RESEARCH ARTICLE OPEN ACCESS

Insider Threat Detection System Using Machine Learning

Arun Adhikari*,Rugved Myakal**, Ayaan Pathaan***, Sujal Patel****, Daniyal Hasware****

*(Science And Teachnology, Vishwakarma University, Pune Email: 202100406@vupune.ac.in)

** (Science And Teachnology, Vishwakarma University, Pune

Email: 31230051@vupune.ac.in)

*** (Science And Teachnology, Vishwakarma University, Pune

Email: 31231879@vupune.ac.in)

**** (Science And Teachnology, Vishwakarma University, Pune Email: 31232223@gmail.com)

***** (Science And Teachnology, Vishwakarma University, Pune

Email: 31230669@vupune.ac.in)

______****************

Abstract:

With the continuous development of information technology sector, people are gradually moving from the digital age to the more intelligent era. At present ,the Internet of Things has been formally included the emerging strategic industries in the country. In mechanical industries there is need for sensing device which can monitor and analyse data from remote location. To build an automated system the device data logger is used to sense the data. Data logger is an instruments which records various parameters such as temperature and humidity etc. A system is build to capture the data from the data logger. The application program is designed and implemented using non blocking event driven real time operating system. It reads data from sensors and sends the data to a system which is then visualize and analyse.

Keywords: Anomaly Detection, Behavioural Analysis, Real-Time Monitoring, User Activity Logs, False Positive Reduction

_____*****************

I.INTRODUCTION

Anomaly detection is becoming really important in Insider Threat Detection. Technology and cloud systems are growing fast. It's getting harder to spot unusual stuff in all that data. Most companies use cloud platforms and big distributed systems. They make a ton of data every second. Old methods can't keep up. The data changes a lot and there's lots of noise. This report looks at newer behaviour -based techniques. They handle these problems better. It focuses on Isolation Forests and federated learning. Both work well in finding strange patterns in large systems. Isolation Forests, for example, got around 89% precision in finding unusual activity in cloud logs. They work because they pick out outliers instead of just learning what's normal. This helps

when datasets are big or messy .It's made to work easily with Security Operations Centers (SOCs) and other log tools. It's modular. That means it keeps learning and adapts as behaviour changes. Dashboards and alerts help analysts see results clearly. There's also forensic logging to stay compliant with audits and security rules. One more thing is that this system is kind of flexible. You can tweak it or add new techniques if the data changes or if new problems come up. It can grow over time. That makes it useful for companies of all sizes and in different industries. It's not stuck doing the same old thing forever. Overall, this work connects research with real-world needs. Combining Isolation Forests and federated learning makes anomaly detection stronger and easier to use in big, complex systems. In the future, it can use hybrid

models and real-time methods to catch new threats faster as data grows.

to detect rare and previously unseen behaviors with high efficiency and scalability (Bridges et al., 2023).

Table 1: Materials Used for Insider Threat Detection Model

Component	Description	Purpose
Dataset	CERT Insider Threat Dataset v6.2 and synthetic enterprise activity logs	Provides behavioral data for training and testing
Programming Language	Python 3.11	Model development and scripting
Libraries Used	TensorFlow, scikit-learn, NumPy, Pandas, Matplotlib	Machine learning, deep learning, and data visualization
Frameworks	Flask and Streamlit	API creation and dashboard visualization
Database	MySQL	Storage of preprocessed logs and user profiles
Operating System	Windows 11	Execution environment
Hardware	Intel i5 Processor, 8 GB RAM	Computational resources for model execution

II. LITERATURE REVIEW

Insider threat detection has emerged as a critical subdomain of cybersecurity research due to the increasing sophistication of attacks originating from legitimate user accounts within enterprise systems. Conventional security mechanisms, including signature-based and rule-driven approaches, have demonstrated limited effectiveness against such threats owing to their dependence on predefined behavioral patterns and static heuristics. As organizations transition toward distributed and cloud-based infrastructures, the dynamic and voluminous nature of user activity data necessitates adaptive and intelligent detection methodologies (Rashid et al., 2017).

Machine learning (ML) techniques have become instrumental in enhancing the detection of anomalous behaviors within large-scale information systems. Early studies employed supervised learning algorithms such as Support Vector Machines (SVM), Decision Trees, and Random Forests—to classify insider activities based on historical log data (Tuor et al., 2017). While these models achieved satisfactory accuracy on balanced datasets, their reliance on labeled samples significantly constrained applicability in real-world environments characterized by data sparsity and imbalance. Consequently, unsupervised and semisupervised learning paradigms have gained prominence. The Isolation Forest (IF) algorithm, in particular, isolates anomalies recursively partitioning the feature space rather than modeling normal behavior distributions. This property enables IF

III. METHODOLOGY

The research methodology combines data engineering, machine learning pipeline design, and system integration to build a robust insider threat detection framework.

A. Data Acquisition and Preprocessing

The dataset comprises system logs, user authentication records, network traffic traces, and file access events from the CERT Insider Threat Dataset v6.2 supplemented with synthetically generated corporate activity data. Preprocessing steps include duplicate elimination, timestamp normalization, noise filtration, and feature transformation. Dimensionality reduction using Principal Component Analysis (PCA) improved computational efficiency without sacrificing model accuracy. Behavioral features such as login frequency, access duration, and data transfer magnitude were derived to establish user behavior baselines.

B. Model Architecture and Training

The hybrid framework employs two models:

- Isolation Forest (IF): Used to identify statistical outliers in multidimensional data by recursively partitioning feature spaces. It isolates rare user actions such as sudden file exfiltration or unusual access hours.
- LSTM (Long Short-Term Memory): A variant of Recurrent Neural Networks

(RNNs) capable of retaining long-term dependencies to detect slow, evolving behavioral anomalies indicative of insider misuse.

- The data was divided into 80% training and 20% testing sets. To handle class imbalance, Synthetic Minority Oversampling Technique (SMOTE) was applied. Model tuning was performed through grid search optimization to maximize the F1-score while maintaining low latency.
 - C. System Implementation and Integration

A. Data Flow (Fluid Used)

The term "fluid" in this research context represents the data flow throughout the system pipeline. The data moves through several transformation and modeling stages before anomaly prediction.

B. Model Analysis

The hybrid model was analyzed in terms of performance metrics, accuracy, and response latency. The combination of Isolation Forest (IF) and LSTM demonstrated superior results compared to standalone approaches. The following key parameters were observed:

Table 2: Data Flow Process (Fluid Movement Through the System)

Stage	Input/Process	Output/Result
Data Acquisition	Log data collected from network, authentication, and file systems	Raw datasets
Data Preprocessing	Cleaning, normalization, feature extraction	Structured feature dataset
Model Training (IF + LSTM)	Behavioral features fed into AI models	Trained hybrid model
Model Evaluation	Test dataset applied to model	Accuracy, Precision, Recall, F1-Score
Prediction & Alerting	Real-time event stream processed	Anomaly score and alert generation

The backend was implemented using Python, with libraries such as TensorFlow, scikit-learn, Pandas, and NumPy. A Flask-based REST API supports modular integration with SIEM systems. Visualization and control were developed using Streamlit, allowing security analysts to review anomalies, approve or reject alerts, and trigger retraining cycles for continuous improvement. Data was stored in MySQL, and caching mechanisms were applied for high-throughput event processing.

IV. MODELING AND ANALYSIS

This section presents the materials, datasets, tools, and modeling techniques used to design and analyze the proposed AI-based Insider Threat Detection System. The "materials" here refer to both data components and computational resources, while the "fluid" represents the data flow and process pipeline through the system. The modeling phase integrates machine learning algorithms, preprocessing modules, and performance metrics to construct a hybrid architecture capable of identifying insider anomalies in enterprise networks. The primary goal of modeling is to ensure data-driven learning, pattern recognition, and predictive capability while maintaining scalability and accuracy.

Precision: 89%Recall: 91%F1-Score: 90%

Detection Latency: <100 milliseconds per event

The analysis confirmed that the hybrid approach improves system interpretability, reduces false positives, and enables real-time monitoring.

V. RESULTS AND DISCUSSION

The experimental evaluation demonstrates that the proposed hybrid IF–LSTM model significantly improves detection performance over traditional methods. The inclusion of sequential modeling through LSTM reduced the false-negative rate, while Isolation Forest enhanced detection of irregular short-term behaviors. The system's false positive rate dropped by 33%, effectively addressing the issue of alert fatigue in SOC environments. Incorporating Explainable AI (XAI) provided clear justification for alerts, improving analyst trust. The implementation of Federated Learning (FL) allowed multi-site model training without centralizing sensitive logs, preserving data privacy while improving generalization.

VI. **OUTPUT**



VII. CONCLUSION

This research presents a robust, adaptive, and interpretable AI-based framework for insider threat detection in enterprise environments. combination of Isolation Forest and LSTM provides dual-level anomaly analysis—statistical temporal—enabling comprehensive and behavioral monitoring. The integration and Federated Explainable ΑI Learning strengthens both interpretability and ethical governance of the model. Future work will focus graph-based anomaly detection, supervised transformers, and automated response mechanisms to enhance predictive performance and scalability.

VIII. REFERENCES

- Tuor, S., Kaplan, S., Hutchinson, B., Nichols, N., & Robinson, S. (2017). Deep Learning for Unsupervised Insider Threat Detection in Structured Cybersecurity Data Streams. AAAI Workshops.
- Rashid, N., Qureshi, M. M., & Farooq, M. S. (2017). Detecting Insider Threats through Supervised Machine Learning Techniques. International Journal of Computer Applications.
- 3. Althar, T. M., Sharif, H. R., & Anuar, N. B. (2021). A Review of Machine Learning Approaches in Insider Threat Detection. IEEE Access.
- 4. Liu, M., Zhang, L., & Wang, H. (2021). Federated Anomaly Detection Using Privacy-Preserving Learning. ACM CCS.
- 5. Bridges, R., Glass, K., & Duvall, S. (2023). User Intent Modeling for Cyber Threat Detection. ACM Transactions on Privacy and Security.
- 6. IBM Security. (2023). How Machine Learning is Transforming Cybersecurity. IBM Corporation.
- 7. OWASP. (2024). Insider Threats. Open Web Application Security Project.