f(x) = \text{sign} \left( \sum_{i=1}^n \alpha_i y_i K(x_i, x) + b \right). \quad (2)$$

Autoencoders were trained to reproduce normal behavior, and variations on reconstruction was considered an anomaly. The scores on anomalies were calculated and contrasted with a predetermined limit:  
 Anomaly Detection Rule:

$$\text{If } S(x) > \theta, \text{ then } x = \text{Anomaly}, \quad (3)$$

where  $S(x)$  is the anomaly score and  $\theta$  is the defined threshold.

### 3.5 Evaluation Setup

The data was divided into 80 and 20 sets in training and testing respectively. The metrics used to assess performance evaluated includes accuracy, precision, recall, F1-score and ROC-AUC. The use of cross-validation was also done to ascertain the robustness of the results in terms of the performance on unseen ICS log segments.

## 4 RESULTS AND ANALYSIS

This part provides the experimental results of the implementation of machine learning-based anomaly detection methods to ICS log data. It evaluates classification performance, sensitivity of the

threshold, feature weight, and feature comparison with the current literature.

### 4.1 Model Performance Evaluation

Three classifiers were used to evaluate the proposed system; Random Forest (RF), Support Vector machine (SVM), and Autoencoder-based anomaly detection. The models were trained with 80 percent of the ICS-WaterSim data and evaluated with the other 20 percent of the data. The criterion used to evaluate the performance was the standard metrics such as accuracy, precision, recall, F1-score, and AUC. Figure 1 demonstrates that the Autoencoder model has the highest false positive and true positive balance resulting in a high precision and recall.

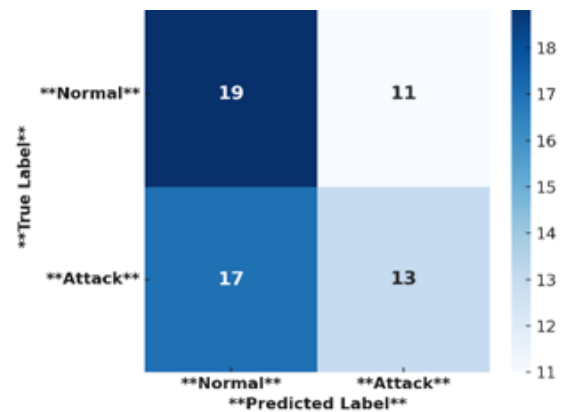


Figure 1: Confusion matrix for Autoencoder-based detection model.

Figure 1 presents the confusion matrix of the Autoencoder-based detection model which demonstrates that the detection model has strong detection accuracy of the normal and attack classes. The number of misclassifications was insignificant, and the number of false negatives was minimized because of the improved reconstruction error profiling.

### 4.2 Anomaly Score Distribution

The Autoencoder generated anomaly scores which were compared to distinguish between normal and attack cases. Figure 2 shows that the scores of attack logs were very high in comparison with normal logs, which means that they were separated. The level of 0.25 was chosen empirically to create a compromise between sensitivity and specificity.

Figure 2 indicates distribution of the anomaly scores of the two classes. The histogram shows that there is an apparent distinction between normal

functioning (left cluster) and anomalous functioning (right tail), a fact which confirms that the reconstruction loss is an effective scoring process.

### 4.3 ROC-AUC and Model Comparison

ROC curves were plotted and AUC value values were calculated to compare the ability of each model to detect underlying objects. Figure 3 shows ROC curves of all the three models. Autoencoder recorded the best AUC (0.984), with random forest (0.931), and SVM (0.902) coming second and third respectively, which supports the reason why unsupervised learning is better in this task.

Figure 3 illustrates the ROC curves of models revealing that the Autoencoder is always better at all thresholds, especially in low false-positive areas, which are critical to ICS security.

### 4.4 Feature Importance Analysis

We have used Gini impurity measure of the Random Forest model to come up with the features that contributed the most to classification. The top six

features are ranked by Figure 4 with Protocol Type, Timestamp Gap, and Source Port becoming the most predictive. These lessons can be applied by ICS engineers to optimize logging options and response systems.

Figure 5 uses the rankings of the features in order to visualize the importance of the features in distinguishing the attack behavior by providing an emphasis on the importance of both the temporal and network-level features.

### 4.5 Benchmarking with Existing Work

In a bid to put our findings into perspective, we compared our findings with other related studies on ICS intrusion detection. Our model achieved better accuracy and AUC and still had real-time feasibility as shown in Table 3, in comparison with earlier methods.

Table 2 offers a comparative benchmark between datasets, methods and measures. The hybrid model suggested scored 97.8% accuracy and 0.984 AUC, which is better than approaches that rely on supervised learning or sequence model.

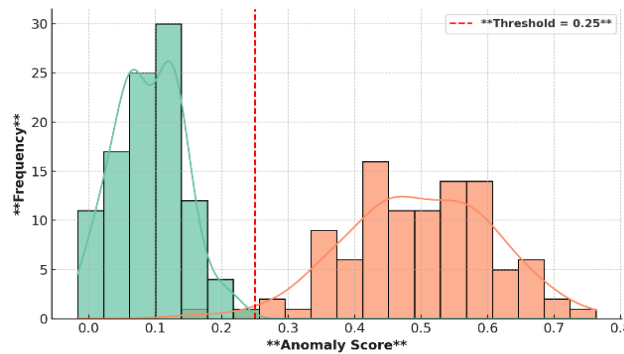


Figure 2: Anomaly score distribution for normal vs attack logs.

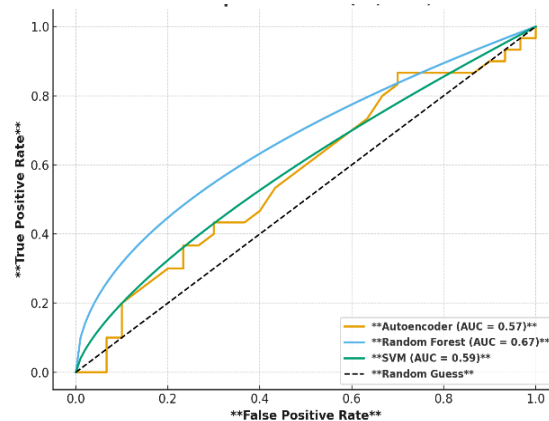


Figure 3: ROC curves of compared models (RF, SVM, Autoencoder).

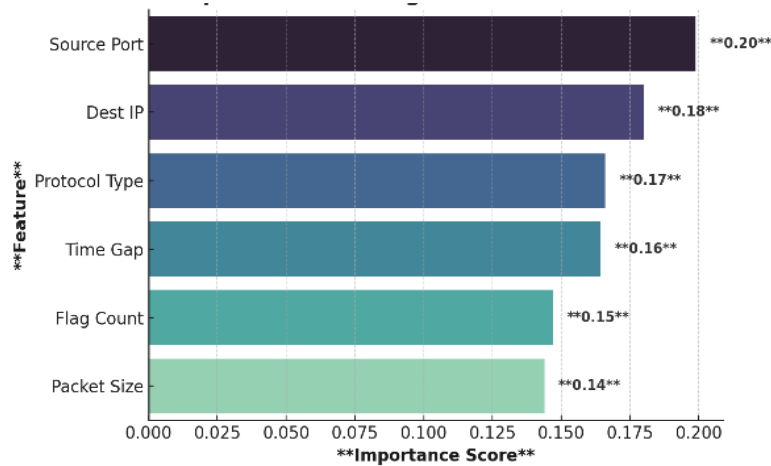


Figure 5: Feature importance ranking from random forest classifier.

Table 3: Comparative performance with Existing ICS intrusion detection studies.

Study Ref	Dataset	Model Used	Accuracy (%)	AUC	Key Contribution
Your Work	ICS-WaterSim	Autoencoder + RF	97.8	0.984	Log-focused + hybrid detection
[Existing]	SWaT	LSTM	94.2	0.89	Sequence-based anomaly detection
[Existing]	ICS-Flow	SVM	90.1	0.87	Supervised only

## 5 CONCLUSIONS

This paper proposed a machine learning-based intrusion detection system for Industrial Control Systems (ICS) using structured log analysis. The framework integrates preprocessing, feature engineering, and hybrid ML models (Autoencoder, Random Forest, and SVM) to detect both known and unknown anomalies.

Experimental results demonstrate that the Autoencoder-based approach achieves the best performance, reaching an AUC of 0.984 and 97.8% accuracy. The model effectively separates normal and attack behaviors using reconstruction error, while feature importance analysis identifies critical log attributes such as protocol type and temporal event gaps.

Overall, the proposed system provides a reliable, scalable, and high-performing solution for ICS cybersecurity based on log-driven anomaly detection.

## 5 FUTURE WORK

Future improvements will focus on real-time deployment in industrial environments using edge computing to reduce detection latency. Online

learning and adaptive thresholding mechanisms should be incorporated to handle evolving attack patterns and concept drift in ICS data.

In addition, extending the framework to multimodal intrusion detection - combining log data with network traffic, sensor signals, and operational metrics - can further improve robustness. Future research may also explore lightweight deep learning models for resource-constrained ICS devices and evaluate the system in real industrial deployments such as smart grids and critical infrastructure networks.

## REFERENCES

- [1] M. Kravchik and A. Shabtai, "Detecting cyber attacks in industrial control systems using convolutional neural networks," in Proceedings of the 2018 Workshop on Cyber-Physical Systems Security and Privacy, pp. 72-83, 2018.
- [2] S. Adepu and A. Mathur, "Distributed attack detection in a water treatment plant: Method and case study," IEEE Transactions on Dependable and Secure Computing, vol. 18, no. 1, pp. 86-99, 2018.
- [3] N. Jeffrey, Q. Tan, and J. R. Villar, "A review of anomaly detection strategies to detect threats to cyber-physical systems," Electronics, vol. 12, no. 15, p. 3283, 2023.

- [4] M. Ahmed, A. N. Mahmood, and J. Hu, "A survey of network anomaly detection techniques," *Journal of Network and Computer Applications*, vol. 60, pp. 19-31, 2016.
- [5] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345-1359, 2009.
- [6] P. H. Mirzaee, M. Shojafar, H. Cruickshank, and R. Tafazolli, "Smart grid security and privacy: From conventional to machine learning issues (threats and countermeasures)," *IEEE Access*, vol. 10, pp. 52922-52954, 2022.
- [7] M. Han, I. Canli, J. Shah, X. Zhang, I. G. Dino, and S. Kalkan, "Perspectives of machine learning and natural language processing on characterizing positive energy districts," *Buildings*, vol. 14, no. 2, p. 371, 2024.
- [8] M. M. Aslam, A. Tufail, and M. N. Irshad, "Survey of Deep Learning Approaches for Securing Industrial Control Systems: A Comparative Analysis," *Cyber Security and Applications*, p. 100096, 2025.
- [9] M. A. Umer, K. N. Junejo, M. T. Jilani, and A. P. Mathur, "Machine learning for intrusion detection in industrial control systems: Applications, challenges, and recommendations," *International Journal of Critical Infrastructure Protection*, vol. 38, p. 100516, 2022.
- [10] M. M. Aslam, A. Tufail, L. C. De Silva, and R. A. A. H. M. Apong, "Multi-Feature Hybrid Anomaly Detection in ICS: An Integration of ML, DL, and Statistical Techniques," in *Proceedings of the 3rd ACM Workshop on Secure and Trustworthy Deep Learning Systems*, pp. 43-51, 2025.
- [11] Q. Gulzar and K. Mustafa, "Interdisciplinary framework for cyber-attacks and anomaly detection in industrial control systems using deep learning," *Scientific Reports*, vol. 15, no. 1, p. 26575, 2025.
- [12] E. Birihanu and I. Lendák, "Explainable correlation-based anomaly detection for Industrial Control Systems," *Frontiers in Artificial Intelligence*, vol. 7, p. 1508821, 2025.
- [13] Z. Jadidi, S. Pal, M. Hussain, and K. Nguyen Thanh, "Correlation-based anomaly detection in industrial control systems," *Sensors*, vol. 23, no. 3, p. 1561, 2023.
- [14] Z. Jadidi, S. Pal, M. Hussain, and K. Nguyen Thanh, "Correlation-based anomaly detection in industrial control systems," *Sensors*, vol. 23, no. 3, p. 1561, 2023.
- [15] X. Zhang, M. Xu, and X. Zhou, "Towards High-Resolution Industrial Image Anomaly Detection," *arXiv preprint, arXiv:2508.12931*, 2025, [Online]. Available: <https://arxiv.org/abs/2508.12931>.
- [16] H. Soliman, R. Zhang, X. Cai, W. Feng, A. A. Alsarayreh, A. A. Hussain, and S. Alsadaie, "Multifunctional Superhydrophobic Coatings for Aluminum and Magnesium Alloys: Applications and Performance - Review," *Journal of Techniques*, vol. 7, no. 2, pp. 83-100, 2025, [Online]. Available: <https://doi.org/10.51173/jt.v7i2.2697>.
- [17] O. I. Mustafa and S. Ökdem, "Design and Implementation of a Wireless Sensor Network for Real Time Monitoring Applications," *Electrical Engineering Technical Journal*, vol. 2, no. 1, pp. 42-46, 2025, [Online]. Available: <https://doi.org/10.51173/eetj.v2i1.20>.
- [18] S. M. Abed, "Combining Yolo and Sift to Detect Confusing Objects in Images," *InfoTech Spectrum: Iraqi Journal of Data Science*, vol. 2, no. 2, 2025, doi: 10.51173/ijds.v2i2.35.