

Kennesaw State University

DigitalCommons@Kennesaw State University

Master of Science in Information Technology
Theses

Department of Information Technology

Fall 12-1-2023

Assessing Blockchain's Potential to Ensure Data Integrity and Security for AI and Machine Learning Applications

Aiasha Siddika
Kennesaw State University

Follow this and additional works at: https://digitalcommons.kennesaw.edu/msit_etd



Part of the [Computer and Systems Architecture Commons](#)

Recommended Citation

Siddika, Aiasha, "Assessing Blockchain's Potential to Ensure Data Integrity and Security for AI and Machine Learning Applications" (2023). *Master of Science in Information Technology Theses*. 16. https://digitalcommons.kennesaw.edu/msit_etd/16

This Thesis is brought to you for free and open access by the Department of Information Technology at DigitalCommons@Kennesaw State University. It has been accepted for inclusion in Master of Science in Information Technology Theses by an authorized administrator of DigitalCommons@Kennesaw State University. For more information, please contact digitalcommons@kennesaw.edu.

Assessing Blockchain's Potential to Ensure Data Integrity and Security for
AI and Machine Learning Applications



A Thesis Presented to
The Faculty of Information Technology Department

by
Aiasha Siddika

Committee Members
Dr. Liang Zhao (Chair)
Dr. Seyedamin Pouriye (Committee Member)
Dr. Xinyue Zhang (Committee Member)

In Partial Fulfillment of Requirements for the Degree
Master of Science in Information Technology

Kennesaw State University
Kennesaw, Georgia

December 2023

Acknowledgments

I would like to take this opportunity to express my sincere gratitude to everyone who has helped me along the way with my thesis. Throughout my academic journey, I have found inspiration, strength, and enlightenment in Allah's unfailing direction. With his blessings, my path has been lighted, enabling me to successfully complete this task.

I would like to express my gratitude to Dr. Liang Zhao, my thesis adviser, for his steadfast support, priceless advice, and endless patience. His knowledge and guidance have greatly influenced the direction of my study.

I owe my family a huge debt of gratitude for their unwavering encouragement, support, and faith in my potential. Their steadfast belief in me has served as a continual source of inspiration. In addition, I would like to thank my friends and colleagues for being a source of support and inspiration throughout my academic journey. The concepts put forward in this thesis have been greatly influenced by their conversations, advice, and support.

Furthermore, I would like to thank all of the study participants and contributors whose assistance and wisdom have been crucial in obtaining the data and information required for this project.

Finally, I want to express my gratitude to everyone who has helped along the way, whether directly or indirectly. Your encouragement, support, and contributions have been crucial to this thesis's successful completion.

Abstract

The increasing use of data-centric approaches in the fields of Machine Learning and Artificial Intelligence (ML/AI) has raised substantial issues over the security, integrity, and trustworthiness of data. In response to this challenge, Blockchain technology offered a promising and practical solution, as its inherent characteristics as a decentralized distributed ledger, coupled with cryptographic processes, offer an unprecedented level of data confidentiality and immutability. This study examines the mutually beneficial connection between Blockchain technology and ML/AI, using Blockchain's inherent capacity to protect against unauthorized alterations of data during the training phase of ML models. The method involves building valid blocks of data from the training dataset and then sending them to the mining process using smart contracts and the Proof of Work (PoW) consensus method. Using SHA256 to produce a cryptographic signature for each data block improves the aforementioned procedure. The public Ethereum blockchain serves as a secure repository for these signatures, whereas a cloud-based infrastructure houses the original data file. Particularly during the training phase of Machine Learning (ML) models, this cryptographic framework is critical in ensuring the data verification procedure. This research investigates the potential collaboration between Blockchain technology and ML/AI, bolstering data quality and trust to enhance data-driven decision-making fortifying the models' ability to provide precise and dependable results.

Table of Contents

Acknowledgements	2
Abstract	3
Table of Contents	4
List of Figures	5
1 Introduction	6
1.1 Motivation and Challenges	9
1.2 Research Aim	10
2 LITERATURE REVIEW	12
3 Methodology and Framework	17
4 Framework Development	21
5 Recommendations and Future Prospects	32
5.1 Insights and Roadmap for the Future	32
5.1.1 Integration of Proof of Federated Learning (PoFL)	32
5.1.2 Collaborative Secure Computing Integration	33
6 Conclusion	34
Bibliography	35

List of Figures

1.1	Advantages of blockchain and AI/ML integration	10
3.1	Proposed system architecture for creating data blocks to train AI/ML models.	18
3.2	Work flow diagram	20
4.1	Application Architecture	21
4.2	User Authentication	23
4.3	Category of Interest Subscription by Users	24
4.4	Subscribed Dataset Dashboard	25
4.5	Dataset Upload	26
4.6	Data download for voting and validation	27
4.7	Applying filter and search	27
4.8	Dataset stored in cloud storage	28
4.9	Dataset Verification and Voting Status Evaluation	29
4.10	Consensus Ledger ensuring transparency and traceability	30
4.11	Data Integrity Verification Process	30

Chapter 1

Introduction

The advent of artificial intelligence (AI) and machine learning (ML) has culminated in a paradigm shift across several sectors, sparking a significant transformation that aligns remarkably with the advancement of 'Industry 4.0.' The aforementioned technologies have not only extensively infiltrated but also fundamentally transformed the framework of several industries, resulting in a surge of inventive advancements that have significantly impacted sectors spanning a wide range. The significant shift is distinguished by the use of automation in various jobs, the advanced skills in recognizing patterns, and the facilitation of decision-making processes led by data¹. The revolution's influence is evident across a wide range of sectors, including conventional manufacturing, logistics, healthcare, and the complex realm of finance. In the field of manufacturing, AI and ML technologies have been used to improve production processes, resulting in increased efficiency and output. The incorporation of these technologies in the field of logistics and shipping has facilitated a more efficient and optimized approach to managing the supply chain. This advancement is primarily driven by the use of predictive analytics and automation. The healthcare industry, which has significant significance, has effectively used AI and ML techniques to facilitate illness diagnosis, evaluate medical imaging, and even forecast patient outcomes.

However, this technological trajectory is not devoid of challenges, as the escalating inte-

gration of AI and ML models into critical applications has raised these significant apprehensions regarding their susceptibility to attacks. There is a growing perception that malicious actors recognize the vulnerability inherent in the training phase. This issue stems from the recognition that the modification of training data has the potential to cause distortions, errors, or compromised conditions throughout the AI/ML models². The consequences of these breaches have significant implications for key areas including as healthcare, banking, and autonomous systems, through algorithmic bias, unequal access, cascading failures, and trade-offs between efficiency and resilience, necessitating further exploration and improved security measures³. The attacks involve manipulating or introducing malicious data into the training dataset on purpose, with the explicit intent of subverting the inherent learning process of the model. Ian et al.⁴ addressed the concept of adversarial attacks used by malicious actors to modify input data in order to deceive the learning process, resulting in potential weaknesses in the performance of the model. This type of attack exploits the model's sensitivity to anomalies, noise, and biased information inherent to the training process. Adversarial attacks encompass the careful insertion of calculated perturbations or modifications into training data, deceiving the model into generating inaccurate predictions⁵. Conversely, poisoning attacks involve infusing a subset of malicious samples into the training dataset to skew the model's decision boundaries, favoring specific categories or outcomes⁶.

As ML advances and pervades diverse sectors, safeguarding the training phase against malicious tampering emerges as a progressively pivotal facet in upholding the dependability and security of ML models. Mitigating these training phase attacks necessitates resilient strategies, such as meticulous preprocessing of data, implementation of anomaly detection methodologies, utilization of adversarial training to bolster model resistance against adversarial disturbances, and ensuring the unclassified and representativeness of training data⁷. Here comes Blockchain, the underlying technology that powers cryptocurrencies like Bitcoin⁸. At its core, Blockchain is a decentralized distributed ledger⁹ that, utilizing

cryptographic techniques, provides confidentiality and immutability. Integrating Blockchain technology in AI/ML ecosystems promises to improve these systems against attacks, leveraging decentralization, immutability, transparency, and cryptographic security to address evolving threats. Blockchain provides a secure framework for monitoring the provenance and modifications of training data, models, and predictions by establishing an immutable and transparent ledger of transactions and data alterations. This enables stakeholders to trace and verify any unauthorized modifications, preserving the AI/ML ecosystem's integrity¹⁰. In addition, the decentralized nature of Blockchain improves collaboration between parties by facilitating the exchange of threat intelligence and enabling proactive responses to emergent attack vectors, as this survey¹¹ explores the crucial merging of Blockchain technology with ML with the goal to provide a safe and decentralized exchange of data and models, while enhancing the efficiency of communication and networking systems. By taking advantage of Blockchain's tamper-resistant properties, we can develop innovative approaches to protect the data required to train ML models, improve the privacy of sensitive data, and establish distributed and trustworthy AI systems.

Blockchain is a decentralized system that utilizes a peer-to-peer (P2P) architecture to enhance the security of data, making it significantly resistant to unauthorized alteration or deletion by malicious individuals. The system is formed of several nodes which are responsible for storing and validating transactions inside blocks. Each block is designed to establish precise connectivity with the previous blocks, so creating a chain-like structure. Miners contribute new blocks to the blockchain network, and the validation and agreement of transactions occur via a decentralized consensus mechanism. This protocol employs techniques such as Proof-of-Work (PoW) and Proof-of-Stake (PoS). This measure guarantees the safeguarding of the blockchain against both internal and external hacking endeavors. Smart contracts are electronic agreement programs that are performed depending on predetermined parameters. These contracts play a significant role in reducing risks and enhancing cost efficiency in many corporate activities¹². As consensus mechanisms for blockchain net-

works, the debate over the suitability of Proof of Stake (PoS) and Proof of Work (PoW) in augmenting data security during the training phase of AI/ML models, PoW stands out as the more pertinent option¹³.

1.1 Motivation and Challenges

The convergence of blockchain technology with the rapidly evolving domains of AI/ML has sparked significant attention and anticipation. In light of the significant advancements achieved in the fields of machine learning (ML) and artificial intelligence (AI), which have brought about profound transformations in several sectors, there has been a notable increase in the use of data-driven approaches. However, this surge has also given rise to some difficulties pertaining to the preservation of data integrity, ensuring security, and maintaining dependability. As the complexity and magnitude of these difficulties increase, it becomes imperative to prioritize the establishment of data trustworthiness within the domain of machine learning and artificial intelligence applications.

The importance of data security and integrity in this particular circumstance cannot be overemphasized enough. The increasing prevalence of data-driven methodologies has led to a heightened need for ensuring protection against unauthorized alterations, which has become a matter of utmost importance. Blockchain technology is becoming recognized as a viable and effective means to mitigate the trust and security concerns that are inherent in the foundational data supporting AI/ML applications depicted in 1.1¹⁴. Through the use of the inherent tamper-resistant characteristics of blockchain technology, novel approaches may be developed to enhance the integrity of crucial data utilized in the training of machine learning models. The implementation of this approach not only serves to augment the protection of sensitive information, but also serves to construct AI systems that are decentralized and deemed trustworthy. Consequently, this development ushers in a new epoch characterized by enhanced dependability and security in technologies powered by data.

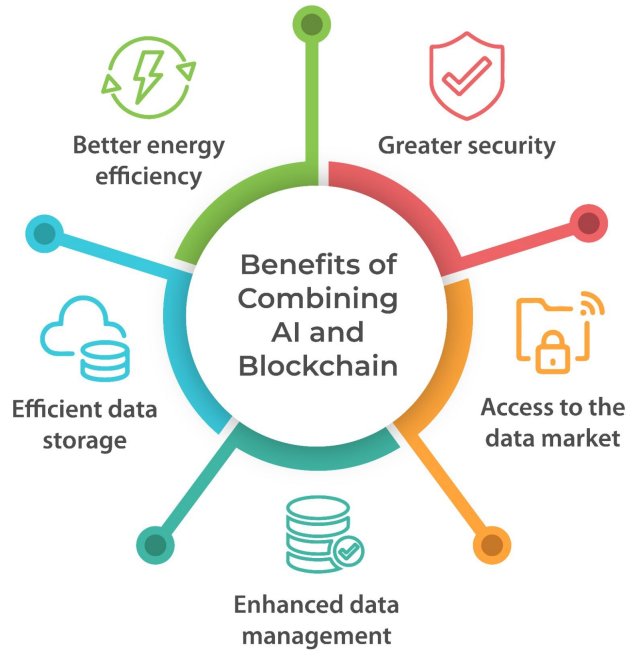


Figure 1.1: *Advantages of blockchain and AI/ML integration*

1.2 Research Aim

The core objective of this study is to examine and illustrate the possible coordination that might arise from the integration of Blockchain technology with ML/AI models. The key focus is on using the inherent attributes of Blockchain, like decentralization, immutability, and cryptographic security, to augment the security, integrity, and reliability of data in the training phase of machine learning models. The research objectives outlined in this study encompass:

- Providing a stable, open, and robust setting for the training of machine learning models, with the ultimate aim of enhancing decision-making processes based on data.
- Creating a novel framework that leverages PoW as a consensus mechanism, SHA256 for cryptographic signature generation, and smart contracts for data block construction.

This will be achieved by tackling significant obstacles to preserving data integrity and se-

curity during the training process.

Chapter 2

LITERATURE REVIEW

This chapter begins a thorough examination of existing scholarly works, exploring the complex convergence of blockchain and AI/ML. A multitude of research articles and studies provide the theoretical and empirical foundations for how Blockchain’s foundational principles have been harnessed to address the inherent vulnerabilities of AI/ML models. Through an in-depth review of the contributions made by pioneers and experts in the field, our objective is to develop a comprehensive retention of the current advancements, pinpoint deficiencies in existing knowledge and lay the groundwork for our significant effort in merging Blockchain technology with AI/ML. This investigation establishes the basis for a comprehensive comprehension of the interconnections and prospective resolutions that emerge at the intersection of these advanced technologies.

The issue of training-set attacks on ML models, where an attacker manipulates the training data to make the model produce desired outcomes is listed by¹⁵. Such attacks are of increasing concern as intelligent agents like spam filters and robots are susceptible to manipulation through their learning capabilities. The study continues by examining possible countermeasures against these assaults on training sets, highlighting the urgent need for strong security measures in the ever-changing domain of intelligent agents equipped with learning capacities. Another paper¹⁶ explores the growing threat of poisoning attacks on ML

models, particularly focusing on linear regression models. In poisoning attacks, adversaries manipulate training data to influence the results of predictive models. A further research¹⁷ introduces a novel method to identify and remove harmful data from training sets used in supervised learning models. It's the first strategy to use data provenance for defense against poisoning attacks. Two versions of the defense are presented for partially trusted and fully entrusted data sets. Experimentally, the method outperforms the baseline, enabling secure use of machine learning models in adversarial environments with reliable provenance data.

With the transition towards machine learning (ML) in computing paradigms, there is an emergence of security risks due to the ability of ML to handle large datasets. This article¹⁸ provides a concise summary of security dangers connected to machine learning, along with the obstacles involved with these threats. Additionally, the document includes examples of vulnerabilities utilizing LeNet and VGGNet on benchmark datasets. In the broader context, R Tomsett et al.¹⁹ highlight the distinctive adversarial and distributed characteristics of the tactical environment, with a specific emphasis on the vulnerability of ML models to adversarial inputs. The review also focuses on the occurrence of model poisoning attacks in distributed learning systems. The growing trend of outsourcing machine learning (ML) training has led to the emergence of potential vulnerabilities that adversaries might exploit²⁰. This is primarily driven by the high costs associated with ML training. One of the notable risks encountered during the training phase is the poisoning attack when attackers seek to compromise machine learning systems by introducing corrupted training data.

A specific emphasis on the contributions of machine learning and blockchain technologies is explored by this chapter²¹ that introduces the dynamic nature of cybersecurity. In recent decades, scholars have extensively investigated many methodologies aimed at enhancing the security of data. Among these techniques, the use of blockchain integrated with machine learning has emerged as a powerful technique that has promise for extensive implementation in practical situations. The objective of this chapter is to critically

examine the impacts of machine learning and blockchain technologies on the improvement of cybersecurity protocols. Another study²² incorporates blockchain technology with AI/ML and presents a significant solution in addressing the security challenges that are inherent in smart cities. This integration offers several features, including immutability, trust-free transactions, and decentralization, which effectively addresses data security and privacy concerns across multiple sectors such as healthcare, agriculture, communication, transport, and smart grids. This review offers a thorough analysis of artificial intelligence and machine learning blockchain solutions for Internet of Things (IoT) communication within the framework of smart cities. A similar study²³ investigates the strategic use of artificial intelligence (AI), blockchain technology, and machine learning to enhance the cybersecurity defenses of neobanks in response to the increasing dangers in this domain. While AI and machine learning enhance threat detection and adaptive responses, blockchain's decentralized structure provides transparency and immutability, making unauthorized access significantly more challenging. The integration of these technologies creates a comprehensive security strategy, enabling proactively addressing emerging risks and prioritizing the confidence of customers and the integrity of data within the digitized financial landscape.

Consensus, a fundamental element of the blockchain network, is achieved through two main approaches: Proof-of-Work (PoW) and Proof-of-Stake (PoS), each offering distinctive advantages germane to the nuanced requirements of AI/ML data management. A recent study²⁴ examines the potential of blockchain technology to reshape various aspects of our lives, highlights the security and performance challenges in PoS consensus protocols, and discusses vulnerabilities. In the context of security, Proof of Work (PoW) is notable for its significant resilience^{25 26 27}, mostly because of the substantial computing effort required for mining. The allocation of significant energy and resources serves as a substantial deterrent against malevolent assaults, hence enhancing the overall security of the system. Moreover, the PoW consensus mechanism effectively facilitates the establishment of decentralization inside the network^{28 29}. The mining process promotes inclusiveness by allowing access to

those possessing the required gear, facilitating a more widespread distribution of power, and mitigating the potential for centralization. PoW has additionally a well-established history, having been used since the emergence of blockchain technology with the introduction of Bitcoin⁸. The persistent and consistent use of this consensus demonstrates its trustworthiness and stability over a significant duration, therefore establishing it as a trustworthy and enduring consensus mechanism within the realm of blockchain technology.

The selection of the cryptographic hash function used to create data blocks is a crucial factor that has broad consequences in the context of blockchain technology applied to ML data. The cryptographic hash function SHA256, which is well-known for its strong security characteristics, is one of the most appealing solutions accessible. Because blockchain and machine learning require large datasets to be handled safely and independently, using SHA256 is particularly advantageous^{30 31 32}. It can ensure data integrity and cryptographic security, speed up processing, and be widely accepted, making SHA256 integration necessary to create a trustworthy and safe blockchain ecosystem that works well with machine learning applications^{33 34}. This paper³⁵ advocates for the adoption of blockchain technology, specifically emphasizing the utilization of SHA256 to ensure entirely secure communication between devices. The proposed technique involves a two-step process: first, authenticating each node using identity-based signatures for secure communication on a blockchain platform; and second, sending blocks through hashing using SHA256 with device identities as public keys. According to theoretical analysis and simulation results, this method performs better than S-LoRaWAN and DLBA-IoT in terms of average detection rate, throughput, scalability, time spent on block authentication, and energy use. The findings of³⁶ underscore how the proposed approach, leveraging SHA256, significantly enhances the security of individual devices and overall network security within the IoT ecosystem. A research highlights the effectiveness of blockchain technology in enhancing patient privacy and security within the context of the Internet of Medical Things³⁷. Traditional privacy techniques prove inadequate, prompting the adoption of blockchain technology with smart contracts

to protect private patient records. Utilizing the SHA256 algorithm for data integrity, the system efficiently creates and secures blocks.

All this existing research primarily concentrates on aspects like fairness and cooperative incentives, often overlooking efficiency and feasibility. Within this perspective, we present a novel blockchain-driven framework aligning with the core roles of machine learning models. From multiple data aggregators collaborate within an aggregator consortium, ensuring efficient, secure, and verifiable data sharing among peers throughout the training phase, a pivotal emphasis is placed on safeguarding data integrity and security.

Chapter 3

Methodology and Framework

The preservation of the reliability of AI/ML models, particularly during their training phase, is crucial to mitigating adversarial attacks, data poisoning, and upholding the confidentiality of sensitive data. Traditional approaches, such as the combination of Proof of Work (PoW) and the SHA256 cryptographic hash function with blockchain technology, are gaining attention for their potential to enhance security and integrity during the training cycle. Combining PoW and SHA256 with blockchain technology, which is known for being decentralized, immutable, and hard to change, could help solve these problems by improving the security and integrity of data during the training cycle.

In the present context, the proposed framework initiates an investigation into the potential collaboration between blockchain, using PoW as a consensus mechanism, and SHA256 to create signatures. This exploration delves into how this integration can augment trust and dependability during the training phase of ML models. Consequently, it paves the way for a novel era characterized by secure, transparent, and resilient AI/ML applications.

Figure 3.1 provides a comprehensive illustration of the structural design of the proposed framework. It clearly outlines the key components, namely the peers that form the blockchain network, who are responsible for the crucial role of validating data. The initial step involves acquiring the input dataset and submitting it to preprocessing. Subsequently,

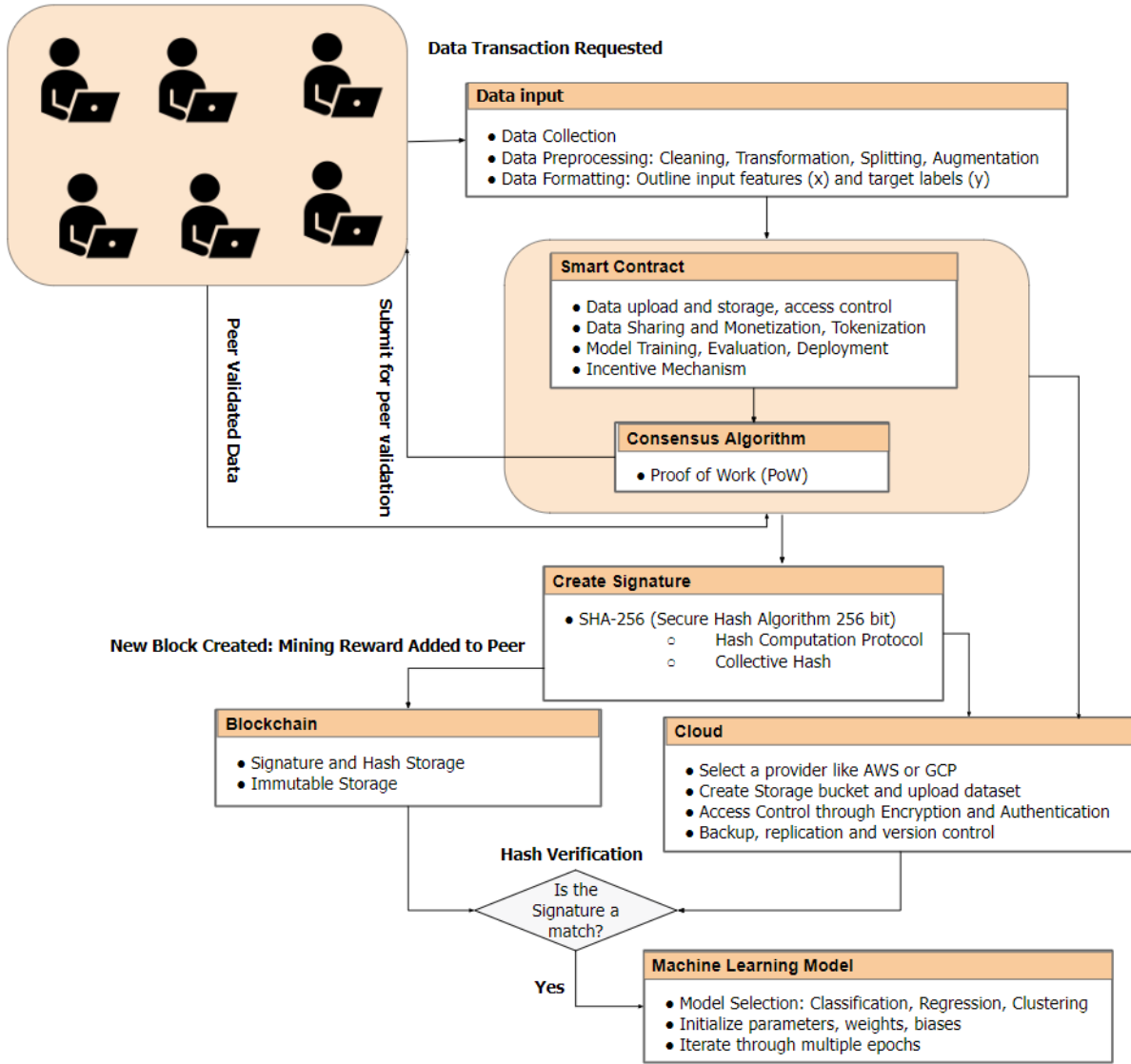


Figure 3.1: Proposed system architecture for creating data blocks to train AI/ML models.

a series of data preparation steps are executed, including thorough data cleaning, transformation, segmentation, and augmentation, all aimed at elevating data quality and utility. Following this, the data is methodically formatted, setting input features (x) and target labels (y), thus structuring the data to optimize its appropriateness for subsequent computational tasks and analytical processes. Once the input data has through a thorough preprocessing step, it is then channeled to the smart contract for the purpose of validation

and the creation of a data block. The smart contract plays a crucial role in verifying the integrity and authenticity of the processed data. Once validated through PoW, the data is then encapsulated into a data block, where it is secured and recorded within the blockchain's immutable ledger. This approach not only ensures the trustworthiness of the data but also strengthens the transparency and tamper-resistant attributes of the blockchain network, making it an ideal choice for safeguarding data in ML models.

After the verification of a block by peers, a signature will be produced via SHA256, resulting in the creation of a collective hash. Following this, the hash will be recorded as a new item in the blockchain ledger. In recognition of their participation, mining rewards will be distributed to the peers who have contributed. Simultaneously, the primary data block, together with its corresponding hash, will be stored securely in a private cloud infrastructure, such as AWS or GCP³⁸. This infrastructure will implement rigorous processes, including data backup, replication, and version control, to ensure the integrity and reliability of the data. Following the validation of the data consistency with the verified hash stored in the ledger, it is prepared for feeding into the ML model. This process establishes a full cycle of data integrity and security within the framework.

Machine Learning models, which heavily depend on securely trained data produced inside a complex framework, provide an unparalleled degree of reliability and robustness in the face of adversarial assaults. The comprehensive incorporation of SHA256 not only enhances the security protocols of the training process but also strengthens the model's resistance to any weaknesses, therefore establishing a resilient framework for the construction of reliable and secure Artificial Intelligence/Machine Learning (AI/ML) systems.

The workflow diagram of 3.2 illustrates the validation procedure for a new dataset inside the blockchain network. The validation process begins with creating the first block header and then proceeds to use the PoW consensus algorithm, in which miners diligently execute computing jobs to attain the necessary threshold for validating the block via the mining process. The miners, who are recognized as peers in the network, access the homepage

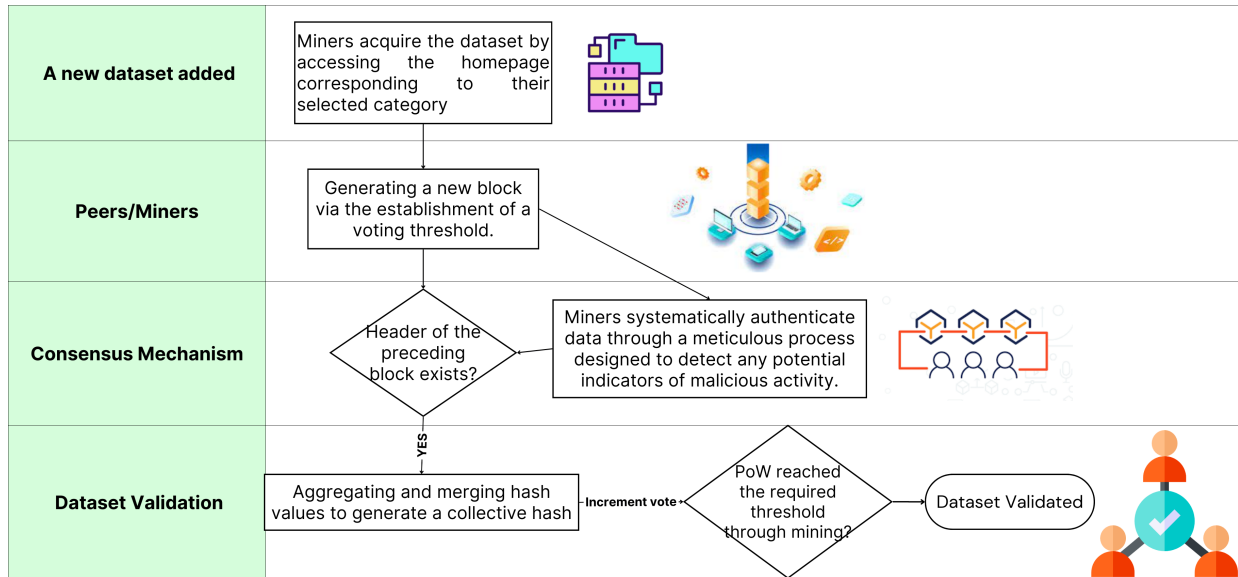


Figure 3.2: *Work flow diagram*

linked to their respective categories in order to get the uploaded dataset. Following acquisition, the dataset undergoes meticulous validation, systematically authenticated by miners to detect any potential indicators of malicious activity or data inconsistencies. The process advances with the aggregation and merging of hash values associated with the validated dataset, generating a collective hash that serves as a unique identifier. The compilation of this collective data contributes to the creation of a new block inside the blockchain, contingent upon the establishment of a voting threshold to secure consensus among miners. Throughout these steps, the workflow highlights the methodical and systematic approach to data validation, placing significant emphasis on the dedication to identifying and resolving potential malicious data alteration in order to maintain the integrity of the blockchain network.

Chapter 4

Framework Development

Through our web application, SecureChainAI, a cutting-edge software dedicated to boosting ML and AI trust and dependability, users can experience a revolution in data security. Our platform harnesses the power of blockchain to guarantee data integrity and security, offering a state-of-the-art solution for safeguarding information. Using a novel method for dataset validation, SecureChainAI combines blockchain technology with crowdsourcing. Through the collective intelligence of a diverse peer network, our platform ensures a meticulous validation process for uploaded datasets, establishing a robust foundation for data integrity.

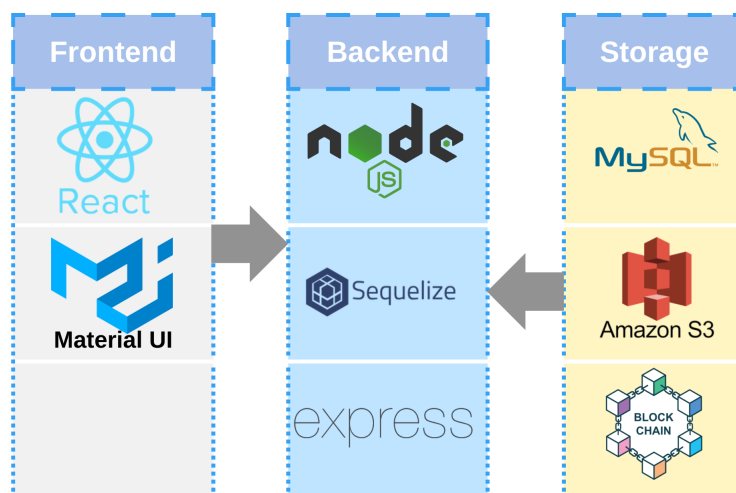


Figure 4.1: *Application Architecture*

By developing this novel framework as the architecture shown in figure 4.1, individuals are able to effortlessly complete the processes of registration, login, dataset uploading, category subscription, dataset validation, and exploration of the extensive blockchain history. These actions are facilitated by an interactive interface that has been specifically created to provide a smooth and user-friendly experience. Presented below is a tabular representation, including the technologies and functionality that have been outlined.

Aspect	Technology	Functionality
Frontend	React, Material UI	Modern and consistent design, enhancing user interface
Backend	NodeJS (Express)	Robust and scalable backend infrastructure
Database	MySQL, Sequelize (ORM)	Simplifies database interactions, ensures a structured and organized database architecture
Storage	AWS S3	Scalable, durable, and accessible storage for uploaded datasets
HTTP Requests	Axios	Streamlines communication between the frontend and backend
Authentication	JWT (JSON Web Token)	Provides secure, stateless authentication mechanism for user identification
Password Security	Bcrypt	Secures user passwords

Table 4.1: *Summary of Technology Overview and Corresponding Features*

The use of these technologies ensures a robust, scalable, and secure platform for dataset validation and approval. The Material UI elevates the user interface with a modern and consistent design; Sequelize simplifies database interactions; Axios streamlines HTTP communication; Bcrypt secures user passwords; and JWT ensures secure authentication. Furthermore, presented below is a concise summary of the functions facilitated by these technologies:

- **User Registration and Authentication:** Users can register and log in to the application, ensuring a personalized and secure experience. JWT (JSON Web Token) is employed

to manage user authentication, providing a stateless, secure mechanism for user identification.

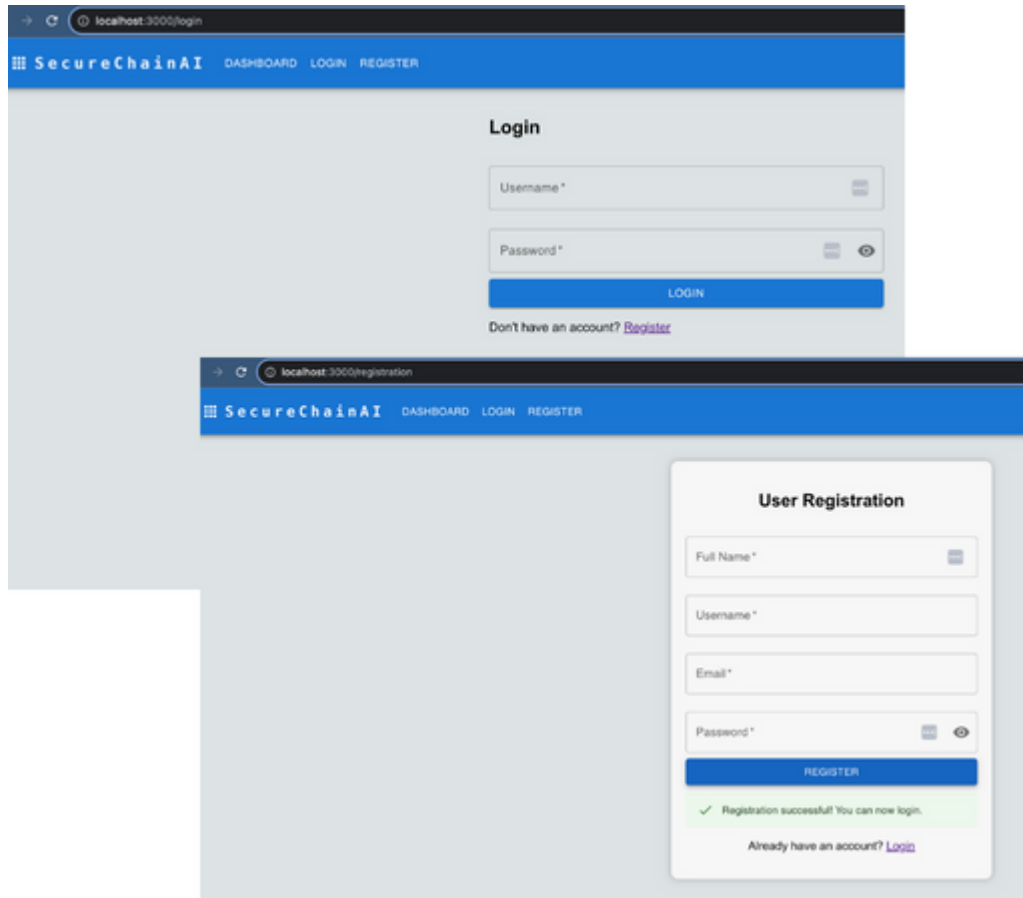


Figure 4.2: *User Authentication*

- **Category Subscription:** After logging in, users can subscribe to specific categories of interest. This feature allows users to focus on datasets that align with their expertise and preferences, creating a more personalized experience, as stated in figure 4.3. The subscription function serves to augment user personalization and ensures the delivery of relevant updates and information according to the user's chosen areas of interest within the system.

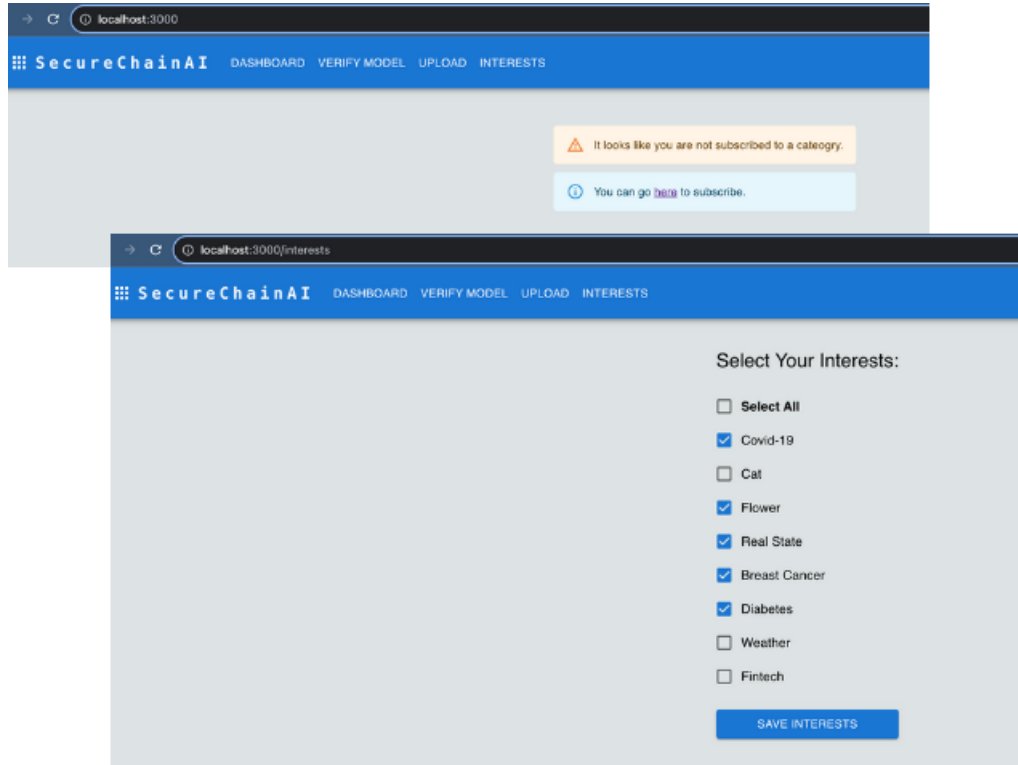


Figure 4.3: *Category of Interest Subscription by Users*

- User Dashboard: Our platform provides customers with an enhanced experience by offering a customized interface. This interface includes a unique dashboard that enables users to effortlessly monitor and follow the current state of datasets they have uploaded in real-time. The customizable dashboard functions as a single center, offering a thorough overview of the validation process for every dataset. By providing visibility into the current status, progress, and any relevant updates, users gain a heightened sense of accountability throughout the dataset validation process. This feature not only improves user interaction but also establishes an open and responsible method for validating datasets, hence promoting confidence and trust in the platform's operations.

Category	Status	File	Uploaded by	Your Vote	Created At	Consensus Ledger
Breast Cancer	voting	breast-cancer-data.csv	Aphrodite Zamora	VALID INVALID	NOV 26, 2023, 5:02 PM	
Diabetes	voting	diabetes-disease.csv	Shihab Jamil	VALID INVALID	NOV 26, 2023, 2:54 PM	
Diabetes	voting	Healthcare-Diabetes.csv	Shihab Jamil	VALID INVALID	NOV 26, 2023, 2:49 PM	
Real State	voting	real-state-price-predictor-data.csv	Aiasha Siddika	VALID INVALID	NOV 25, 2023, 10:02 PM	
Breast Cancer	approved	breast-cancer-data.csv	Zena Bridges	MISSED	NOV 25, 2023, 4:39 PM	SHOW JSON
Breast Cancer	approved	properties.csv	Aiasha Siddika	MISSED	NOV 23, 2023, 5:04 PM	SHOW JSON
Real State	approved	real-state-dataset.zip	Aiasha Siddika	MISSED	NOV 22, 2023, 5:11 PM	SHOW JSON
Diabetes	approved	diabetes (1).csv	Shihab Jamil	MISSED	NOV 22, 2023, 5:08 PM	SHOW JSON

Figure 4.4: *Subscribed Dataset Dashboard*

- **Dataset Management:** Users can upload datasets, which are stored in AWS S3 buckets, ensuring scalability, durability, and accessibility. Axios is utilized to handle HTTP requests, facilitating smooth communication between the frontend and backend. Users possess the capacity to both upload datasets and download them, facilitating a smooth and reciprocal exchange of data. The user-friendly interface, as seen in Figures 4.5, figure 4.6 and figure 4.7, allows users to effectively arrange datasets according to chronological order, hence promoting a systematic and methodical approach to data administration. Furthermore, users have the convenient ability to verify the voting status of any dataset, promoting transparency within the validation process. The diverse range of capabilities provided by this system enables users to effectively manage and monitor datasets at any stage of their lifespan, therefore equipping them with a complete set of tools.

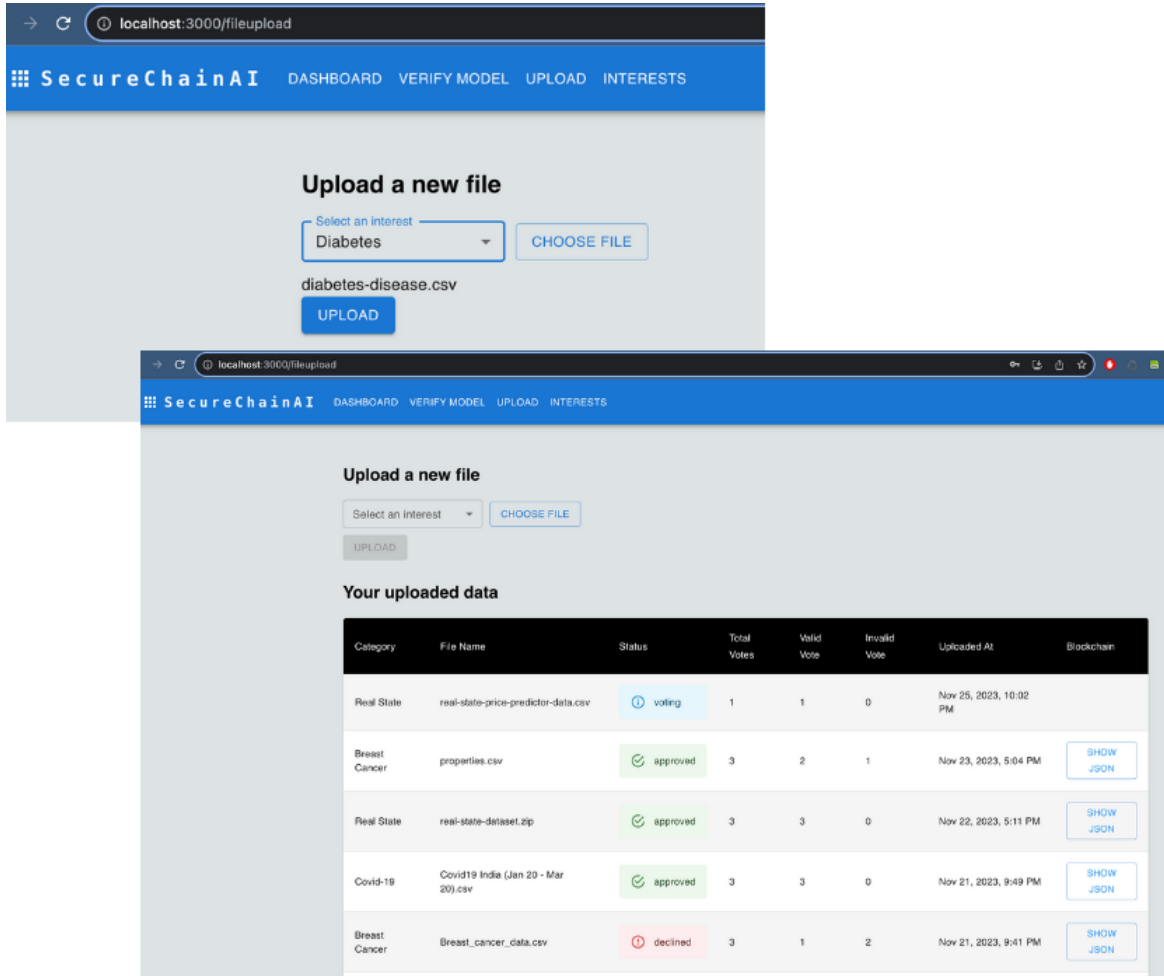


Figure 4.5: *Dataset Upload*

- **Data Validation, Voting and Monitoring System:** The purpose of the Data Validation, Voting, and Monitoring System is to facilitate the active participation of users who have subscribed to relevant categories in a collaborative process. The individuals possess the capability to acquire datasets and evaluate their reliability, engaging in a collaborative decision-making structure. The homepage of the system presents a well-organized presentation of datasets in a tabular format. These datasets are carefully filtered to match the user's subscriptions, resulting in a concise and user-friendly summary.

The implementation of a monitoring system, using NodeJS with Express, effectively

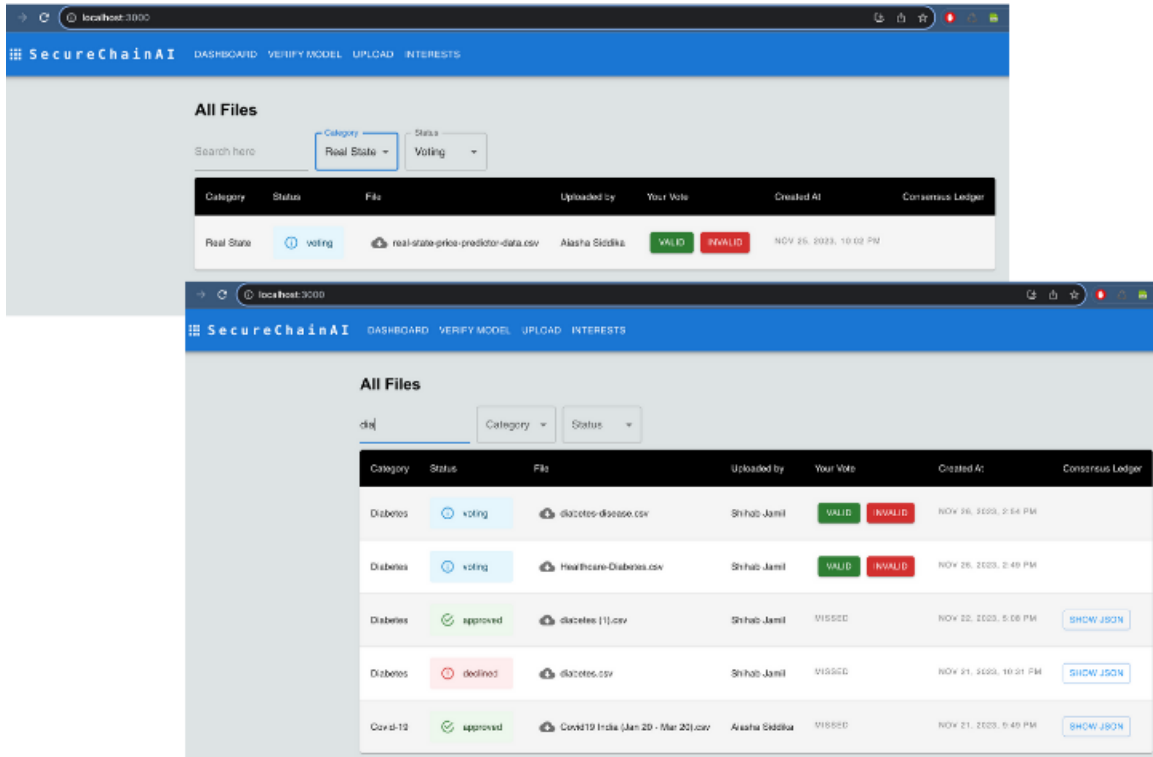


Figure 4.6: Data download for voting and validation

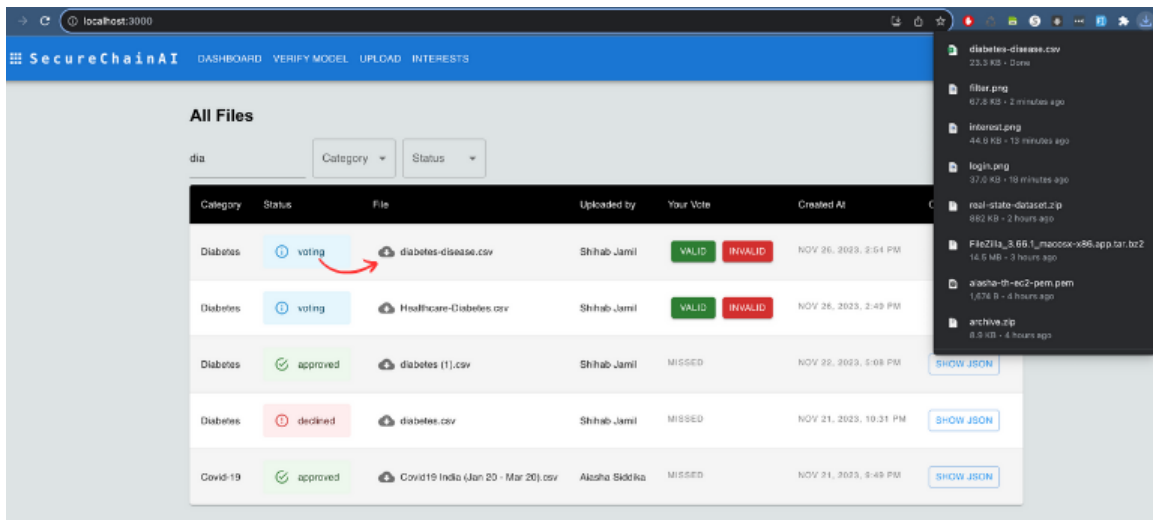


Figure 4.7: Applying filter and search

ensures transparency and efficiency by closely overseeing the whole process of dataset validation. The technology automatically categorizes datasets as accepted or denied

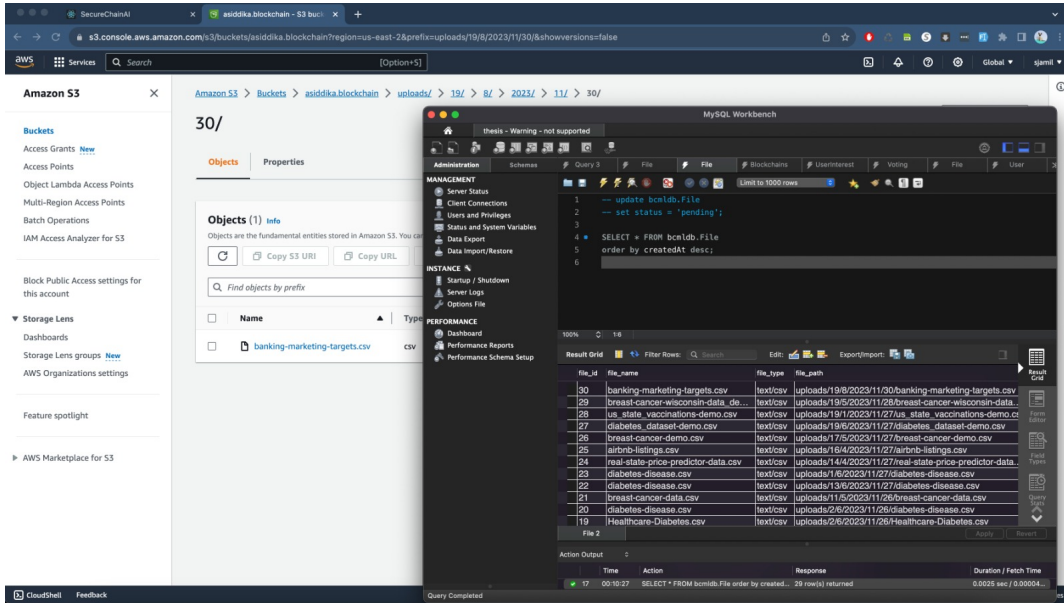


Figure 4.8: Dataset stored in cloud storage

after they have received a certain number of votes. The use of this automated decision-making mechanism plays a crucial role in promoting transparency, optimizing workflow, and improving efficiency in the process of validating datasets. As a result, it enhances the overall efficacy of the system.

- **Blockchain Integration:** After the dataset has been authorized, a SHA256 hash is algorithmically constructed by analyzing its contents. The hash, which functions as a distinct and safe representation of the dataset's information, is then stored with great attention to detail. The storing method described above results in the creation of a full blockchain history. This history serves as an unchangeable and transparent record of the dataset's development and the trip it undergoes for approval. Users are granted the capacity to explore the historical records of this blockchain, enhancing the degree of transparency and traceability for every authorized dataset. This, in turn, fosters confidence in the data's authenticity and integrity.
- **Integrity Verification:** To ensure data integrity, users can verify datasets before feeding

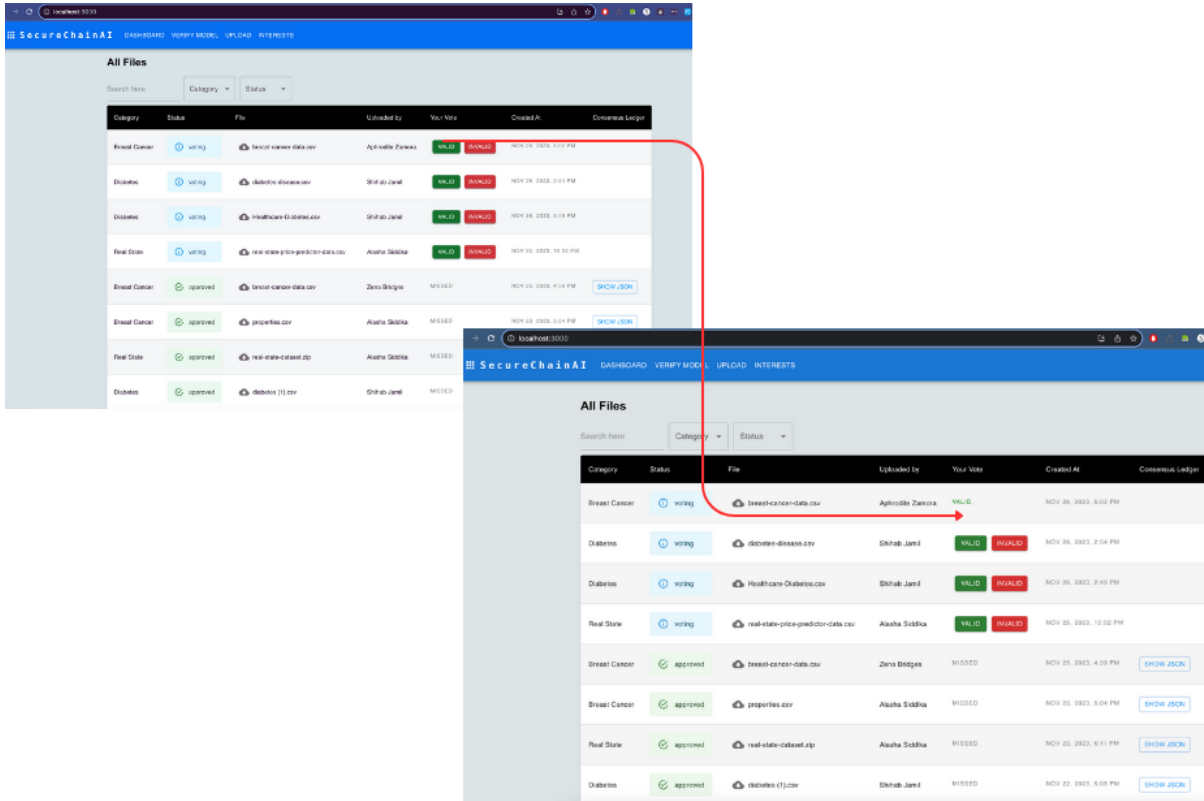


Figure 4.9: Dataset Verification and Voting Status Evaluation

them into machine learning models. The app compares the stored hash (on approval) with a run-time hash generated from the content of the file stored in S3, indicating whether the data is pure or compromised shown in figure 4.11.

After validation is complete, SecureChainAI produces a unique fingerprint or hash, which becomes an unchangeable proof of validity for the dataset. This cryptographic seal identifies the dataset as a trustworthy and validated source of information and acts as an unchangeable proof of its validity. By using a decentralized public Ethereum blockchain to record every dataset’s complete history, ensuring transparency and traceability. This includes participant actions and timestamps. As part of our commitment to data security, SecureChainAI thoroughly verifies the integrity of datasets before incorporating them into machine learning pipelines. By comparing the hash that is created with the one that is saved, this dual-check process makes sure that the dataset

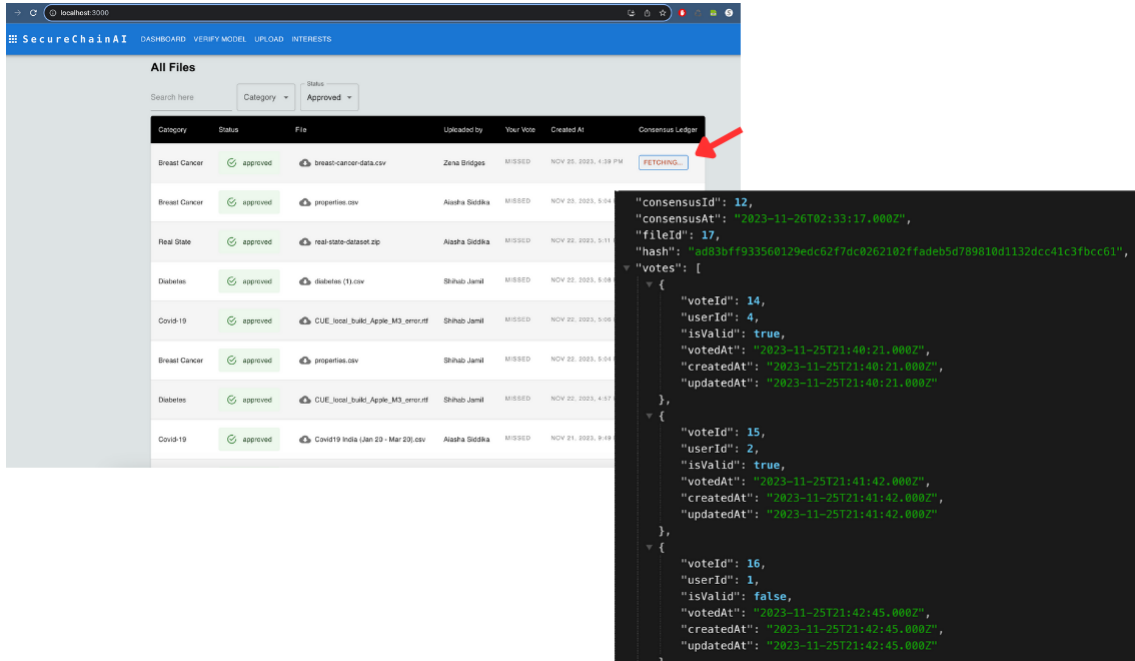


Figure 4.10: Consensus Ledger ensuring transparency and traceability

is safe and unaltered over its entire life on the cloud.

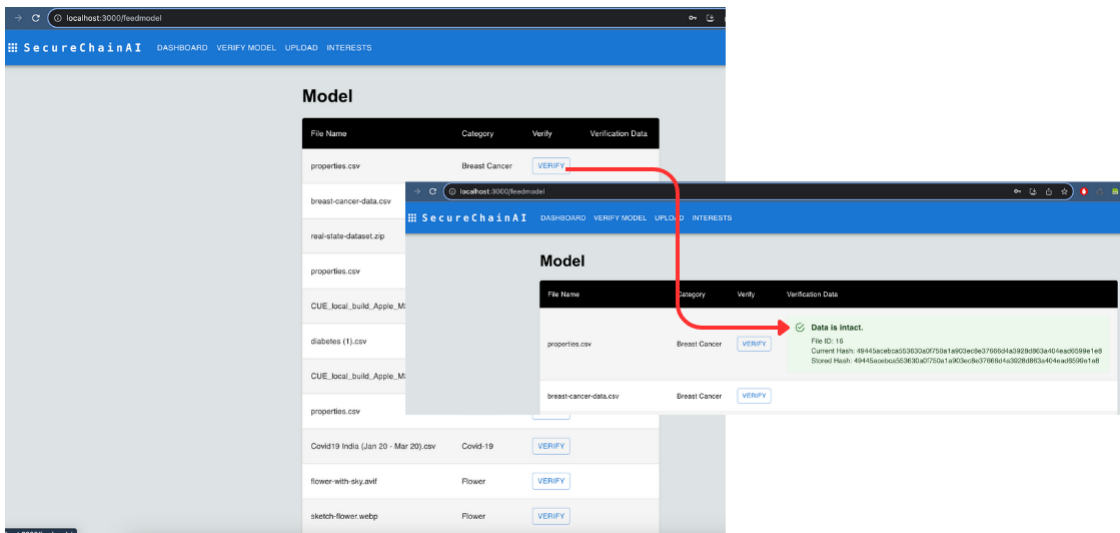


Figure 4.11: Data Integrity Verification Process

Therefore, the implementation of SecureChainAI occupies a leading position in the field of data security, using blockchain technology and crowdsourcing methods to ensure thor-

ough validation of datasets. This application significantly transforms the trustworthiness of AI/ML, establishing new benchmarks for ensuring data integrity and reliability.

Chapter 5

Recommendations and Future Prospects

Our forthcoming research endeavors will entail the comprehensive examination of this integrated framework across a multitude of industry domains, including but not limited to healthcare and finance, with the objective of appraising its applicability, resilience, and scalability within authentic, real-world contexts. Furthermore, the objective is to develop user interfaces that exhibit qualities of intuitiveness and user-friendliness, with the purpose of enhancing seamless interactions with decentralization.

5.1 Insights and Roadmap for the Future

5.1.1 Integration of Proof of Federated Learning (PoFL)

Though our current research has shed light on existing consensus mechanisms like PoW, Proof of Federated Learning (PoFL) presents a potential field for further research. PoFL is a noteworthy alternative that introduces a decentralized and collaborative learning paradigm introduced by Konečný et al.³⁹. It enables cooperative calculations for model updates

without disclosing raw data by integrating cryptographic approaches like Secure Multi-Party Computation (SMPC)⁴⁰. Because PoFL prioritizes both collaborative model improvement and data safety, it is a viable option for improving ML/AI data security. When FL is used⁴¹, people who share resources (often called miners or peers) work together to train models without implying a central storage system. Orchestrating this decentralized process is the responsibility of the pool manager.

5.1.2 Collaborative Secure Computing Integration

As a key element of ensuring that participants can collaboratively perform computations on the model while preserving the privacy of their individual data, the utilization of SMPC in the initial phase could safeguard privacy in the computation process. The aggregation process, integral to the collective advancement of the model, maintains the privacy of individual data sources. The ensuing modifications are systematically recorded in a secure and transparent manner within blockchain ledgers, making them resilient to tampering. As we conclude our exploration into model security, a noteworthy future recommendation is the incorporation of the Secure Model Verification (SMV) step, proposed by the author⁴², seamlessly integrated into the Cloud Service Provider (CSP). This procedure's careful design guarantees that model predictions are carried out within the CSP with the highest confidentiality and integrity.

Additionally, integrating the Function Secret Sharing (FSS) protocol⁴³ with SMPC offers an intriguing path for future research. With this protocol, both the hosting firm and the miner may remain anonymous, which is an attractive proposition. The inclusion of FSS greatly enhances the overall privacy and security of the environment where model training occurs by protecting each participant's inputs and models throughout the collaborative learning process. Investigating how the FSS protocol is implemented and affected inside the SMPC framework seems like a wise course of action for improving security protocols in the context of cooperative learning settings.

Chapter 6

Conclusion

This research highlights the potential of Blockchain technology to revolutionize the training cycle of ML models with enhanced data security. Through the exploration of fundamental research inquiries, we have successfully illuminated the merits, quantifiable gains, and approaches to overcome obstacles in the process of implementation. Our comprehensive framework not only fortifies ML model security but also fosters transparency and trust in blockchain-based, decentralized learning environments, forging a robust synergy by integrating PoW and SHA256. As technological advancements continue, the integration of blockchain technology with the training of ML models is a noteworthy development that underscores our dedication to protecting digital assets within the ever-evolving cyber environment.

Bibliography

- [1] Rahul Rai, Manoj Kumar Tiwari, Dmitry Ivanov, and Alexandre Dolgui. Machine learning in manufacturing and industry 4.0 applications. *International Journal of Production Research*, 59(16):4773–4778, 2021. doi: 10.1080/00207543.2021.1956675.
- [2] Marcus Comiter. Attacking artificial intelligence. *Belfer Center Paper*, 8:2019–08, 2019.
- [3] Victor Galaz, Miguel A. Centeno, Peter W. Callahan, Amar Causevic, Thayer Patterson, Irina Brass, Seth Baum, Darryl Farber, Joern Fischer, David Garcia, Timon McPhearson, Daniel Jimenez, Brian King, Paul Larcey, and Karen Levy. Artificial intelligence, systemic risks, and sustainability. *Technology in Society*, 67:101741, 2021. ISSN 0160-791X. doi: <https://doi.org/10.1016/j.techsoc.2021.101741>. URL <https://www.sciencedirect.com/science/article/pii/S0160791X21002165>.
- [4] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples, 2015.
- [5] Xianmin Wang, Jing Li, Xiaohui Kuang, Yu-an Tan, and Jin Li. The security of machine learning in an adversarial setting: A survey. *Journal of Parallel and Distributed Computing*, 130:12–23, 2019.
- [6] Henry Chacon, Samuel Silva, and Paul Rad. Deep learning poison data attack detection. In *2019 IEEE 31st International Conference on Tools with Artificial Intelligence (ICTAI)*, pages 971–978, 2019. doi: 10.1109/ICTAI.2019.00137.
- [7] Ali Sayghe, Yaodan Hu, Ioannis Zografopoulos, XiaoRui Liu, Raj Gautam Dutta, Yier Jin, and Charalambos Konstantinou. Survey of machine learning methods for detecting false data injection attacks in power systems. *IET Smart Grid*, 3(5):581–595, 2020.

- [8] Satoshi Nakamoto. Bitcoin: a peer-to-peer electronic cash system. <https://bitcoin.org/bitcoin.pdf>, 2008.
- [9] Claudia Antal, Tudor Cioara, Ionut Anghel, Marcel Antal, and Ioan Salomie. Distributed ledger technology review and decentralized applications development guidelines. *Future Internet*, 13(3), 2021. ISSN 1999-5903. doi: 10.3390/fi13030062. URL <https://www.mdpi.com/1999-5903/13/3/62>.
- [10] Khaled Salah, M. Habib Ur Rehman, Nishara Nizamuddin, and Ala Al-Fuqaha. Blockchain for ai: Review and open research challenges. *IEEE Access*, 7:10127–10149, 2019. doi: 10.1109/ACCESS.2018.2890507.
- [11] Yiming Liu, F. Richard Yu, Xi Li, Hong Ji, and Victor C. M. Leung. Blockchain and machine learning for communications and networking systems. *IEEE Communications Surveys Tutorials*, 22(2):1392–1431, 2020. doi: 10.1109/COMST.2020.2975911.
- [12] Muhammad Shafay, Raja Wasim Ahmad, Khaled Salah, Ibrar Yaqoob, Raja Jayaraman, and Mohammed Omar. Blockchain for deep learning: review and open challenges. *Cluster Computing*, 26(1):197–221, 2023.
- [13] Yakov Vainshtein and Ehud Gudes. Use of blockchain for ensuring data integrity in cloud databases. In *Cyber Security Cryptography and Machine Learning: 5th International Symposium, CSCML 2021, Be'er Sheva, Israel, July 8–9, 2021, Proceedings 5*, pages 325–335. Springer, 2021.
- [14] Turing. Does artificial intelligence impact blockchain technology?, 2022. URL <https://www.turing.com/kb/does-artificial-intelligence-impact-blockchain-technology>. Accessed on November 21, 2023.
- [15] Shike Mei and Xiaojin Zhu. Using machine teaching to identify optimal training-

- set attacks on machine learners. In *Proceedings of the aaai conference on artificial intelligence*, volume 29, 2015.
- [16] Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, and Bo Li. Manipulating machine learning: Poisoning attacks and countermeasures for regression learning. In *2018 IEEE symposium on security and privacy (SP)*, pages 19–35. IEEE, 2018.
- [17] Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, and Jaehoon Amir Safavi. Mitigating poisoning attacks on machine learning models: A data provenance based approach. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pages 103–110, 2017.
- [18] Faiq Khalid, Muhammad Abdullah Hanif, Semeen Rehman, and Muhammad Shafique. Security for machine learning-based systems: Attacks and challenges during training and inference. In *2018 International Conference on Frontiers of Information Technology (FIT)*, pages 327–332. IEEE, 2018.
- [19] Richard Tomsett, Kevin Chan, and Supriyo Chakraborty. Model poisoning attacks against distributed machine learning systems. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, volume 11006, pages 481–489. SPIE, 2019.
- [20] Zhibo Wang, Jingjing Ma, Xue Wang, Jiahui Hu, Zhan Qin, and Kui Ren. Threats to training: A survey of poisoning attacks and defenses on machine learning systems. *ACM Computing Surveys*, 55(7):1–36, 2022.
- [21] Rohini Patil, Monika Mangla, and Smita Bansod. Scope of machine learning and blockchain in cyber security. In *Intelligent Approaches to Cyber Security*, pages 35–53. Chapman and Hall/CRC.

- [22] Badal Gami, Manav Agrawal, Deepak Kumar Mishra, Danish Quasim, and Pawan Singh Mehra. Artificial intelligence-based blockchain solutions for intelligent healthcare: A comprehensive review on privacy preserving techniques. *Transactions on Emerging Telecommunications Technologies*, 34(9):e4824, 2023.
- [23] A Shaji George. Securing the future of finance: How ai, blockchain, and machine learning safeguard emerging neobank technology against evolving cyber threats. *Partners Universal Innovative Research Publication*, 1(1):54–66, 2023.
- [24] Gabriel Antonio F Rebello, Gustavo F Camilo, Lucas CB Guimaraes, Lucas Airam C de Souza, Guilherme A Thomaz, and Otto Carlos MB Duarte. A security and performance analysis of proof-based consensus protocols. *Annals of Telecommunications*, pages 1–21, 2021.
- [25] Marc Pilkington, Rodica Crudu, and Lee Gibson Grant. Blockchain and bitcoin as a way to lift a country out of poverty-tourism 2.0 and e-governance in the republic of moldova. *International Journal of Internet Technology and Secured Transactions*, 7(2): 115–143, 2017.
- [26] Rahela Lokman Hemashrif. An analysis of decentralized finance and its applications. Master’s thesis, uis, 2021.
- [27] Mengjiang Liu, Qianhong Wu, Yiming Hei, and Dawei Li. Blockchain-based licensed spectrum fair distribution method towards 6g-envisioned communications. *Applied Sciences*, 13(16):9231, 2023.
- [28] Mousa Mohammed Khubrani and Shadab Alam. Blockchain-based microgrid for safe and reliable power generation and distribution: A case study of saudi arabia. *Energies*, 16(16):5963, 2023.
- [29] George R Lucas. *Ethics and cyber warfare: the quest for responsible security in the age of digital warfare*. Oxford University Press, 2017.

- [30] Rahul Johari, Vivek Kumar, Kalpana Gupta, and Deo Prakash Vidyarthi. Blossom: Blockchain technology for security of medical records. *ICT Express*, 8(1):56–60, 2022.
- [31] P Velmurugadass, S Dhanasekaran, S Shasi Anand, and V Vasudevan. Enhancing blockchain security in cloud computing with iot environment using ecies and cryptography hash algorithm. *Materials Today: Proceedings*, 37:2653–2659, 2021.
- [32] Sheping Zhai, Yuanyuan Yang, Jing Li, Cheng Qiu, and Jiangming Zhao. Research on the application of cryptography on the blockchain. In *Journal of Physics: Conference Series*, volume 1168, page 032077. IOP Publishing, 2019.
- [33] Raffaele Martino and Alessandro Cilardo. Designing a sha-256 processor for blockchain-based iot applications. *Internet of Things*, 11:100254, 2020.
- [34] Rizwan Malik, Hammad Raza, Muhammad Saleem, et al. Towards a blockchain enabled integrated library management system using hyperledger fabric: Using hyperledger fabric. *International Journal of Computational and Innovative Sciences*, 1(3):17–24, 2022.
- [35] Reza Fotohi and Fereidoon Shams Aliee. Securing communication between things using blockchain technology based on authentication and sha-256 to improving scalability in large-scale iot. *Computer Networks*, 197:108331, 2021.
- [36] KN Devika and Ramesh Bhakthavatchalu. Parameterizable fpga implementation of sha-256 using blockchain concept. In *2019 International Conference on Communication and Signal Processing (ICCSP)*, pages 0370–0374. IEEE, 2019.
- [37] Ibrahim Shawky Farahat, Waleed Aladrousy, Mohamed Elhoseny, Samir Elmougy, and Ahmed Elsaid Tolba. Improving healthcare applications security using blockchain. *Electronics*, 11(22):3786, 2022.
- [38] Manish Saraswat and R.C. Tripathi. Cloud computing: Comparison and analysis of cloud service providers-aws, microsoft and google. In *2020 9th International Conference*

System Modeling and Advancement in Research Trends (SMART), pages 281–285, 2020.
doi: 10.1109/SMART50582.2020.9337100.

- [39] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- [40] Youyang Qu, Md Palash Uddin, Chenquan Gan, Yong Xiang, Longxiang Gao, and John Yearwood. Blockchain-enabled federated learning: A survey. *ACM Computing Surveys*, 55(4):1–35, 2022.
- [41] Guanming Bao and Ping Guo. Federated learning in cloud-edge collaborative architecture: key technologies, applications and challenges. *Journal of Cloud Computing*, 11(1):94, 2022.
- [42] Aditya Pribadi Kalapaaking, Ibrahim Khalil, and Xun Yi. Blockchain-based federated learning with smpc model verification against poisoning attack for healthcare systems. *IEEE Transactions on Emerging Topics in Computing*, pages 1–11, 2023. doi: 10.1109/TETC.2023.3268186.
- [43] Aditya Pribadi Kalapaaking, Veronika Stephanie, Ibrahim Khalil, Mohammed Atiquzaman, Xun Yi, and Mahathir Almashor. Smpc-based federated learning for 6g-enabled internet of medical things. *IEEE Network*, 36(4):182–189, 2022.