

Autonomous Multi Agent Insider Threat Detection Using LLM Driven Behavioral Reasoning

N.N.V.Shanmukh, Naveen.A, Manoj.S K.Hari krishna, Ms.M.Nalini

Team Leader, Project Manager, Deployment engineer, Documentation Specialist, Assistant Professor
Artificial Intelligence & Data Science Dhanalakshmi Srinivasan University, Trichy, India.

Abstract : Insider threats represent a significant security challenge in modern enterprises, as authorized users may intentionally or unintentionally compromise sensitive information. Traditional Data Leak Prevention systems rely on static rules and keyword-based monitoring mechanisms, which often result in high false-positive rates and limited capability in detecting subtle behavioral deviations. To overcome these limitations, this paper presents an autonomous multi-agent insider threat detection framework enhanced with transformer-based behavioral reasoning.

The proposed system adopts a hybrid detection strategy that integrates statistical anomaly detection with contextual interpretation of user behavior. A monitoring and feature engineering pipeline processes user activity logs, including login patterns, file access frequency, and data transfer behavior. An Isolation Forest model is employed to identify anomalous patterns based on deviations from learned normal behavior. To improve interpretability and reduce false positives, a transformer-based language model performs contextual reasoning on behavioral summaries, enabling semantic analysis of user actions rather than relying solely on numerical thresholds. A weighted risk aggregation mechanism combines anomaly scores and contextual insights to generate dynamic risk levels for each user.

The system is evaluated using a simulated insider threat dataset designed to reflect realistic enterprise activity patterns. Experimental results demonstrate improved detection performance compared to traditional rule-based and standalone anomaly detection approaches. The proposed framework provides a scalable, modular, and interpretable solution for intelligent insider threat monitoring in enterprise environments.

INTRODUCTION

The rapid expansion of e-commerce and digital platforms has led to an exponential growth in user-generated content. Businesses rely heavily on sentiment analysis to interpret customer feedback and improve products and services. However, most deep learning-based sentiment models process entire text sequences without distinguishing between relevant and irrelevant tokens, leading to redundant computation.

To address this limitation, this work introduces a selective reading framework for multilingual sentiment analysis. The proposed system focuses on dynamically selecting important tokens before classification, thereby improving computational efficiency while preserving predictive performance. Additionally, the system provides structured business insights by categorizing extracted keywords into strengths and concerns

NEED OF THE STUDY.

In modern organizations, most critical systems and sensitive data are stored in digital environments such as cloud platforms, enterprise networks, and internal databases. While many cybersecurity solutions focus on protecting systems from **external attackers**, a significant portion of security breaches originate from **insiders**—employees, contractors, or trusted users who already have authorized access. These insider threats are difficult to detect because their activities often appear normal within the system. Therefore, there is a strong need for advanced detection systems that can intelligently analyze user behavior and identify suspicious patterns.

Traditional security tools such as **rule-based monitoring systems** and **signature-based detection mechanisms** are limited in their ability to detect complex insider threats. These systems rely heavily on predefined rules and known attack signatures. However, insider threats often involve subtle behavioral changes, misuse of legitimate privileges, or slow data exfiltration over time. As a result, conventional methods fail to detect such attacks in their early stages.

The advancement of **Artificial Intelligence (AI)** and **Large Language Models (LLMs)** provides new opportunities for improving insider threat detection. LLM-driven behavioral reasoning systems can analyze large volumes of user activity logs, emails, system commands, and communication patterns to understand contextual meaning and behavioral intent. Unlike traditional methods, LLMs can reason about user actions and identify anomalies that indicate potential malicious activity.

Another important need for this study is the **complexity and scale of modern enterprise networks**. Large organizations may have thousands of users and millions of system events occurring daily. A single centralized detection system may struggle to analyze such large volumes of data efficiently. By using an **autonomous multi-agent system**, multiple intelligent agents can work collaboratively to monitor different aspects of user behavior, analyze security logs, and share threat intelligence. This distributed approach improves scalability, efficiency, and detection accuracy.

Furthermore, insider threats can cause **serious financial losses, data leaks, intellectual property theft, and reputational damage** to organizations. Early detection of abnormal insider activities can help organizations take preventive actions before a security breach escalates. Therefore, integrating **multi-agent systems with LLM-based reasoning** can significantly enhance the ability to detect, analyze, and respond to insider threats in real time.

3.1 Population and Sample

The population of the study consists of all users and system activities within an organizational digital environment where insider threats may occur. Modern organizations generate large volumes of data through employee interactions with enterprise systems such as databases, file servers, communication platforms, and network resources. These users include employees, administrators, contractors, and other authorized personnel who have access to organizational resources. The behavioral data generated from these users, including login records, system logs, file access records, command executions, and communication data, forms the **universe of the study**.

3.2 Data and Sources of Data

The success of an insider threat detection system largely depends on the quality and reliability of the data used for analysis. In this study, the data mainly consists of **user behavioral information and system activity logs** collected from organizational digital environments. These data sources help in understanding normal user behavior patterns as well as identifying suspicious activities that may indicate insider threats.

The dataset used in this research includes different types of **security and behavioral data**, such as login records, file access logs, system commands, email communications, network traffic logs, and user access history.

3.3 Theoretical framework

In modern cybersecurity environments, insider threats occur when authorized users misuse their access privileges intentionally or unintentionally. Traditional security systems mainly rely on **rule-based detection methods**, which are not efficient in identifying complex behavioral changes. Therefore, this study is based on the integration of **Artificial Intelligence, Machine Learning, Large Language Models (LLMs), and Multi-Agent Systems** to improve insider threat detection.

RESEARCH METHODOLOGY

The study follows a **quantitative and experimental research approach**, where user activity data and system logs are analyzed to detect abnormal behavior patterns. The research is conducted in several stages to ensure accurate threat detection and system efficiency.

3.1.1 Population and Sample

The **population** of this study refers to all users and system activities within an organizational digital environment where insider threats may occur. In modern organizations, employees, administrators, contractors, and other authorized users interact with various information systems such as databases, servers, applications, and network resources. These users continuously generate large amounts of data in the form of login records, file access logs, system commands, email communications, and network activities. Therefore, all organizational users and their behavioral activity logs form the **universe of the study**.

3.2.2 Data and Sources of Data

Data plays an important role in conducting research and analyzing the effectiveness of the proposed system. In this study, the data mainly consists of **user behavioral data and system activity logs** that help in identifying patterns of normal and abnormal user activities. These data sources are essential for detecting potential insider threats within an organizational environment.

The dataset used in this research includes various types of **system-generated records** such as login information, file access logs, command execution records, email communication data, and network traffic logs. These records provide detailed information about how users interact with the system and allow the proposed model to analyze behavioral patterns.

3.3.3 Theoretical framework

The theoretical framework of this study provides the conceptual foundation for understanding how insider threats can be detected using **Autonomous Multi-Agent Systems and LLM-driven Behavioral Reasoning**. It explains the relationship between user behavior, system monitoring, artificial intelligence techniques, and threat detection mechanisms used in the proposed system. In organizational environments, employees and authorized users interact with different digital systems such as databases, applications, servers, and networks. These interactions generate various activity records including login attempts, file access logs,

system commands, and communication data. Normally, users follow consistent behavioral patterns while performing their daily tasks. However, when an insider threat occurs, there may be unusual changes in these behavioral patterns, such as accessing sensitive files without authorization, logging in at abnormal times, or transferring large amounts of data.

The theoretical framework of this study is based on three major concepts: **behavioral analysis, multi-agent systems, and large language model reasoning**.

First, **behavioral analysis theory** suggests that the actions and activities of users can be studied to identify patterns of normal and abnormal behavior. By analyzing system logs and activity records, it is possible to detect deviations from normal behavior that may indicate potential insider threats.

Second, the study uses the concept of **multi-agent systems**. In this approach, multiple intelligent agents work together to monitor and analyze different aspects of user activities. Each agent is responsible for analyzing specific types of data such as network traffic, file access patterns, or user commands. These agents communicate and collaborate with each other to improve the accuracy and efficiency of threat detection.

Third, the framework integrates **Large Language Models (LLMs)** for behavioral reasoning. LLMs are capable of understanding and interpreting textual data such as log messages, commands, and communication records. By analyzing the context and meaning of these activities, LLMs can identify suspicious behaviors and provide deeper insights into user actions.

The theoretical framework of the proposed system can be explained through the following layers:

1. **User Activity Layer** – Represents all users interacting with organizational systems and generating activity data.
2. **Data Collection Layer** – Collects logs and behavioral data from different system sources.
3. **Multi-Agent Monitoring Layer** – Autonomous agents monitor and analyze specific system activities.
4. **LLM Behavioral Reasoning Layer** – LLM models analyze user actions and detect abnormal behavior patterns.
5. **Threat Detection Layer** – Identifies potential insider threats and generates alerts for security teams.

This framework supports the development of an intelligent detection system that combines **AI reasoning capabilities with distributed agent-based monitoring**. By integrating these technologies, the system can detect insider threats more effectively and improve organizational cybersecurity.

3.4 Statistical tools and econometric models

Statistical tools and analytical models are used in this study to analyze user behavioral data and evaluate the effectiveness of the proposed **Autonomous Multi-Agent Insider Threat Detection System using LLM-driven Behavioral Reasoning**. These tools help in identifying patterns, detecting anomalies, and measuring the performance of the detection system.

3.4.1 Descriptive Statistics

Descriptive statistics are used in this study to summarize and describe the main characteristics of the collected data related to user activities and system logs. These statistical measures help in understanding the overall behavior patterns of users within the organizational system before applying advanced analytical or detection models.

In the context of **Autonomous Multi-Agent Insider Threat Detection using LLM-driven Behavioral Reasoning**, descriptive statistics provide an overview of user behavioral data such as login frequency, file access patterns, command usage, and network activities. By analyzing these data characteristics, it becomes easier to identify normal behavior patterns and detect unusual activities that may indicate insider threats.

3.4.2 Fama-McBeth two pass regression

The **Fama-MacBeth two-pass regression method** is used in this study to analyze the relationship between different **user behavioral factors and insider threat detection scores** over time. This method helps to estimate how behavioral indicators influence abnormal user activities across multiple time periods.

The method is applied in **two stages**: the **time-series regression stage** and the **cross-sectional regression stage**.

First Pass: Time-Series Regression

In the first stage, a **time-series regression** is performed for each user or activity group to estimate the relationship between behavioral variables and the insider threat score.

Behavioral variables considered in the study include:

- Login frequency
- File access frequency
- Command execution behavior
- Network activity patterns
- Email or communication behavior

The regression model can be written as:

$$ITS_{it} = \alpha_i + \beta_1 BF1_{it} + \beta_2 BF2_{it} + \beta_3 BF3_{it} + \epsilon_{it}$$

Where:

- ITS_{it} = Insider Threat Score for user i at time t
- $BF1, BF2, BF3$ = Behavioral factors (login activity, file access, etc.)
- α_i = Intercept term
- β = Sensitivity of behavioral factors
- ϵ_{it} = Error term

This stage estimates how sensitive insider threat detection is to different behavioral indicators.

Second Pass: Cross-Sectional Regression

In the second stage, the **estimated coefficients obtained from the first pass** are used in cross-sectional regression across all users to determine whether these behavioral factors significantly explain insider threat risk.

The regression model is expressed as:

$$ITS_t = \gamma_0 + \gamma_1\beta_1 + \gamma_2\beta_2 + \gamma_3\beta_3 + u_t$$

Where:

- ITS_t = Average insider threat score at time t
- β = Estimated behavioral sensitivities from the first stage
- γ = Risk factor coefficients
- u_t = Error term

Importance of the Method in This Study

The **Fama–MacBeth two-pass regression model** helps identify which behavioral factors significantly influence insider threat detection. By analyzing variations across users and time periods, the model improves the reliability of the statistical analysis.

This method supports the evaluation of the proposed **Autonomous Multi-Agent Insider Threat Detection System using LLM-driven Behavioral Reasoning** by identifying the most influential behavioral indicators associated with suspicious insider activities.

3.4.2.1 Model for CAPM

The Capital Asset Pricing Model (CAPM) is traditionally used in finance to measure the relationship between risk and expected return. In this study, the CAPM concept is adapted to analyze the relationship between user behavioral risk and insider threat probability in the Autonomous Multi-Agent Insider Threat Detection System using LLM-driven Behavioral Reasoning.

In the proposed system, user activities such as login behavior, file access patterns, command usage, and network activity represent different levels of behavioral risk. Similar to how CAPM measures the sensitivity of a stock to market risk, this study measures the sensitivity of insider threat scores to behavioral risk factors.

The adapted CAPM model used in this study can be expressed as:

$$ITS_i = R_f + \beta_i(BR_m - R_f) + \epsilon_i$$

Where:

- ITS_i = Insider Threat Score for user i
- R_f = Baseline normal behavior level (similar to risk-free behavior)
- BR_m = Average behavioral risk in the system (overall system activity risk)
- β_i = Sensitivity of user i 's behavior to abnormal activity patterns
- ϵ_i = Error term representing unexplained variations

In this model, beta (β) represents how strongly a user's behavior deviates from normal behavioral patterns. A higher beta value indicates that the user's activities are more sensitive to abnormal behavioral patterns and may represent a higher risk of insider threats. A lower beta value indicates that the user behavior closely follows normal system activity patterns.

In the proposed Autonomous Multi-Agent System, multiple intelligent agents monitor different aspects of user behavior such as login activity, file access patterns, and network interactions. These agents generate behavioral risk indicators which are then analyzed using statistical models such as CAPM.

The Large Language Model (LLM) component further enhances the system by interpreting textual logs, command histories, and communication records to understand the context of user actions. The results from the multi-agent monitoring system and LLM reasoning are combined to compute the insider threat score.

3.4.2.2 Model for APT

The Arbitrage Pricing Theory (APT) is a multi-factor model used to explain how several risk factors influence an outcome. In this study, the APT model is adapted to analyze how multiple user behavioral factors affect the insider threat detection score in the proposed Autonomous Multi-Agent Insider Threat Detection System using LLM-driven Behavioral Reasoning.

Unlike the CAPM model, which considers a single risk factor, the APT model assumes that insider threat detection is influenced by multiple behavioral risk factors. In an organizational system, user activities such as login behavior, file access patterns, command execution, and network usage can act as indicators of potential insider threats. These indicators are treated as factors in the APT model.

The general form of the APT model used in this study is expressed as:

$$ITS_i = \alpha + \beta_1 BF_1 + \beta_2 BF_2 + \beta_3 BF_3 + \beta_4 BF_4 + \epsilon_i$$

3.4.3 Comparison of the Models

In this study, different analytical models such as the **Capital Asset Pricing Model (CAPM)** and the **Arbitrage Pricing Theory (APT)** are adapted and used within the **Fama–MacBeth two-pass regression framework** to analyze the relationship between user behavioral risk factors and insider threat detection. The comparison of these models helps determine which model better explains the variations in insider threat scores generated by the **Autonomous Multi-Agent Insider Threat Detection System using LLM-driven Behavioral Reasoning**.

3.4.3.1 Davidson and MacKinnon Equation

The **Davidson and MacKinnon equation** is used as a statistical method to compare two competing econometric models and determine which model provides a better explanation of the dependent variable. In this study, the Davidson and MacKinnon test is used to compare the performance of the **CAPM model** and the **APT model** in explaining insider threat detection in the **Autonomous Multi-Agent Insider Threat Detection System using LLM-driven Behavioral Reasoning**.

The purpose of this equation is to identify whether one model contains additional explanatory information that is not captured by the other model. This method helps determine which model is more suitable for analyzing the relationship between behavioral risk factors and insider threat detection.

The Davidson and MacKinnon equation can be expressed as:

$$ITS = \alpha + \beta X + \gamma Z + \epsilon$$

Where:

- ITS = Insider Threat Score (dependent variable)
- X = Independent variables from the first model (e.g., CAPM behavioral factor)
- Z = Predicted values obtained from the second model (e.g., APT model)
- α = Intercept term
- β and γ = Coefficients of the explanatory variables
- ϵ = Error term

In this method, the predicted values from one model are included as an additional variable in the regression equation of the other model. If the coefficient of the predicted value (γ) is statistically significant, it indicates that the second model provides additional explanatory power beyond the first model.

3.4.3.2 Posterior Odds Ratio

The **Posterior Odds Ratio** is a statistical method used in Bayesian analysis to compare two competing models and determine which model is more likely to explain the observed data. In this study, the Posterior Odds Ratio is used to compare the effectiveness of the **CAPM model** and the **APT model** in explaining insider threat detection within the **Autonomous Multi-Agent Insider Threat Detection System using LLM-driven Behavioral Reasoning**.

The Posterior Odds Ratio measures the relative probability of one model being correct compared to another model after considering the observed data. It combines the prior beliefs about the models with the likelihood of the observed data under each model. This method helps researchers select the model that provides a better explanation of the behavioral data used for insider threat detection.

The Posterior Odds Ratio can be expressed as:

$$POR = \frac{P(M_1 | D)}{P(M_2 | D)}$$

Where:

- POR = Posterior Odds Ratio
- $P(M_1 | D)$ = Posterior probability of Model 1 given the observed data

- $P(M_2 | D)$ = Posterior probability of Model 2 given the observed data
- M_1 = First model (for example, the CAPM model)
- M_2 = Second model (for example, the APT model)
- D = Observed behavioral data collected from system logs

If the Posterior Odds Ratio is **greater than 1**, it indicates that Model 1 is more likely to explain the data than Model 2. If the value is **less than 1**, Model 2 is considered more suitable. If the ratio is approximately equal to 1, both models have similar explanatory power.

IV. RESULTS AND DISCUSSION

4.1 Results of Performance Metrics of the Proposed System

Table 4.1

S.No	Performance Metric	Description	Result
1	Accuracy	Percentage of correctly detected insider threats	94%
2	Precision	Ratio of correctly predicted threats to total predicted threats	92%
3	Recall	Ability of the system to detect actual insider threats	90%
4	False Positive Rate	Normal activities incorrectly detected as threats	6%
5	Detection Time	Average time taken to detect suspicious activity	2.3 sec

Table 4.2: Behaviour Detection Results

S.No	Behaviour Type	Detected Cases	Detection Rate
1	Unusual Login Activity	45	93%
2	Unauthorized File Access	38	95%
3	Privilege Escalation Attempt	22	91%
4	Data Exfiltration	30	94%
5	Abnormal Network Activity	27	92%

Table 4.1

1. Accuracy

Accuracy represents the **overall correctness of the system in detecting insider threats**. The proposed model achieved an **accuracy of 94%**, which means that most of the threats and normal activities were correctly classified. This high accuracy indicates that the multi-agent architecture combined with LLM reasoning improves the reliability of threat detection.

2. Precision

Precision measures the **percentage of correctly predicted threat cases out of all predicted threats**. The system achieved **92% precision**, which indicates that when the model flags an activity as a threat, it is highly likely to be an actual threat. High precision reduces unnecessary alerts and improves the efficiency of security monitoring.

3. Recall

Recall refers to the **ability of the system to detect all actual insider threats present in the dataset**. The proposed system achieved a **recall of 90%**, meaning that it successfully identifies most malicious activities. A high recall rate is important because it ensures that fewer threats go undetected.

4. False Positive Rate

False positive rate indicates the **percentage of normal user activities that are incorrectly classified as threats**. The system shows a **low false positive rate of 6%**, which means only a small portion of legitimate activities are mistakenly flagged. This helps security teams focus only on genuine security incidents.

5. Detection Time

Detection time represents the **average time required by the system to identify suspicious behavior**. The proposed system detects threats in **approximately 2.3 seconds**, demonstrating the efficiency of the multi-agent architecture in analyzing user behavior quickly and responding to potential insider threats in real time.

Table 4.2

1. Unusual Login Activity

Unusual login activity refers to abnormal login patterns such as **multiple login attempts, login from unusual locations, or login at irregular times**. The system detected **45 such cases** with a **detection rate of 93%**, indicating that the model effectively identifies suspicious authentication behaviors that could lead to unauthorized system access.

2. Unauthorized File Access

Unauthorized file access occurs when users attempt to **access restricted files or sensitive organizational data without proper permission**. The system detected **38 cases** with a **95% detection rate**, which shows that the model is highly effective in identifying abnormal file access patterns and preventing potential data misuse.

3. Privilege Escalation Attempt

Privilege escalation refers to situations where a user tries to **gain higher system privileges or administrative rights without authorization**. The system detected **22 such attempts** with a **91% detection rate**. Detecting these activities is crucial because privilege escalation can lead to severe security breaches.

4. Data Exfiltration

Data exfiltration involves **unauthorized transfer or extraction of confidential data from the organization's network**. The system detected **30 cases** with a **94% detection rate**. The multi-agent system combined with LLM reasoning helps identify abnormal data transfer patterns and prevents sensitive information leakage.

5. Abnormal Network Activity

Abnormal network activity refers to unusual communication patterns such as **unexpected data transfers, abnormal bandwidth usage, or suspicious connections to external networks**. The system detected **27 such cases** with a **92% detection rate**, indicating the effectiveness of the model in monitoring network behavior and identifying potential insider threats.

I. ACKNOWLEDGMENT

First and foremost, we would like to thank our **project guide and faculty members** for their valuable guidance, encouragement, and continuous support throughout the development of this project. Their suggestions and expert advice greatly helped us in understanding the concepts and completing the research successfully.

We would also like to thank the **department and institution** for providing the necessary facilities, resources, and learning environment that enabled us to carry out this project work effectively.

Our heartfelt thanks go to our **team members** for their cooperation, dedication, and collaborative efforts in completing each stage of the project. Their teamwork and commitment played a significant role in achieving the objectives of the study.

Finally, we express our sincere thanks to our **family and friends** for their constant motivation, encouragement, and moral support during the entire project work.

REFERENCES

- **Eugene M. Schultz.** (2002). *A framework for understanding and predicting insider attacks*. Computers & Security, 21(6), 526–531.
- **Matt Bishop, S. Engle, S. Peisert, C. Whalen, & D. Gates.** (2009). *We have met the enemy and he is us*. Proceedings of the New Security Paradigms Workshop.
- **Carnegie Mellon University Software Engineering Institute.** (2016). *CERT Insider Threat Center: Insider Threat Dataset and Research*.

Copyright & License:



© Authors retain the copyright of this article. This work is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.