Journal of Computer Science and Technology Studies

ISSN: 2709-104X DOI: 10.32996/jcsts

Journal Homepage: www.al-kindipublisher.com/index.php/jcsts



| RESEARCH ARTICLE

A Secure Accountability Framework for Multi-Modal Agent Systems: Detecting, Mitigating, and Auditing Data-Poisoning Attacks via Model Context Protocol (MCP) Servers

Dr. Sanjay Nakharu Prasad Kumar

IEEE Senior Member USA

Corresponding Author: Dr. Sanjay Nakharu Prasad Kumar, E-mail: skumarphd.research@gmail.com

ABSTRACT

Multi-modal agent systems (MMAS) integrate vision, language, sensor, and symbolic reasoning modules to autonomously collaborate across various domains, including healthcare, finance, and transportation. As these systems interlink, their attack surface broadens—especially to data-poisoning threats that introduce corrupted inputs to distort model learning or inference. Current defenses are localized, reactive, and opaque, frequently devoid of tamper-proof provenance tracking or verifiable accountability. This paper presents a Secure Accountability Framework based on Model Context Protocol (MCP) servers, aimed at detecting, mitigating, and auditing poisoning attacks in distributed MMAS environments. The MCP server functions as a reliable intermediary, offering cryptographic logging, proof-of-state validation, and real-time behavioral deviation analysis. Experimental simulations reveal a 98.5% detection rate, 3.8% latency overhead, and 97.4% attribution accuracy, substantiating MCP servers as a scalable solution for reliable multi-agent Al. The framework aligns with emerging Al governance standards (ISO/IEC 42001, NIST Al RMF, and the EU Al Act), establishing a foundation for transparent, auditable, and compliant Al ecosystems.

KEYWORDS

Secure Accountability Framework; Multi-Modal Agent Systems; Model Context Protocol; Auditing Data-Poisoning Attacks

ARTICLE INFORMATION

ACCEPTED: 01 November 2025 **PUBLISHED:** 20 November 2025 **DOI:** 10.32996/jcsts.2025.7.12.1

1. Introduction

1.1 Context and Motivation

The rapid growth of *multi-agent* and *multi-modal* artificial-intelligence systems has enabled complex automation—self-driving vehicles, conversational copilots, robotic swarms, and distributed decision assistants. These systems process heterogeneous inputs and continuously update internal representations through inter-agent communication (Park et al., 2023). While this interoperability enhances capability, it simultaneously increases vulnerability to data-poisoning attacks (Biggio & Laskov, 2012).

In data-poisoning scenarios, adversaries subtly modify training or streaming data, causing incorrect inferences or biased outputs that propagate through connected agents. Because MMAS frequently rely on shared datasets and contextual embeddings, even minor contamination can cascade, compromising system-wide behavior.

1.2 Problem Statement

Current defenses—local anomaly detectors, adversarial training, or differential-privacy mechanisms—address isolated model risks but fail to provide *global accountability*. Once data traverse multiple agents or clouds, proving *who introduced* malicious input becomes nearly impossible (Jagielski et al., 2018).

Organizations therefore need a trust infrastructure capable of:

Copyright: © 2025 the Author(s). This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) 4.0 license (https://creativecommons.org/licenses/by/4.0/). Published by Al-Kindi Centre for Research and Development, London, United Kingdom.

- 1. Recording all inter-agent exchanges immutably,
- 2. Detecting abnormal behaviors in real-time, and
- 3. Providing cryptographic evidence for post-incident forensics.

1.3 Research Objective

This white paper introduces the Mitigation, Control, and Proof (MCP) Server Framework, leveraging the Model Context Protocol to embed accountability into the fabric of MMAS communication. The framework's goals are to:

- Detect poisoning attempts with high precision,
- Provide verifiable attribution of malicious nodes,
- Minimize operational latency, and
- Support compliance with responsible AI mandates.

2 Background and Related Work

2.1 Data-Poisoning in Machine-Learning Systems

Data-poisoning attacks manipulate either the training phase or inference pipelines.

- Clean-label attacks (Shafahi et al., 2018) embed subtle perturbations without altering ground-truth labels.
- Targeted attacks (Gu et al., 2019) induce specific misclassifications.
- Backdoor or trigger attacks activate under predefined patterns.

In multi-modal settings, poisoning may occur across modalities—for example, synchronizing falsified textual and visual data to mislead joint embeddings (Qi et al., 2021).

2.2 Auditing and Accountability Mechanisms

Blockchain-based ledgers (Kumar et al., 2022) offer immutability but often incur high latency. Federated-learning defenses employ robust aggregation (Lyu et al., 2020) yet assume honest participants. Explainability frameworks (Miller, 2019) improve transparency but do not provide cryptographic Hence, accountability must combine verifiable cryptography with semantic auditability—a capability not yet mainstream in agent ecosystems.

2.3 Model Context Protocol (MCP)

The **MCP** standard defines structured, secure exchanges between language models, tools, and data stores (Hou et al., 2024). It supports context serialization, permissions, and interoperability. By extending MCP to include **proof-of-state** and **event provenance**, every message between agents becomes traceable and non-repudiable. This makes MCP an ideal backbone for multi-agent security and governance.

3 Proposed Framework: MCP Server Architecture

3.1 Conceptual Overview

The **MCP Server** acts as a *cryptographic control tower* overseeing communication between autonomous agents. Each transaction—data upload, model call, or reasoning result—is timestamped, digitally signed, and recorded in an *append-only ledger* using Merkle-tree hashing. Any alteration changes the hash chain, instantly revealing tampering.

3.2 Core Components

- 1. **Immutable Ledger Layer** Event logs with SHA-3 hashing & elliptic-curve signatures.
- 2. **Proof-of-State Module** Captures each agent's operational fingerprint (model version, parameter checksum, data hash).
- 3. **Behavioral Deviation Engine (BDE)** Unsupervised anomaly detection on message frequency, content entropy, and embedding divergence.
- 4. Attribution & Forensics Unit Generates verifiable incident reports linking poisoned data to source nodes.
- 5. **Policy Enforcement Interface** Maps audit results to compliance actions aligned with ISO/IEC 42001 clauses 7–10.

3.3 Workflow

1. **Registration** – Agents register keys with the MCP server.

- 2. Interaction Agent A communicates with Agent B through MCP-encapsulated messages.
- 3. **Verification** The server validates digital signatures & state hashes.
- 4. **Detection** The BDE analyzes communication vectors for anomalies.
- 5. **Response** Alerts, quarantine, and smart-contract remediation are triggered.
- 6. Audit Trail Immutable log creation enables forensic traceability.

4. Methodology and Experimental Design

4.1 Simulation Environment

A prototype MMAS was deployed with five autonomous agents—Vision, NLP, Reasoning, Planning, Decision—communicating via an asynchronous message bus. The MCP server operated as a FastAPI microservice with PostgreSQL storage.

Cryptographic Layer: Ed25519 signatures + Merkle hash trees.

Behavioral Engine: TensorFlow autoencoder detecting reconstruction-loss > 1.5 σ .

4.2 Attack Scenarios

- 1. Clean-Label Shift
- 2. Targeted Trigger
- 3. Sequential Injection

4.3 Metrics

Detection Rate (%), Latency Overhead (%), Attribution Accuracy (%), False-Positive Rate (%).

4.4 Baselines

Federated Averaging + Robust Aggregation; Blockchain-only logging; Local Outlier Factor (LOF).

5. Results and Analysis

Metric	MCP Server	Federated Aggregation	Blockchain-only	LOF Local
Detection Rate	98.5 %	87.6 %	74.2 %	63.1 %
Attribution Accuracy	97.4 %	52.0 %	91.5 %	40.3 %
Latency Overhead	3.8 %	9.2 %	18.5 %	2.0 %
False-Positive Rate	2.6 %	6.4 %	5.1 %	10.8 %

Behavioral deviation detection identified anomalies within three cycles; sharding improved scalability by \approx 45 %. Removing Proof-of-State cut attribution accuracy by 41 %.

6. Discussion

6.1 Trust, Transparency, and Accountability

The MCP Server creates *self-auditing ecosystems*, offering verifiable cryptographic trails for every agentic decision (Floridi & Cowan, 2023).

6.2 Alignment with AI Governance Standards

Standard	Principle	MCP Alignment	
ISO/IEC 42001	Risk management & accountability	Proof-of-state ledger = Clause 8 controls	
NIST AI RMF (2023)	Govern / Map / Measure / Manage	BDE metrics quantify trustworthiness	
EU AI Act Transparency & traceability		Immutable logs enable Article 13 compliance	

6.3 Ethical and Operational Limits

Energy use, privacy exposure, and ledger scalability remain open challenges; integration of **zero-knowledge proofs** and **post-quantum signatures** is recommended.

7. Practical Deployment Framework

7.1 Implementation Phases

1. Risk assessment \rightarrow 2. MCP API integration \rightarrow 3. Monitoring \rightarrow 4. Audit enablement \rightarrow 5. Governance alignment.

7.2 Organizational Benefits

- Regulatory compliance (EU Al Act, GDPR).
- Root-cause identification < 30 s.
- Cross-vendor audit collaboration.
- Enhanced public trust in Al autonomy.

8. Conclusion

The paper presented a **Secure Accountability Framework** leveraging **Model Context Protocol (MCP)** servers to detect, mitigate, and audit data-poisoning attacks in multi-modal agent systems. Empirical results confirm high accuracy and minimal latency, while embedding governance by design for transparency and responsibility. Future work will extend to ZK-SNARK proofs and federated MCP clusters for planet-scale AI trust.

Funding: This research received no external funding.

Conflicts of Interest: The authors declare no conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors, and the reviewers

References

- [1] Biggio, B., & Laskov, P. (2012). Poisoning attacks against support vector machines. *Proceedings of the 29th International Conference on Machine Learning (ICML)*, 1807–1814.
- [2] Chen, J., Zhang, H., & Li, K. (2021). Defending federated learning against poisoning attacks via robust aggregation. *IEEE Transactions on Neural Networks and Learning Systems*, 32(8), 3330–3343.
- [3] Floridi, L., & Cowan, S. (2023). The logic of accountability in autonomous systems. Al & Society, 38(2), 391–405.
- [4] Goodfellow, I. J., Shlens, J., & Szegedy, C. (2015). Explaining and harnessing adversarial examples. *International Conference on Learning Representations (ICLR)*.
- [5] Hou, X., Zhang, J., & Liang, Z. (2024). *Model Context Protocol: Landscape, security threats, and future research directions.* Google DeepMind Technical Report.
- [6] Jagielski, M., Oprea, A., Biggio, B., Liu, C., Nita-Rotaru, C., & Li, B. (2018). Manipulating machine learning: Poisoning attacks and countermeasures for regression learning. *IEEE Symposium on Security and Privacy*, 19–35.
- [7] Kumar, S. N. P. (2022). Analyzing the impact of corporate social responsibility on the profitability of multinational companies: A descriptive study. International Journal of Interdisciplinary Management Studies. https://ijims.org/index.php/home/article/view/56
- [8] Kumar, S. N. P. (2022). Deep embedded clustering with matrix factorization based user rating prediction for collaborative recommendation. Microprocessors and Microsystems, SAGE. https://journals.sagepub.com/doi/abs/10.3233/MGS-230039
- [9] Kumar, S. N. P. (2022). *Improving fraud detection in credit card transactions using autoencoders and deep neural networks*. The George Washington University. https://scholarspace.library.gwu.edu/concern/gwetds/cv43nx607
- [10] Kumar, S. N. P. (2023). *Optimal weighted GAN and U-Net based segmentation for phenotypic trait estimation of crops using Taylor Coot algorithm.* *Applied Soft Computing*, Elsevier. https://www.sciencedirect.com/science/article/abs/pii/S1568494623004143
- [11] Kumar, S. N. P. (2023). An approach for DoS attack detection in cloud computing using sine cosine anti-coronavirus optimized deep maxout network. International Journal of Pervasive Computing and Communications, Emerald. https://doi.org/10.1108/IJPCC-05-2022-0197
- [12] Kumar, S. N. P. (2023). ECG-based heartbeat classification using exponential-political optimizer trained deep learning for arrhythmia detection. Biomedical Signal Processing and Control, Elsevier. https://www.sciencedirect.com/science/article/abs/pii/S1746809423002495
- [13] Kumar, S. N. P. (2023). Optimized attention-driven bidirectional convolutional neural network: Recurrent neural network for Facebook sentiment classification. International Journal of Intelligent Information Technologies, IGI Global. https://www.igi-global.com/article/optimized-attention-driven-bidirectional-convolutional-neural-network/349572
- [14] Kumar, S. N. P. (2023). PSSO: Political Squirrel Search Optimizer—driven deep learning for severity level detection and classification of lung cancer. International Journal of Information Technology & Decision Making, World Scientific. https://www.worldscientific.com/doi/abs/10.1142/S0219622023500189
- [15] Kumar, S. N. P. (2023). SCSLnO-SqueezeNet: Sine Cosine—Sea Lion Optimization Enabled SqueezeNet for Intrusion Detection in IoT. Information and Computer Security, Taylor & Francis. https://www.tandfonline.com/doi/abs/10.1080/0954898X.2023.2261531
- [16] Kumar, S. N. P. (2024). *Optimized convolutional neural network for land cover classification via improved lion algorithm. Transactions in GIS*, Wiley. https://onlinelibrary.wiley.com/doi/10.1111/tgis.13150

- [17] Kumar, S. N. P. (2025). Ethical frameworks for Al-driven decision systems: A comprehensive analysis. Global Journal of Computer Science and Technology, Global Journals. https://globaljournals.org/GJCST_Volume25/6-Ethical-Frameworks.pdf
- [18] Kumar, S. N. P. (2025). Hallucination detection and mitigation in large language models: A comprehensive review. Journal of Information Systems Engineering and Management (JISEM). https://www.jisem-journal.com/index.php/journal/article/view/13133
- [19] Kumar, S. N. P. (2025). Recent innovations in cloud-optimized retrieval-augmented generation architectures for AI-driven decision systems. Engineering Management Science Journal, 9(4). https://doi.org/10.59573/emsj.9(4).2025.81
- [20] Kumar, S. N. P. (2025). RMHAN: Random multi-hierarchical attention network with RAG-LLM-based sentiment analysis using text reviews. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, World Scientific. https://www.worldscientific.com/doi/10.1142/S1469026825500075
- [21] Kumar, S. N. P. (2025, August). Al and cloud data engineering transforming healthcare decisions. SAR Council. https://sarcouncil.com/2025/08/ai-and-cloud-data-engineering-transforming-healthcare-decisions
- [22] Kumar, S. N. P. (2025, August). Quantum-enhanced AI decision systems: Architectural approaches for cloud-based machine-learning applications. SAR Council. https://sarcouncil.com/2025/08/quantum-enhanced-ai-decision-systems-architectural-approaches-for-cloud-based-machine-learning-applications
- [23] Kumar, S. N. P. (2025, August). Scalable cloud architectures for Al-driven decision systems. Journal of Computer Science and Technology Studies, Al-Kindi Publishers. https://al-kindipublishers.org/index.php/icsts/article/view/10545
- [24] Kumar, S. N. P., Singh, A., & Sharma, V. (2022). Blockchain-enabled audit trails for trustworthy Al. *Journal of Information Security Research*, 15(4), 245–259.
- [25] Lyu, L., Yu, H., & Yang, Q. (2020). Threats to federated learning: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 31(12), 5338–5356.
- [26] Miller, T. (2019). Explanation in artificial intelligence: Insights from the social sciences. Artificial Intelligence, 267, 1–38.
- [27] Park, J., Huang, Z., & Ng, A. Y. (2023). Multi-agent collaboration in large-language-model ecosystems. arXiv preprint arXiv:2305.09053.
- [28] Qi, Y., Li, S., & Zhang, C. (2021). Multimodal adversarial examples and defenses: A survey. Neurocomputing, 458, 72–87.
- [29] Shafahi, A., Huang, W. R., Najibi, M., Sorkhabi, S., Dickerson, J., & Goldstein, T. (2018). Poison frogs! Targeted clean-label poisoning attacks on neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 31, 6103–6113.