



International Journal of Multidisciplinary Research in Science, Engineering and Technology

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)



Impact Factor: 8.206

Volume 8, Issue 11, November 2025

| www.ijmrset.com | Impact Factor: 8.206 | ESTD Year: 2018 |



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

Ethical and Technical Framework for AI-Driven Insider Threat Detection

Jyothiba Salunkhe¹, Karthik M²

Assistant Professor, Dept. of Computer Science & Application, The Oxford College of Science, Bangalore, India¹ PG Student [MCA], Dept. of Computer Applications and Science, The Oxford College of Science, Bangalore, India²

ABSTRACT: The rapid adoption of Artificial Intelligence (AI) in cybersecurity has revolutionized threat detection mechanisms, yet its use in insider threat detection raises complex ethical, legal, and technical challenges. Traditional approaches focus on algorithmic efficiency while neglecting privacy, fairness, and accountability. This research proposes a comprehensive Ethical and Technical Framework for AI-Driven Insider Threat Detection (ETF-AITD), integrating principles of ethical AI with advanced detection architectures. The framework introduces four interlinked layers: Data Governance, AI Transparency, Human Oversight, and Adaptive Mitigation. Using a mixed-method approach—literature synthesis, system modeling, and simulated evaluation—this study explores how ethical principles can be systematically encoded within technical implementations. Findings reveal that incorporating fairness constraints, privacy-preserving analytics, and explainable AI (XAI) reduces false positives by 18% and improves trust perception among stakeholders by 35% in simulated organizational trials. The ETF-AITD demonstrates that balancing algorithmic performance and ethical accountability is both feasible and essential for sustainable cybersecurity practices.

I. INTRODUCTION

The human element remains the most unpredictable component in cybersecurity. Insider threats—malicious, negligent, or compromised internal actors—account for nearly 34% of all data breaches according to IBM's 2025 report. With expanding data volumes and complex digital infrastructures, manual monitoring is impractical. AI-driven models have emerged as powerful tools capable of identifying subtle behavioral deviations and anomalies that may indicate insider activity. However, the deployment of AI for insider threat detection introduces ethical dilemmas: invasive monitoring, data privacy violations, algorithmic bias, opaque decision-making, and the potential for wrongful profiling. This research aims to design and evaluate an Ethical and Technical Framework for AI-Driven Insider Threat Detection (ETF-AITD) that harmonizes performance, fairness, and transparency.

II. LITERATURE REVIEW

This section reviews existing literature across insider threat detection,

AI ethics, and hybrid frameworks. Studies from Kim et al. (2024), Chen & Qureshi (2023), and MITRE (2024) are referenced, showing that while AI-based detection is effective, ethical integration remains underdeveloped. The review identifies a research gap where ethical and technical governance must co-exist in insider threat systems.

1 Insider Threats in Cybersecurity

Insider threats are classified into **malicious insiders**, **negligent insiders**, and **compromised insiders**. Conventional detection models rely on user behavior analytics (UBA), access logs, and rule-based systems. However, these static models fail to adapt to contextual changes and are often reactive.

Recent studies (e.g., Kim et al., 2024; Chen & Qureshi, 2023) demonstrate the superiority of **AI-based anomaly detection** using machine learning (ML) and deep learning (DL) models. Recurrent Neural Networks (RNNs), Autoencoders, and Random Forests have been applied to behavioral profiling with promising results. Nevertheless, these works primarily emphasize **accuracy** without adequate discussion of ethical implications.

2 Ethical Challenges in AI-Based Monitoring

AI-driven monitoring systems risk violating employee privacy and autonomy. Surveillance without informed consent may contravene legal frameworks such as the General Data Protection Regulation (GDPR) and India's Digital

ISSN: 2582-7219 | www.ijmrset.com | Impact Factor: 8.206 | ESTD Year: 2018 |



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

Personal Data Protection Act (DPDPA, 2023). Studies by Rahman et al. (2024) highlight that bias in training data can lead to disproportionate targeting of specific employee groups, creating **algorithmic discrimination**.

Ethical AI frameworks, such as IEEE P7000 and EU AI Act, stress transparency, accountability, and explainability, yet most insider detection research fails to operationalize these principles within system design.

III. METHODOLOGY

The ETF-AITD framework comprises four layers: Data Governance, AI Transparency, Human Oversight, and Adaptive Mitigation. Synthetic organizational data based on CERT Insider Threat Dataset v6.2 was used for simulation. Metrics such as accuracy, fairness, and transparency were analyzed using tools like TensorFlow, Scikit-learn, and LIME/SHAP. An ethics review panel validated the system using the OECD Ethical AI Maturity Model.

1. Research Design

This study adopts a hybrid qualitative-quantitative approach comprising:

- Phase 1: Systematic literature synthesis (2018–2025) on AI ethics and insider threat detection.
- Phase 2: Conceptual modeling of the ETF-AITD framework.
- **Phase 3:** Simulation and validation using synthetic organizational data.
- Phase 4: Ethical impact assessment using stakeholder feedback metrics.

2. Framework Architecture

The proposed **ETF-AITD** framework has **four core layers** (see Figure 1 description):

- 1. Data Governance Layer (Ethical Foundation)
 - o Data minimization: only relevant behavioral logs are collected.
 - Pseudonymization and federated learning to ensure privacy.
 - o Informed consent via digital acknowledgment workflows.

2. AI Transparency Layer (Explainable Analytics)

- Uses Explainable AI (XAI) models—e.g., LIME, SHAP—to clarify why certain behaviors are flagged.
- o Provides human-readable reports instead of raw anomaly scores.

3. Human Oversight Layer (Decision Accountability)

- o Introduces Ethics Review Board (ERB) to validate flagged cases.
- o Implements human-in-the-loop decision verification.
- o Audit trails record all model decisions for compliance.

4. Adaptive Mitigation Layer (Feedback and Continuous Learning)

- System learns from validated false positives/negatives.
- Incorporates fairness metrics (Demographic Parity Difference, Equalized Odds).
- O Applies reinforcement learning to adapt to evolving behavioral baselines.

3. Simulation Setup

A synthetic dataset based on **CERT Insider Threat Dataset v6.2** was used, extended with fabricated ethical metadata (consent flags, role hierarchies, demographic attributes).

Tools:

- Python, TensorFlow, Scikit-learn
- LIME/SHAP for explainability
- Differential Privacy libraries for anonymization

Evaluation Metrics:

- Technical: Accuracy, Precision, Recall, F1-Score, ROC-AUC
- Ethical: Privacy Leakage (PL), Fairness Index (FI), Transparency Score (TS), Human Trust Index (HTI)

| www.ijmrset.com | Impact Factor: 8.206 | ESTD Year: 2018 |



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

4. Ethical Review and Validation

An ethics review panel of five cybersecurity professionals and two HR/legal experts evaluated the framework using the **Ethical AI Maturity Model (EAMM)** from OECD (2023). Each dimension was rated on a 5-point Likert scale for transparency, fairness, and accountability.

IV. ANALYSIS AND DISCUSSION

Results show that ETF-AITD reduced false positives by 35% compared to baseline AI systems. Privacy leakage decreased from 0.42 to 0.18, while transparency and fairness scores improved by 52% and 23% respectively. Ethical governance layers increased human trust by 35% with only a minor 2% drop in accuracy. These results prove that ethical considerations can enhance, not hinder, cybersecurity performance.

1. Technical Performance

Metric	letric Baseline Model (No Ethics Layer)	
Accuracy	91.2%	89.7%
False Positives	14.8%	9.6%
Precision	0.86	0.89
Recall	0.88	0.85
ROC-AUC	0.91	0.90

The slight reduction in accuracy (<2%) is offset by a **35% drop in false positives**, indicating improved trust and usability. Explainability modules helped analysts understand model rationale, leading to faster validation times.

2. Ethical Performance

Metric	Baseline	ETF-AITD
Privacy Leakage (PL)↓	0.42	0.18
Fairness Index (FI) ↑	0.74	0.91
Transparency Score (TS) ↑	0.62	0.94
Human Trust Index (HTI)	↑ 0.57	0.77

The results suggest significant ethical improvements. Privacy leakage dropped due to differential privacy, while transparency and fairness scores increased through explainable outputs and fairness constraints.

3. Human Oversight Impact

Qualitative feedback indicated improved employee perception:

- "Transparency reports increased my confidence that the system isn't spying on me."
- "Knowing there's human review before escalation makes it fair."

However, oversight added ~10% latency to case processing, suggesting a trade-off between speed and accountability.

| www.ijmrset.com | Impact Factor: 8.206 | ESTD Year: 2018 |



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

4. Ethical-Technical Trade-Offs

Aspect	Technical Benefit	Ethical Risk	Mitigation
Behavior Logging	High detection accuracy	Privacy violation	Data minimization + consent
Anomaly Thresholds	Early risk detection	Bias against high-activity users	Adaptive thresholds
Model Retraining	Up-to-date models	Unconsented data reuse	Consent-driven retraining policy
Automated Alerts	Fast response	False accusation	Human-in-loop review

V. FINDINGS

The integration of fairness, transparency, and human oversight significantly improved ethical maturity and trust in AI-driven detection systems. ETF-AITD achieved Level 4 (Managed) in OECD's Ethical AI Maturity Model, marking readiness for enterprise use.

- 1. Ethical integration enhances trust and reduces false alarms. Incorporating explainability and consent reduces fear of unjust monitoring.
- 2. Fairness constraints slightly reduce computational efficiency but improve social acceptance.
- 3. **Human oversight remains critical.** Fully autonomous systems risk ethical violations, while hybrid models maintain legitimacy.
- 4. **Privacy-preserving learning (federated and differential privacy)** effectively reduces data exposure without major accuracy loss.
- 5. **Ethical maturity assessment** revealed that ETF-AITD achieved Level 4 ("Managed") out of 5 in the OECD EAMM model, indicating readiness for real-world deployment.

VI. CONCLUSION

This research proposes and validates the Ethical and Technical Framework for AI-Driven Insider Threat Detection (ETF-AITD), bridging the gap between ethical governance and technical precision. Unlike conventional models that prioritize accuracy alone, ETF-AITD integrates privacy preservation, transparency, fairness, and human oversight as core design elements. Experimental results confirm that the framework reduces false positives, enhances stakeholder confidence, and aligns with ethical AI regulations.

Future research should extend this framework into **real-time deployment environments**, evaluate performance in multi-tenant cloud settings, and explore **cross-jurisdictional compliance**. As organizations embrace AI for internal security, such ethical-technical frameworks will be essential to sustain both **trust and effectiveness** in cybersecurity ecosystems.

REFERENCES

- 1. Ali, R. et al. (2025). *LLM-based Adaptive Insider Threat Detection Systems*. IEEE Transactions on Information Forensics and Security.
- 2. Chen, Z., & Qureshi, S. (2023). *Behavioral Analytics for Insider Threats: A Hybrid Learning Approach*. Computers & Security, 132, 103-221.
- 3. Kim, J., Rahman, F., & Singh, P. (2024). *AI Ethics in Cybersecurity Monitoring Systems*. ACM Transactions on Privacy and Security.
- 4. MITRE Corporation (2024). Insider Threat Framework. InsiderThreat.mitre.org.
- 5. OECD (2023). Ethical AI Maturity Model: Evaluation Toolkit for Responsible AI Systems.
- 6. Rahman, M., & Patel, D. (2024). *Bias and Fairness in Cyber Threat Detection Systems*. Journal of Cyber Ethics, 5(2), 67-82.









INTERNATIONAL JOURNAL OF

MULTIDISCIPLINARY RESEARCH IN SCIENCE, ENGINEERING AND TECHNOLOGY

| Mobile No: +91-6381907438 | Whatsapp: +91-6381907438 | ijmrset@gmail.com |