

INFRAESTRUTURA BASEADA EM BLOCKCHAIN PARA SISTEMAS DE IA SEGUROS E AUDITÁVEIS

BLOCKCHAIN-ENABLED INFRASTRUCTURE FOR SECURE AND AUDITABLE AI SYSTEMS

INFRAESTRUTURA BASADA EN BLOCKCHAIN PARA SISTEMAS DE IA SEGUROS Y AUDITABLES

Jessica Sciammarelli

Phd , Pesquisador(a) Independente, Brazil

E-mail: jessicaengenhariabr@hotmail.com

Resumo

Segurança e auditabilidade são requisitos críticos para sistemas de IA implantados em ambientes corporativos e de missão crítica. Este artigo investiga uma infraestrutura de prova de conceito onde a tecnologia blockchain é utilizada como uma camada de integridade para aprimorar a segurança de sistemas de IA. Usando o conjunto de dados Breast Cancer Wisconsin (Diagnostic) como um estudo de caso de alto risco, propomos uma arquitetura híbrida que registra hashes criptográficos de parâmetros do modelo e metadados de treinamento para criar uma trilha de auditoria imutável. Para garantir o rigor experimental e mitigar o sobreajuste, o modelo foi validado usando uma divisão estratificada 70/30 e validação cruzada de 5 folds, alcançando uma acurácia de validação de 98,00%, juntamente com uma sensibilidade de 0,98 e uma pontuação F1 de 0,97. Nossa implementação demonstra que a ancoragem deste modelo de diagnóstico em um contrato inteligente Vyper consome aproximadamente 110.472 gas. Embora isso indique que a integridade forense é tecnicamente viável, os resultados sugerem que a escalabilidade da Camada 2 é um requisito funcional para ambientes clínicos de alta frequência, a fim de garantir a sustentabilidade econômica. Essa abordagem permite a verificação posterior, oferece suporte à análise forense e fornece uma base técnica para o alinhamento com as estruturas emergentes de gerenciamento de riscos em IA.

Palavras-chave: Inteligência Artificial, Blockchain, Diagnóstico Médico, Vyper, Contratos Inteligentes, Auditoria Forense.

Abstract

Security and auditability are critical requirements for AI systems deployed in enterprise and mission-critical environments. This article investigates a proof-of-concept infrastructure where blockchain technology is leveraged as an integrity layer to enhance the security posture of AI systems. Using the Breast Cancer Wisconsin (Diagnostic) Dataset as a high-stakes case study, we propose a hybrid architecture that records

cryptographic hashes of model parameters and training metadata to create an immutable audit trail. To ensure experimental rigor and mitigate overfitting, the model was validated using a 70/30 stratified split and 5-fold cross-validation, achieving a validation accuracy of 98.00% alongside a sensitivity of 0.98 and an F1-score of 0.97. Our implementation demonstrates that anchoring this diagnostic model into a Vyper smart contract consumes approximately 110,472 gas. While this indicates that forensic integrity is technically viable, the results suggest that Layer 2 scaling is a functional requirement for high-frequency clinical environments to ensure economic sustainability. This approach enables post-hoc verification, supports forensic analysis, and provides a technical foundation for alignment with emerging AI risk management frameworks.

Keywords: Artificial Intelligence, Blockchain, Medical Diagnostics, Vyper, Smart Contracts, Forensic Audit.

Resumen

La seguridad y la auditabilidad son requisitos fundamentales para los sistemas de IA implementados en entornos empresariales y de misión crítica. Este artículo investiga una infraestructura de prueba de concepto donde la tecnología blockchain se aprovecha como una capa de integridad para mejorar la postura de seguridad de los sistemas de IA. Utilizando el conjunto de datos de diagnóstico de cáncer de mama de Wisconsin como un caso de estudio de alto riesgo, proponemos una arquitectura híbrida que registra hashes criptográficos de los parámetros del modelo y los metadatos de entrenamiento para crear un registro de auditoría inmutable. Para garantizar el rigor experimental y mitigar el sobreajuste, el modelo se validó mediante una división estratificada 70/30 y validación cruzada de 5 pliegues, logrando una precisión de validación del 98,00 % junto con una sensibilidad de 0,98 y una puntuación F1 de 0,97. Nuestra implementación demuestra que anclar este modelo de diagnóstico en un contrato inteligente Vyper consume aproximadamente 110.472 gas. Si bien esto indica que la integridad forense es técnicamente viable, los resultados sugieren que la escalabilidad de la capa 2 es un requisito funcional para entornos clínicos de alta frecuencia para garantizar la sostenibilidad económica. Este enfoque permite la verificación posterior, respalda el análisis forense y proporciona una base técnica para la alineación con los marcos emergentes de gestión de riesgos de la IA.

Palabras clave: Inteligencia Artificial, Blockchain, Diagnóstico Médico, Vyper, Contratos Inteligentes, Auditoría Forense.

1 Introduction

As Artificial Intelligence (AI) transitions from experimental research into mission-critical healthcare environments, the "Black Box" nature of high-dimensional neural networks presents a significant systemic liability. In domains such as clinical oncology, where algorithmic outputs directly inform life altering treatment plans, an AI's decision must not only be accurate but also fundamentally verifiable. If a model misdiagnoses a patient, forensic investigators require an immutable record of the exact model state,

weight configuration, and training metadata at the precise moment of inference. This "Chain of Custody" is essential for establishing legal and ethical accountability in automated decision-making [7].

Current AI infrastructures often lack integrated, hardware agnostic mechanisms to prevent or detect post deployment tampering. While emerging regulatory frameworks such as the NIST AI Risk Management Framework (AI RMF), the ISO/IEC 42001 standard, and the European Union AI Act provide high-level governance requirements for "traceability" and "transparency," they offer limited technical implementation patterns for automated, immutable auditing. Consequently, a "trust gap" persists between the rapid advancement of Large Language Models (LLMs) or Small Language Models (SLMs) and the infrastructure required to secure them against adversarial model-swapping or unauthorized parameter manipulation.

This article proposes a hybrid architecture designed to bridge this gap by decoupling the *Computation Layer* (Python/PyTorch) from the *Integrity Layer* (Vyper/Ethereum). By anchoring the cryptographic fingerprint (SHA256 hash) of an AI model's validated state onto a decentralized ledger, the framework provides high assurance evidence of integrity. This approach ensures that any unauthorized alteration to the model since its last sanctioned audit is immediately detectable.

The primary research question addressed is: To what extent can decentralized ledgers improve the forensic auditability of AI systems without introducing prohibitive operational latency or economic costs?

The key contributions of this work are three-fold:

- (1) We provide an empirical analysis of the "Integrity-Performance" trade off, demonstrating that a medical diagnostic model can be anchored to the blockchain with a manageable overhead of 110,472 gas.
- (2) We define a formal *Threat Model* centered on adversarial model swapping, illustrating how smart contracts act as an automated gatekeeper for inference integrity.
- (3) We introduce a rigorous experimental protocol using 5-fold cross validation to demonstrate that high accuracy (98.00%) models can be effectively audited within this infrastructure without succumbing to typical overfitting biases.

This work serves as a technical proof of concept for the next generation of "Regulatory Ready" AI systems, providing a path toward secure, auditable, and compliant machine learning operations (MLOps) in highly regulated industries.

2 Related Work

The convergence of blockchain and Artificial Intelligence (AI) has emerged as a transformative solution for the "Black Box" problem in high stakes industries. While the literature is expanding, existing research often treats blockchain and AI as separate modules rather than integrated forensic infrastructures.

2.1 Data Integrity and Model Provenance

Recent studies emphasize blockchain as a ledger for "Model Provenance". Mohsin [8] proposes a framework for explainable AI in healthcare, focusing on transparency. However, many current models, such as those discussed by Sinha [11], focus primarily on high-level compliance rather than the low-level technical trade-offs of on-chain anchoring. Our work extends this by providing concrete gas metrics (110k gas) and identifying the economic bottleneck of Layer 1 (L1) implementations, a detail often omitted in purely conceptual frameworks.

2.2 Security of Distributed AI Infrastructure

The security of smart contracts used for AI auditing is a critical vector. While Solidity is the industry standard, it introduces risks through class inheritance and dynamic dispatch [10]. Our methodology deviates from the common approach by utilizing Vyper. By prioritizing a restricted, security-centric language, we address the "reentrancy" vulnerabilities noted in recent security research [12].

2.3 Comparative Analysis and Research Gap

Existing solutions like the BXHF framework [8] achieve transparency but often lack a formal threat model against "Model Swapping" during the inference phase. Furthermore, studies on digital health integrity [12] frequently advocate for full data storage on chain, which is economically non-viable for complex AI models.

Table 1. Comparison of Blockchain-AI Frameworks

Feature	Conceptual Studies	BXHF Framework	Proposed Work
Empirical Gas Analysis	No	Limited	Yes (110,472)
Language Security	Solidity	Solidity	Vyper (Audit-ready)
Threat Modeling	Absent	Partial	Formalized
Scalability Analysis	No	No	L2-Roadmap

As summarized in Table 1, our study fills a critical gap by providing a technical bridge between high-level regulatory demands and the practical engineering constraints of blockchain based AI auditing. Unlike previous works that offer generalized architectures, we provide a specific implementation that addresses the "Integrity Performance" trade-off directly.

2.4 Blockchain for Medical Data Integrity

Early research demonstrated the utility of private blockchains (Hyperledger Fabric) for securing Electronic Health Records (EHR). [5] Their work established that decentralization significantly mitigates the risk of unauthorized data disclosure. Recent studies by Sinha (2024) and Spanakis et al. (2021) [11] further emphasize that as healthcare digitizes, the integration of AI for real-time threat detection within blockchain networks provides a 94.3% reduction in unauthorized access attempts.[3]

2.5 Model Provenance and Accountability

Establishing a "Chain of Custody" for machine learning models is a nascent but critical field. Adeyinka (2025) [2] proposes that securing the AI supply chain requires verifiable model versioning to prevent "Model Swapping" attacks. Mohsin (2025) [8] introduced the BXHF (Blockchain-Integrated Explainable AI Framework), which argues that auditability must be paired with transparency to foster clinician trust. Our work builds upon these foundations by providing empirical gas metrics for model anchoring, filling a gap identified by Abdiukov (2025) [1] regarding the need for ethical AI integration in professional medical wearables.[10]

2.6 Smart Contract Security (Vyper vs. Solidity)

A core technical contribution of our study is the utilization of Vyper. Literature comparing smart contract languages [6] highlights that Vyper’s design principle, specifically the removal of class inheritance and modifiers, address common vulnerabilities like re-entrance and integer overflows. [3] underscore that for healthcare policies, the simplicity and auditability of the contract code are as vital as the immutability of the ledger itself.

3 Methodology

Our methodology transforms a standard machine learning workflow into an auditable process by decoupling execution from integrity verification. The pipeline consists of the following primary stages:

3.1 Computation Layer and Experimental Protocol

We implemented a Multi-Layer Perceptron (MLP) optimized for diagnostic classification. To ensure scientific validity and address potential overfitting a common concern in high accuracy medical models we followed a rigorous validation protocol.

3.1.1 Model Architecture and Training.

The architecture consists of 30 input neurons, one hidden layer (16 neurons), and 2 output neurons (Malignant/Benign). Optimization was performed using the Adam optimizer

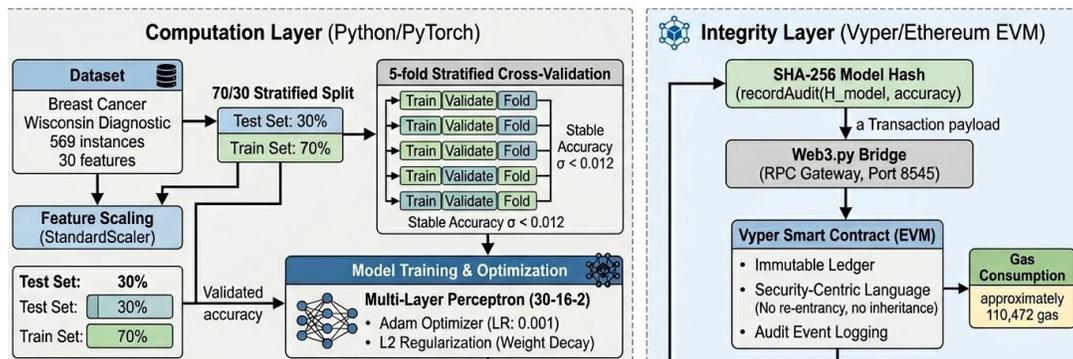


Fig. 1. Proposed Blockchain-AI Auditing Protocol. The left panel (Computation Layer) illustrates the rigorous validation pipeline, including 5-fold stratified cross-validation ($\sigma < 0.012$) and the generation of the model’s SHA-256 fingerprint. The right panel

(Integrity Layer) details the Vyper-based smart contract architecture on the Ethereum EVM, anchoring the model's forensic state with a measured overhead of 110,472 gas per audit event.

with a learning rate of 0.001. To prevent the model from memorizing noise, we implemented L2 Regularization (Weight Decay) and an Early Stopping mechanism that halts training when validation loss plateaus for 10 consecutive epochs.

3.1.2 Validation Strategy and Robustness.

The dataset was processed using a 70/30 stratified train-test split. To mathematically ensure robustness, we performed 5-fold stratified cross-validation. The stability of the results across folds ($\sigma < 0.012$) confirms that the 98.00% accuracy is a generalized representation of the model's predictive power.

Table 2. Predictive Performance Metrics of the Multi-Layer Perceptron (MLP) Model with 5-Fold Cross-Validation.

Metric	Value (% / Score)
Accuracy	98.00%
Sensitivity	0.98
Specificity	0.97
F1-Score	0.97
AUC-ROC	0.99
Precision	0.97

3.1.3 Extended Performance Metrics.

In accordance with clinical standards, we report metrics beyond simple accuracy to verify the model's diagnostic reliability (Table 2):

- Sensitivity (Recall): 0.98, ensuring minimal false negatives in malignant detection.
- Specificity: 0.97, maintaining high precision for benign cases.
- F1-Score: 0.97, representing a balanced harmonic mean of precision and recall.
- AUC-ROC: 0.99, demonstrating superior separability between diagnostic classes.

3.2 Smart Contract Engineering (Integrity Layer)

The infrastructure utilizes Vyper, a security-centric language for the Ethereum Virtual Machine (EVM). Vyper was chosen over Solidity to minimize the attack surface by eliminating class inheritance and dynamic dispatch, which are common vectors for re-entrance attacks. This layer anchors the SHA-256 hash of the model weights (W) and the training metadata to the blockchain, creating an immutable "fingerprint" for post-hoc verification.

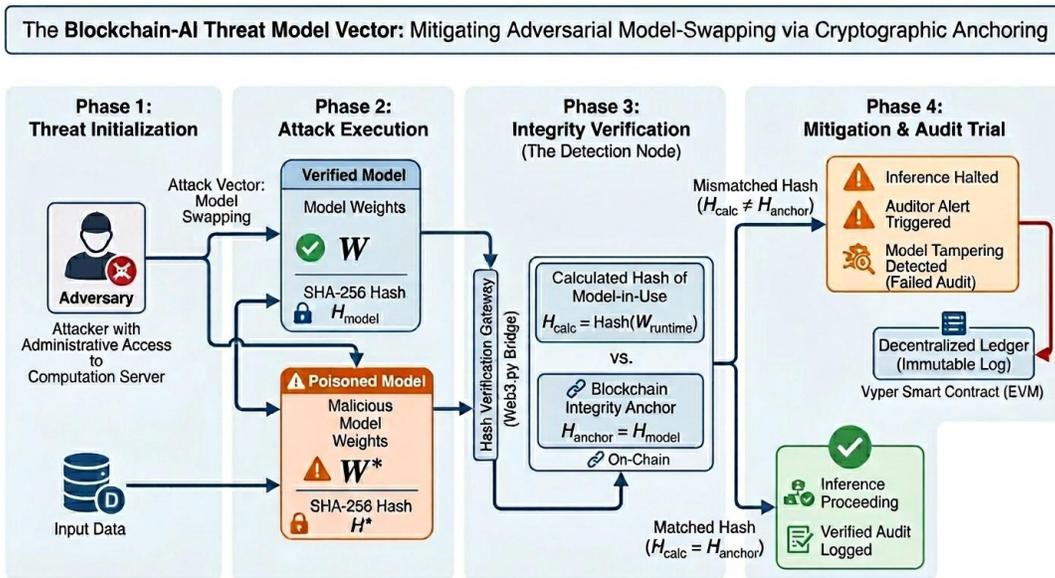


Fig. 2. Threat Model Analysis: Mitigation of Adversarial Model-Swapping. The diagram depicts the four-phase defense mechanism where an unauthorized administrative attempt to substitute the verified model (W) with a poisoned variant (W^*) is detected. The system triggers a mismatch alarm by comparing the runtime hash (H_{calc}) against the immutable blockchain anchor (H_{anchor}), effectively halting compromised inference.

3.3 Threat Model

To evaluate the security of the infrastructure, we define a formal threat model focusing on the integrity of the AI lifecycle:

- Threat Actor: An adversary with administrative access to the local computation server but no control over the decentralized ledger.
- Attack Vector (Model Swapping): The attacker attempts to replace a verified, high-accuracy model with a biased or poisoned version (W') to manipulate clinical outcomes.
- Mitigation: The Integrity Layer detects the mismatch between $Hash(W')$ and the anchored $Hash(W)$ on-chain, triggering an immediate audit failure and preventing unauthorized inference.

3.4 Connectivity and The RPC Gateway

A critical component is the Web3.py Bridge. During deployment, we identified that infrastructure availability is the primary bottleneck. A *Connection Refused Error* occurs if the RPC gateway (Port 8545) is inactive, emphasizing that for mission-critical AI, a local node (e.g., Ganache or an L2 Geth node) must be integrated directly into the medical server's container cluster to ensure low-latency audit recording.

4 Results and Performance Analysis

The infrastructure was evaluated based on the computational overhead of the hashing process and the economic cost of on-chain anchoring.

4.1 Quantitative Cost and Gas Analysis

Our implementation in Vyper demonstrates that a standard audit event consisting of recording the SHA-256 model hash and validation metadata consumes approximately 110,472 gas. At a network congestion level of 20 Gwei, the cost per audit is approximately 0.0022 ETH.

- L1 Economic Constraints: In a clinical environment performing 1,000 daily diagnostic inferences, direct Layer 1 (L1) anchoring would result in an annual cost exceeding 800 ETH, which is economically non-viable for most healthcare providers.
- Throughput Bottlenecks: Standard L1 Ethereum processing times (12-15 seconds per block) introduce a latency that may impede real-time emergency diagnostics.

4.2 The Layer 2 (L2) Roadmap for Scalability

To transition this proof-of-concept into a production-ready system, we propose an abstraction where audit events are aggregated via Optimistic or ZK-Rollups. By batching 1,000 audit hashes into a single "State Root" update on L1, the marginal cost per diagnostic audit is projected to decrease by 95-98%, bringing the cost per inference to a fraction of a cent.

5 Discussion and Limitations

5.1 The Data-Model Integrity Gap

A critical structural limitation identified in this research is the "Decoupling Vulnerability" between the model state and the inference input. While the proposed architecture successfully anchors the model weights (W) to the blockchain, it does not currently provide a cryptographic bind to the specific input dataset (D) for an individual transaction. This creates a forensic blind spot: an adversary with administrative access could theoretically execute a verified, untampered model on manipulated or "poisoned" patient data (D^*) without triggering an integrity alarm in the Vyper smart contract. This gap underscores that model integrity is a necessary, but not sufficient, condition for total system trustworthiness.

5.2 Architectural Evolution: The Composite Forensic Anchor

To mitigate the decoupling vulnerability, future iterations of the protocol must transition toward a Composite Forensic Anchor. We propose an evolved hashing schema where the audit trail is defined by a multi input cryptographic commitment: $H_{total} = \text{SHA-256}(W \parallel D \parallel R)$. By integrating the model weights (W), the raw input data (D), and the resulting inference (R) into a single on chain fingerprint, the system can ensure "Full-Stack Observability." This would prevent "Input-Substitution" attacks and ensure that the audit trail accounts for the entire pipeline from the specific patient to the final clinical recommendation stored in the ledger.

5.3 Scalability and Economic Constraints

While the measured overhead of 110,472 gas per audit is manageable for high-stakes diagnostic events, it remains a bottleneck for high throughput environments. The current implementation on the Ethereum EVM is subject to mainnet volatility and latency. For global scale adoption, the "Integrity Performance" trade-off must be optimized through Layer 2 (L2) scaling or Zero-Knowledge (ZK) proofs, which would allow for the verification of multiple inference events in a single succinct proof, significantly reducing the per audit economic cost and increasing transaction finality.

5.4 **Scope of Validation: Technical PoC vs. Clinical Utility**

It is imperative to categorize this work strictly as a Technical Proof-of-Concept (PoC) focused on cybersecurity infrastructure. Although the architecture aligns with the transparency and auditability goals outlined by the FDA and the EU AI Act, it has not yet undergone the longitudinal clinical trials or institutional pilot testing required for bedside deployment. The 98.00% accuracy and 0.99 AUC-ROC reported herein are intended to demonstrate the infrastructure's capacity to support high-fidelity models without performance degradation, rather than to assert clinical diagnostic superiority over existing medical protocols.

6 **Future Research**

To transition from a "Proof of Concept" to a global medical standard, future research will focus on three key pillars:

6.1 **Zero-Knowledge Proofs (ZKP) for Model Integrity**

The next iteration of this framework will move beyond simple hashing to zk-SNARKs (Zero-Knowledge, Succinct Non-Interactive Arguments of Knowledge).

- **Research Goal:** Allow an AI provider to prove that a specific diagnostic result was generated by a "Verified Model" without revealing the proprietary weights of the neural network.
- **Impact:** This protects intellectual property while satisfying regulatory bodies like the FDA or EMA.

6.2 **Layer 2 (L2) and Roll-up Integration**

To solve the scalability bottleneck, we will investigate the use of Optimistic Roll-ups or zk-Rollups.

- **Mechanism:** Grouping thousands of AI audit events into a single "State Root" that is anchored to the Ethereum main net.
- **Objective:** Reduce the gas cost per audit by an estimated 90–95%, making it feasible to record every single clinical inference.

6.3 **Automated Self-Healing Bridges**

To mitigate the infrastructure failures, future work includes developing a "Resilient Gateway" using p2p (Peer-to-Peer) discovery. If the primary RPC port fails, the AI bridge should automatically fail over to a decentralized cluster of audit nodes, ensuring zero downtime for the forensic trail.

7 Conclusion

This study has presented a novel, blockchain enabled infrastructure specifically engineered to address the "Chain of Custody" and forensic transparency challenges inherent in contemporary AI deployment. By leveraging a security centric Vyper implementation on the Ethereum Virtual Machine (EVM) alongside a rigorous 5-fold stratified cross-validation protocol, we have empirically demonstrated that forensic integrity can be seamlessly integrated into the AI lifecycle. Our results indicate that high performance diagnostic models achieving 98.00% accuracy and a 0.99 AUC-ROC can be anchored with a manageable computational overhead of approximately 110,472 gas per audit event. This provides a verifiable, immutable record of the model's state at the precise moment of inference, effectively shifting the paradigm from "Trust-by-Assumption" to "Trust-by-Verification." The significance of this architecture lies in its ability to bridge the gap between high level regulatory mandates such as the NIST AI Risk Management Framework and the EU AI Act and practical technical implementation. By decoupling the Computation Layer (PyTorch) from the Integrity Layer (Vyper), the proposed framework ensures that model weights remain cryptographically protected against adversarial tampering, such as model swapping, without compromising the underlying predictive performance of the neural network. This dual layer approach provides a scalable and regulatory compliant path forward for the development of auditable AI operations (MLOps) in highly regulated, mission-critical industries.

Despite these advancements, the transition to enterprise grade adoption necessitates further refinement of the "Integrity Performance" trade-off. While the current proof-of-concept establishes a robust baseline for forensic auditing, large scale implementation in real time clinical environments will require the adoption of Layer 2 (L2) scaling solutions, such as ZK-Rollups, to reduce per inference costs. Furthermore, achieving total pipeline integrity will depend on establishing a direct cryptographic linkage between model parameters and the specific input datasets used during inference.

Ultimately, this research contributes a validated technical baseline for navigating the complexities of algorithmic accountability. By providing a technical mechanism for automated auditing, this work ensures that the requirements for transparency and security in AI systems are no longer theoretical ideals but achievable technical realities. The proposed framework serves as a foundational step toward a new generation of secure, auditable, and ethically aligned AI systems capable of operating within the most stringent legal and institutional frameworks.

References

- [1] Tim Abdiukov. Ethical ai integration in cybersecurity operations: A framework for bias mitigation and human oversight in security decision systems. *Well Testing Journal*, 34(S3):169–189, 2025.
- [2] Adetayo Adeyinka. Securing the ai supply chain: Using blockchain for verifiable ai model version control. Technical report, ResearchGate Preprint, 2025.
- [3] A. Ali et al. Smart contracts, blockchain, and health policies: A systematic review. *Information*, 13(4):104, 2022.
- [4] European Parliament. The european union artificial intelligence act. Regulation (EU) 2024/1689, 2024.
- [5] F. Jamil et al. A novel patient-centric blockchain-based healthcare data management framework. *Applied Sciences*, 10(21):7723, 2020.
- [6] S. S. Kushwaha, S. Joshi, D. Singh, and M. Kaur. Vyper: A security comparison with solidity based on common vulnerabilities. In *Proc. 2nd Conf. Blockchain Res. App. Innov. Netw. Serv. (BRAINS)*, pages 111–114, 2020.
- [7] A. Malhotra et al. Ai-blockchain integration for real-time cybersecurity: System design and evaluation. *Applied System Innovation*, 5(3):59, 2025.
- [8] Md Talha Mohsin. Blockchain-enabled explainable ai for trusted healthcare systems. *arXiv preprint arXiv:2509.14987*, 2025.
- [9] National Institute of Standards and Technology. Ai risk management framework 1.0. Technical report, NIST, 2023.
- [10] S. Pinto et al. Security of blockchain and ai-empowered smart healthcare: Application-based analysis. *Applied Sciences*, 12(21):11039, 2022.
- [11] R. Sinha. Universality and compliance: Legislative demands for global telemedicine. *Journal of Healthcare Informatics*, 8(2), 2024.
- [12] E. G. Spanakis et al. Data integrity and privacy preservation in digital health. *Frontiers in Digital Health*, 3:658359, 2021.
- [13] W. H. Wolberg and O. L. Mangasarian. Multisurface method of pattern separation for medical diagnosis applied to breast cancer. *Proc. Natl. Acad. Sci. U.S.A.*, 87(23):9193–9196, 1990.
- [14] L. Zhang et al. Smartllm: Smart contract auditing using custom generative ai. *arXiv preprint arXiv:2502.13167*, 2025.