# Decentralization of Big Data with Blockchain

**Venkata Naga Sai Kiran Challa**

**Abstract:** *Based on these concepts, this paper investigates how Blockchain technology can be integrated with Big Data and Machine Learning as a decentralized storage and processing. This paper seeks to capitalize on Blockchain technology's decentralized model for better and more secure data handling and Machine Learning for improved data analysis. Some of the discoveries are the opportunities for future automation by smart contracts, the increased privacy of groups' data by federated learning, and Machine Learning algorithms for more efficient anomaly detection. It solves critical problems of centralized data systems like a centralized point of failure, vulnerability to cyberattacks, and slow data handling. The findings of this research also have significant implications for managing big data using a new and better approach that is more effective and efficient in terms of scalability and security. Nevertheless, integrating these technologies for data management and analysis may present challenges in complexity, computational overhead, and scalability; however, the advantages of applying these technologies have great potential for future developments.*

**Keywords:** Blockchain, Big Data, Machine Learning, decentralized storage, data security

## 1. Introduction

**Background**

Big Data is a term used to describe the large amount of data produced in real - time from various sources like social media, sensors, and transactional systems. This is high - volume data, often called big data, due to its velocity and variety. It makes it valuable in providing insights to help carry out business and organization operations for the accomplishment of a business decision. Big Data Analysis facilitates the possibilities of finding patterns, trends, and connections, which, when used to make good decisions, can be of help in the enhancement of production processes, customer treatment, and the achievement of business advantage. Nonetheless, Big Data is a much broader area of coverage, and therefore, it becomes difficult to solve for storage, processing, and security. Centralized ownership does not cope with the amount of traffic by growing its capacities and is therefore vulnerable to clogging and hacking attempts. It is however still not easy to maintain data's 'innocence' and 'confidentiality' in the present world plagued with increasing cyber threats. Blockchain integration and Machine Learning are two of the best solutions to these challenges that are currently available in the advancement of technology. Blockchain makes the data structure distributed and immutable, making it more reliable and open to the public. At the same time, machine learning is a tool for making prognoses based on the received data. In combination, these technologies can contribute to the radical overhaul of Big Data management, storage, and utilization processes while doing away with the shortcomings of the current systems of Big Data implementation that are based on centralization.
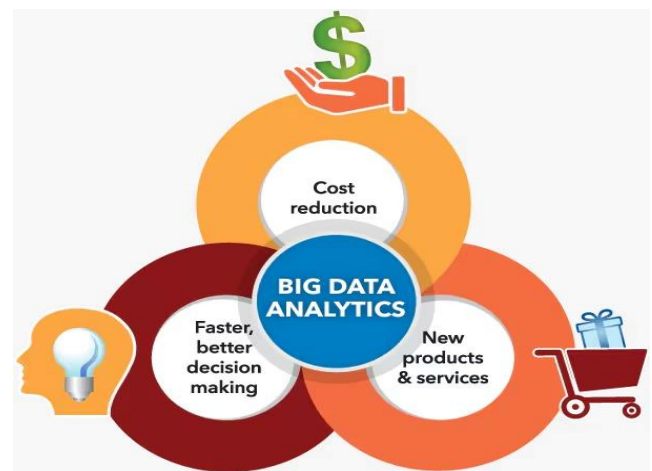


**Figure 1:** Big Data Analytics

**Blockchain Technology**

Blockchain is a distributed and tamper - proof database that allows the recording of numerous varieties of transactions across many computers in a way that cannot be changed. Cryptographic principles are used; it is open to public viewing, and it is not owned or controlled by any one individual. It is an architecture that guarantees that when a specific transaction is written, it cannot be changed without changing all subsequent blocks, which would require the consensus of most of the network. Blockchain is based on the cryptographic security, openness, and distribution of records in a network. The data's privacy is maintained by using cryptographic algorithms in that only the right people are allowed to make alterations to the data or have access to it. Sustainability is another benefit because copies of the ledger are distributed to all the participants in the network, and it is possible to prevent the distortion of the records. System decentralization means the absence of a central control point, which may lead to the occurrence of system failure and enhance reliability. This technology has been embraced in other fields apart from the financial markets, where it began as the means of transaction in cryptocurrencies; it has found its use in areas such as manufacturing and supply chain, the medical field, and banking, among others, where the issue of traceability of transactions is of great importance. Blockchain gives a high level of security and transparency to transactions because it allows the creation of a history of transactions that cannot be altered, which makes it suitable for improving the Big Data system. This is accompanied by consensus algorithms like the Proof of Work (PoW) or the Proof of Stake

(PoS) that guarantee that the network approves all transactions to minimize fraudulent activities. These mechanisms protect the data from being changed by any malicious person because changing a single block would push the consensus for all the subsequent blocks, and it is computationally infeasible. Additionally, Blockchain incorporates smart contracts that contract with code programmed directly into the contract, which can self - execute, thus requiring less intervention from third parties, leading to enhanced functionality and cost - effectiveness.



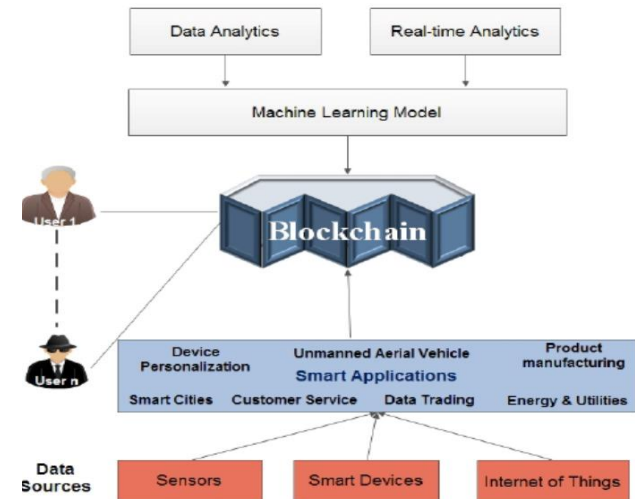**Figure 2:** Proof of Work (PoW) vs the Proof of Stake (PoS)

### 1.1 Rationale

On this basis, integrating Big Data with Blockchain solves several issues inherent to Big Data management, including security, data integrity, and processing speed. Big Data, in its true sense, is characterized by the ability to handle large volumes of data with high velocity and in a variety of formats. They could be more flexible and centralized, thus posing serious challenges when it comes to integration and security threats. Blockchain eliminates most of these problems by decentralizing the data, so it is not concentrated in single nodes, making the system more secure and robust. In addition, the ledger feature of Blockchain makes it possible to guarantee that once data is stored in the system, it cannot be altered again, a factor that makes it have robust data security. Extending this framework to include Machine Learning (ML) leads to improved Big Data analysis and insights by applying ML's predictive strengths. They refer to the ability of the ML algorithms to analyze large volumes of data in search of patterns and trends to support smart contracts. Automated contract execution, whereby agreements are coded on the Blockchain to execute specific actions once the analysis results are out, helps streamline the operations and minimize the need for human interjection. Such an approach offers enhanced security, effectiveness, and transparency when managing Big Data.

### 1.2 Objectives and Scope

The purpose of this study is to examine primarily whether Big Data could be combined with Blockchain and Machine Learning to solve current issues with data storage. This entails considering the existing technologies, proposing integration plans, and analyzing their benefits and risks. The areas of concern in the research include the following. Firstly, it will discuss the typical issues that arise when dealing with big data and the traditional approaches to storing and processing data,

which will help to define more concrete pain points that Blockchain and ML can solve. Secondly, the research will explore the essential background of Blockchain technology, focusing on the concept, characteristics, working, and use in strengthening data protection and accuracy. Thirdly, it will discuss different machine learning approaches and how they may be implemented with big data, particularly about implementing machine learning and blockchain technologies. The research will also outline and assess the distributed storage arrangements and smart contract platforms that enable these technologies.



**Figure 3:** How Machine Learning can be used with Blockchain Technology

This study will help in determining the application and usability and the possible advantages of the system, such as security, efficiency, and transparency, that may be encountered in the process that the system may possess disadvantages, such as complexity, performance, and scalability. This approach is rather extensive, as it is designed to give a clear idea of how these technologies can be used together to redefine the approach to Big Data.

## 2. Literature Review

It is important to note that Big Data Storage plays a crucial role in data analysis and processing. The previous methods of handling Big Data are primarily based on distributed computing platforms like Apache Hadoop and Apache Spark. In Hadoop, a distributed file system (HDFS) is used along with a MapReduce model, which enables it to work on big data in a distributed manner across a cluster of interconnected computers. It partitions the data into smaller parts and processes these parts in parallel, offering scalability and resilience to faults (White, 2012). Apache Spark, in contrast, has in - memory computing solutions that significantly boost their capabilities. It employs a DAG scheduler, a query optimizer, and a physical execution engine to process complex data efficiently (Zaharia et al., 2016). All of them are designed to accommodate various forms of data processing, batch and stream, and interactive querying, thus making them handy for Big Data processing. However, these systems do come with data security, integrity, and latency problems. Centralized data storage models are a problem since they focus on having one extensive data storage system, which

increases the risk of point failures and cyber - attacks (Grolinger et al., 2013).

Furthermore, Big Data is characterized by a vast amount, rate, and range of data that require sophisticated methods to handle and analyze data proficiently. For example, Blockchain is a modern approach that addresses such issues because it decentralizes data storage. At the same time, Machine Learning is another modern tool that might help address these issues because it boosts the data analysis function. Implementing these technologies in Big Data systems can enhance data management's overall structure, safety, and effectiveness.

## Blockchain Technology

Blockchain was initially associated only with cryptocurrencies and financial services; however, in the modern world, it has found application in diverse fields such as supply chain, healthcare, and others. Blockchain provides complete traceability in the supply chain from manufacturing to product ownership, thereby reducing counterfeit products (Kamath, 2018). For example, IBM Food Trust is a system that employs Blockchain technology to track products in the food industry, from farm to consumer, to guarantee the food's quality and origin. Blockchain in healthcare helps maintain the patient's arable and easily shareable with records other healthcare providers without compromising the patient's privacy (Angraal et al., 2017). MedRec is a Blockchain - related system that enables patients to decide who to provide or not provide with their records to promote data privacy and sharing. In the finance sector, Blockchain has the capability to minimize fraud and increase transaction speed by maintaining a transparent and secure record of the transactions (Peters & Panayi, 2016). Smart contracts in Blockchain technology enable the execution of contracts and the handling of disputes without third parties, making the processes more affordable and faster. For instance, Ethereum's platform allows for creating and implementing smart contracts, thus offering decentralized applications or apps that function independently of any central authority. However, Blockchain has drawbacks, including scalability, energy consumption, and regulatory issues. There are still some concerns regarding Blockchain, and efforts to fix these problems are still being investigated, with methods such as PoS and sharding proposed to improve the efficiency and scalability of Blockchain (Buterin, 2017). Blockchain, Big Data, and Machine Learning could establish a highly safe and effective data management and sharing system.

## Machine Learning & Smart Contracts

Machine Learning (ML) is a subfield of Artificial Intelligence that allows systems to learn from and make decisions about data. ML is typically used for classification, regression, clustering, anomaly detection, and reinforcement learning, which can be grouped into three broad categories: supervised, unsupervised, and reinforcement learning (Murphy, 2012). In supervised learning, algorithms such as Support Vector Machines (SVM) and neural networks utilize tagged training data for classification or prediction.
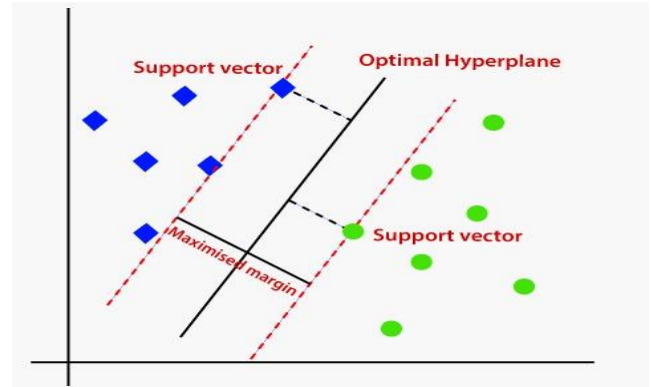


**Figure 4:** Support Vector Machine (SVM)

In contrast, the later learning category, k - means clustering, searches for patterns in data that have not been tagged. Learning, known as reinforcement learning, including Q - learning, generates ideal actions based on trial and error (Sutton & Barto, 2018). Blockchain - based smart contracts refer to self - executing contracts that contain code written directly into the contractual language to execute particular functions. They are used in many spheres of life, such as finances, supply chains, and legal systems, to settle contracts or enforce laws without mediators (Szabo, 1997). For instance, in the financial industry, smart contracts are applied to execute trades and clear them; this shortens the time and minimizes the expenses. ML augments smart contracts by making them capable of decision - making based on data analytics, thus making them more effective. For instance, ML algorithms can involve analyzing market trends and automatically initiating smart contract reactions, such as buying or modifying clauses due to current data. The combination of ML with Blockchain and smart contracts can produce constantly evolving and self - executing systems, thereby offering more excellent value and minimizing the need for human interference.

## Previous Research

Various reviews of the literature has made it clear that previous work has addressed various facets of incorporating Blockchain with other technologies while stressing that Blockchain possesses the capability to revolutionize how data is processed and protected. For example, in the case of a Blockchain - based distributed cloud storage system, Storj, the authors prove that Blockchain allows the distribution of data storage and makes the overseers less dependent on central storage servers (Wilkinson et al., 2014). In their paper about Blockchain, et al. (2016) outlined how it can be used to deal with IoT, which is essentially the process through which information can be transferred securely between connected gadgets. In the healthcare theme, Blockchain has been deemed to secure EHRs and enable sufficient sharing of information between the givers of healthcare (Azaria et al., 2016). Also, the integration of Blockchain with Machine Learning has enhanced the capability of data analysis. For instance, Federated Learning, which is a type of multiple - node node that does not doesn't retain the original data, removes privacy concerns (McMahan et al., 2017). Researchers have explained how blockchain can be used to safeguard FL processes, ensure the veracity of data, and foster trust between the involved parties (Nawari & Ravindran, 2019). However, more emerging research on leveraging

Blockchain with Big Data and ML still needs to be done. Current studies in this area aim to design composite systems that build on the power of every individual technology to combat the difficulties that face handling, protecting, and analyzing large - scale data. This integration anticipates bringing about authoritarian systems that can manage the challenges of current data environments and open the path for advanced developments.

## Big Data
### Definition and Characteristics
Big Data is defined by three main characteristics, often called the 3Vs. This means that in Big Data, there are three V's: Volume, velocity, and variety.

- In real - time, volume refers to the mega volumes of data from different sources, including social media, sensors, transactions, and more (Gartner, 2018). The scale of this data varies and can be petabytes or exabytes.
- Velocity explains the rate at which data is created and analyzed. It encompasses the rate at which data is produced, processed, and streamed, which is helpful in applications requiring real - time real - time processing (Kitchin, 2014). For instance, data analysis in the trading of stocks and other related commodities entails processing time in milliseconds.
- Variety is about the types of data that are present, including tabular data (similar to databases), semi - structured data (such as XML and JSON), and unstructured data, which are text, images, videos, and social media (Gandomi & Haider, 2015). The fact that the data can be of various formats and origins increases integration hurdles.

## Current Methods
Modern data storage and Big Data analysis approaches are based on complex distributed computing platforms and cloud solutions. One of the well - known frameworks is Apache Hadoop, which utilizes the Hadoop Distributed File System (HDFS) to distribute the data across several servers and improve fault tolerance and scalability (Vavilapalli et al., 2013). Apache Spark is gaining popularity because of in - memory computing, which is far faster than disk - based computing, as seen by Zaharia et al. (2016). Spark's batch and real - time data processing capacity is ideal for Big Data use. Besides these frameworks, NoSQL databases such as Apache Cassandra or MongoDB are capable of having a flexible schema design and horizontal scaling, which can be very important for the management of Big Data that is heterogeneous and constantly changing in nature (Chang et al., 2008). These databases handle structured, semi - structured, and unstructured big data optimized for throughput and real - time analytics.

Some popular cloud service providers are Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP), which which highly available and highly scalable Big Data solutions.

**Table 1:** AWS vs Azure vs GCP



For data warehousing, Amazon has Redshift, and for data storage, it has S3 whereas for data analysis, Google has BigQuery, which allows SQL queries to be run on a large set of data quickly (Sadalage & Fowler, 2012). These platforms are compatible with several Big Data tools, including access to machine learning and AI services, offering a complete set of analytical features. Other computing models, including edge computing, are essential for data processing nearer to the source. This is especially helpful for the Internet of Things (IoT), where real - time is crucial for data processing (Satyanarayanan, 2017). Apache Kafka is used in real - time data streaming, enabling the continuous ingestion of data and subsequent data processing in a distributed environment (Kreps et al., 2011). One of the greatest strengths of Kafka is

its capacity to process large volumes of data while maintaining low latency rates, which is critical for Big Data today. In unison, these methods advance how Big Data can be governed and analyzed, overcoming the drawbacks of centralized systems and responding to the modern requirements of the data - driven world.

## Challenges
It is necessary to underline several crucial and complex problems related to the management of centralized data storage systems. The first challenge is the SPOF challenge, where the center node can cause massive damage through data loss or service disruption (Grolinger et al., 2013). Other disadvantages include high costs because additional and more

costly infrastructure is needed to hold more significant volumes of data. Besides, the centralization of the controls results in easy vulnerability by hackers, and this Will complicate data protection. A failure to protect the clients' data puts the concerned persons or firms at a loss and other related consequences. It is impossible to solve these problems with the help of modern centralized and not very reliable solutions in data storage and processing. Separate systems also face latency issues since data has to be transferred between the central location and the concerned department, which might result in increased time consumption (Khan et al., 2014). This is especially true if the application is in real - time and in areas where time is critical. Additionally, some large companies, especially the ones based in the European Union that implement GDPR, require assistance in this matter because centralized systems require strict regulation of access and storage locations. Blockchain alleviates these problems in that data is distributed across the network, so it is faster and follows the regulations.

## Blockchain Technology

### Explanation of Blockchain

Blockchain is an open - source distributed ledger technology used to store transactions in peer - to - peer networks. Every exchange requires the formation of a block, and these blocks are linked one after another in a line to give a checkered record of all the transactions taking place (Nakamoto, 2008). Once a transaction is made on the blockchain, the record cannot be changed again, and if there is a need to change something, the entire blocks must be changed; this will require consensus of the details by the majority of nodes in the network (Yaga et al., 2018). Such a feature is considered Blockchain, and it is used in applications that may require factors like integrity and data security.

Another component that is incorporated in the blockchain involves Consensus Algorithms, including Proof of work (PoW), Proof of stake (PoS), and Byzantine fault tolerance (BFT), which is used in coordinating nodes in a distributed network (Dinh et al., 2018). In proof of work, transactions, as well as the arrangement of blocks in a block, are certified through computational work that consumes generous energy. In this case, PoS selects validators better through the staking process and consumes even less energy. The major is that with BFT, the system becomes dependable and can include faults and ill - intentioned nodes. This makes it advisably fitting for permissioned blockchains whereby nodes that are assured are in a position to be trusted.
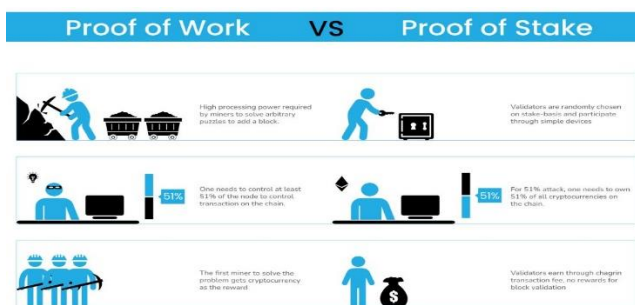


**Figure 5:** Proof of Work and Proof of Stake (PoS)

### How Blockchain Works
Blockchain is a decentralized system with every participant connected on a peer - to - peer basis, and all have a copy of the entire database. It is important to note that consensus mechanisms validate the transactions, including Proof of Work (PoW) and Proof of Stake (PoS). In PoW, nodes (miners) work to solve some cryptographic issues and create new blocks to be incorporated into the chain, guaranteeing the network's security through computational work (Dinh et al., 2018). Meanwhile, PoS chooses validators based on the number of native coins they possess and are willing to "stake" for the network. Hence, they do not require much computational power (King & Nadal, 2012). These mechanisms help to guarantee that all nodes have a copy of the same ledger and thus avoid situations when a user spends the same bitcoin twice.

### Benefits of Blockchain
Blockchain has been identified to have several significant advantages that ensure that it forms a good, stable, and reliable solution for data management. It has cryptographic and consensus properties that improve its security points to hacking and other related cyber - attacks (Zheng et al., 2017). Another benefit of the network is transparency; every transaction is visible to other participants, which increases trust and responsibility. While decentralization removes hierarchy and, with it, the 'single point of failure' that could be exploited by malicious parties (Underwood, 2016), Additionally, since Blockchain can record several transactions without having to be changed, it brings about data accuracy, which is vital in financial and supply chain industries among others. Altogether, these attributes improve the dependability, protectedness, and effectiveness of the data management systems and thereby remove many of the shortcomings of the centralized schemes.

## Machine Learning

### Overview
Machine Learning is a subset of Artificial Intelligence where the models use data to make decisions. It is applied in many fields, such as data mining, diagnosis, and prognosis.

### Role in Data Analysis
There is no doubt that today's machine learning (ML) models play crucial roles in analyzing large sets of data and detecting patterns that can help improve decision - making activities. There are various types of ML data analysis, such as supervised learning, unsupervised learning, reinforcement learning, and others, which help to analyze data ide, ntify patterns, and make forecasts. For instance, supervised learning algorithms like linear regression and neural networks are used to forecast the results from the given data, which is significant in finance and health (Murphy, 2012). Cluster analysis, which belongs to unsupervised learning, helps put together data sets into meaningful clusters that are useful in market segmentation analysis and customer identification (Aggarwal, 2015). Trial and error active reinforcement learning, in which the agent tries to propose the best mode to proceed from the state with the help of the environment, is applied for the complexity and dynamism like robots and games (Sutton & Barto, 2018). The accumulated data and enhanced data processing allow for making quick decisions,

improving the functioning of various organizational business processes, and optimizing the sphere of predictive maintenance. They can also foresee and identify flags pointing out cases of fraud or operational issues, thus increasing security and efficiency (Chandola et al., 2009). Due to this, the incorporation of ML is crucial in modern data analysis as it improves the reliability and efficiency of the predictions made in organizations.

## Integration with Blockchain

Superimposing Machine Learning (ML) on the Blockchain structure can enhance the efficiency of the data processing and predict quality. L Thus, L requires an accurate and validated training set, whereas blockchain is an unalterable, decentralized method of record keeping (Yaga et al., 2018). So, introducing the recording of the data in the Blockchain enables one to ensure that the data used for the ML training is reliable and can be amended only by specific subjects. In addition, Blockchain can implement decentralized ML further, using Federated Learning, where the models are shared between the nodes. In contrast, the data themselves are not shared and, therefore, are kept secure and private (McMahan et al., 2017). This is especially beneficial when particular regions must be protected, for instance, in the health sector or the financial field. The blockchain can also enhance the ML models' transparency as it records and stores the training and updating processes required for compliance with the regulations. Moreover, smart contracts on the Blockchain can perform tasks, such as deploying and updating the ML models, triggering the model, and sending the prediction to the specified conditions (Christidis & Devetsikiotis, 2016). This integration not only enhances the security of ML models but also optimizes the steps of attesting and managing ML applications to facilitate the development of more secure and efficient AI systems.
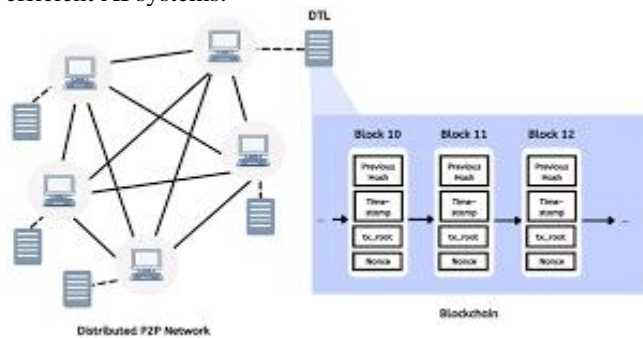


**Figure 6:** Federated Learning and Blockchain

## Smart Contracts

### Definition and Functionality

Smart contracts are digital contracts with predefined agreement terms baked into code and executable on Blockchain. These contracts offer certainty of the terms and execute the terms agreed upon without third parties facilitating them, hence fostering trust among the parties involved (Szabo, 1997). Every smart contract is stored on a decentralized ledger and thus is not reversible once executed. Smart contracts are self - executed, implying that the conditions of the contract will be met and the terms implemented without the possibility of fraud and the need for third - party interference (Buterin, 2017). This automation is attained through decentralized code processing, which

spreads the contract code in the network of Blockchain nodes to ensure effectiveness and reliability.

## Use Cases and Benefits

Smart contracts are not limited to the financial sector; they can be used in any business, such as supply chain, contracts, and more. Smart contracts can also be applied in the financial sector, where they help eliminate third parties in transactions, implying that the cost of each transaction will be considerably lowered (Peters & Panayi, 2016). For instance, trading in derivatives and securities can be done through an automated system where conditions such as price fluctuation elicit an automatic settlement. Concerning supply chain management, smart contracts provide an overhaul in transparency and accountability because each transaction and movement of products is archived in the Blockchain, thus providing stakeholders with a holistic view of the entire supply chain process (Kamath, 2018). This can significantly minimize fraud and other wastages of resources often associated with manual processes. In legal contexts, smart contracts can also facilitate the execution of contractual conditions as the system will only allow for specific actions that comply with the set conditions in the contract. Smart contracts bring many advantages such as improved productivity, decreased expenses, robust protection, and openness. Smart contracts remove the intermediaries and help improve the efficiency of the processes while simultaneously offering a reliable means of automating extensive, complicated procedures.

## Blockchain integration and applying ML

Combining smart contracts with Blockchain and ML can lead to intelligent systems adapting to real - time data and model changes. Due to their immutability and decentralization properties, smart contracts can directly implement ML - based decisions without using third parties to secure Blockchain - based transactions (Christidis & Devetsikiotis, 2016). For instance, ML models can use big data from IoT sensors to identify equipment failures in predictive maintenance. These predictions can set off other smart contracts that cause maintenance to be arranged, replacements to be purchased, or technicians to be informed, thus averting downtime and its costs (Kouicem et al., 2018). In the financial service industry, ML algorithms can consider market patterns and engage in trading through smart contracts in line with specific terms and conditions to facilitate effective and accurate transaction processing (Puschmann, 2017). This integration also improves the overall governance, as every update to the ML model or intelligent contract is logged on the Blockchain, facilitating auditability and adherence to the regulatory frameworks. Moreover, integrating the two underlying technologies, ML and smart contracts, can lead to more dynamic systems that change according to the data they receive, offering more targeted and effective services.

## Data Processing and Storage

### Ensuring Data Quality

This is particularly important in the case of historical data, which is used in most data analysis processes to draw conclusions and make predictions. The data quality management process includes data cleaning, which deals with handling duplication and missing values; data normalization, which entails ensuring that data is standardized; and data

validation, which involves checking the accuracy of data before entering it into the system. Such tools include Apache NiFi and Talend Data Integration, which are excellent for data ingestion, transformation, and validation (Oussous et al., 2018). These processes help ensure that the data fed into Machine Learning (ML) models is clean and of high quality so that the predictions and insights from the data are accurate.

## Decentralized Storage Solutions

InterPlanetary File System (IPFS) is another decentralized distributed storage platform for storing huge data off - chain and addressing scalability, while Blockchain technology guarantees data integrity. Regarding data storage, it is possible to use a distributed file system where data is split into several segments. It is located in several nodes, which minimizes the probability of its loss and increases availability (Benet, 2014). Blockchain can capture cryptographic hashes of these data fragments, providing a check - sum - like mechanism for each piece of data and an immutable audit trail. It combines the core strengths of decentralized storage and Blockchain and is therefore highly suitable for Big Data solutions.

## Decentralized Machine Learning

### Federated Learning

Federated Learning is a distributed model training method where multiple nodes are trained on different subsets of the data without passing through the raw data. This method provides high privacy and security as the data is not transferred to any central site; only the model update is transferred and combined to form a global model (McMahan et al., 2017). It is beneficial when handling sensitive information that should not be shared with the public, such as health and business. Frameworks such as TensorFlow Federated can help deploy this decentralized ML approach, enhancing data protection while providing collaborative learning.

### Training and Prediction

One of the ways of employing blockchain technology is by storing historical data and the result of the analysis in the blockchain, enhancing its security and immutability. Such data can be used to train the ML models. Thus, by utilizing Blockchain, we can guarantee the quality and origin of the training data, which is further vital to building effective and accurate predictive models (Yaga et al., 2018). The trained models can then be used to make further predictions on other data sets. For example, in the case of a predictive maintenance scenario, the data gathered on the performance of the equipment over time can be stored on the Blockchain and then fed into a model to make future failure predictions for the equipment to encourage timely maintenance actions.

## Smart Contract Automation
### Automating Actions

Smart contracts are digital contracts based on the code, meaning that once the terms of an agreement are coded, the contract automatically executes the terms. These contracts operate on Blockchain and are capable of performing actions in accordance with ML projections, including the activation of maintenance procedures in case of equipment malfunction. For example, in the manufacturing environment, IoT sensors

can gather information on the performance of machines, and the ML algorithm can then use the data to identify future equipment failures. They can then automatically generate service calls for maintenance, request spare parts, or even notify the technicians, thus reducing the time it takes to address the problems and enhance productivity (Christidis & Devetsikiotis, 2016).
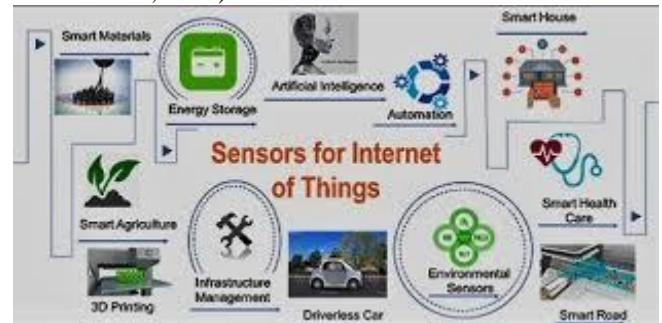


**Figure 7:** How Internet of Things (IoT) Sensors Work

## Dynamic Contracts

Smart contracts are a development of conventional contracts in that they use algorithms to execute contractual agreements; dynamic, smart contracts are a further evolution of smart contracts in that they use ML algorithms to learn. For instance, in the insurance industry, auto insurance using ML models can predict the behavior of drivers in real - time and change the insurance premium in response to driver's behavior and other factors such as weather conditions (Puschmann, 2017). These contracts can also change their terms over time using the latest data, whereas the pricing models in traditional contracts are often flexible and may need to be fairer. The combination of the utilization of ML in smart contracts makes it possible to develop responsive self - learning solutions that can make adjustments depending on the surrounding environment and offer individualized services.

## Enhancing Security and Privacy

### Privacy Preservation

Data security and privacy are paramount while maintaining the Blockchain, achieved through privacy preservation techniques. Computation on encrypted data is another innovative cryptographic technique used in secure multi - party computation and homomorphic encryption. MPC enables several participants to compute a function over their shares without revealing their data to any other participant; this is appropriate in collaborative contexts (Bogdanov et al., 2008). On the other hand, homomorphic encryption enables computations to be made on encrypted values and produce a result that, when decrypted, is the same as that of computations on actual values. They ensure that data of high - security concern is still beneficial for analytical purposes.

### Anomaly Detection

There is a possibility of using machine learning algorithms to identify anomalies in blockchain networks and prevent compromise on the system's security and integrity. Anomaly detection entails searching for patterns in the data that are different from the system's expected behavior, which can be used to identify fraudsters or hackers (Chandola et al., 2009). Real - time detection of such anomalies can be achieved by using techniques like clustering, statistical methods, and

neural networks, thereby providing early warning signals to the system to prevent risks. When ML - based anomaly detection is combined with Blockchain's unalterable record - keeping system, organizations can develop systems that protect them from various illicit and make the results trustworthy.

## Benefits of Integration

- The integration of Big Data, Blockchain, and Machine Learning (ML) offers several substantial benefits: The integration of Big Data, Blockchain, and Machine Learning (ML) offers several substantial benefits:
- Enhanced Security and Data Integrity: Blockchain guarantees the integrity of records in databases while making them resistant to alteration. On. Conversely, ML algorithms improve data security by providing means to identify irregularities and potential threats in real time.
- Improved Data Analysis and Insights: Big data analysis is one of the critical strengths of the ML models as it allows organizations to gain insights into various aspects of their business and make informed decisions.
- Increased Automation and Efficiency: Thus, the use of ML in smart contracts automatically executes processes to minimize human intervention.
- Transparency and Trust: Blockchain brings transparency and trust in the decision - making process of the ML models by maintaining a record in the distributed ledger system of blockchain.
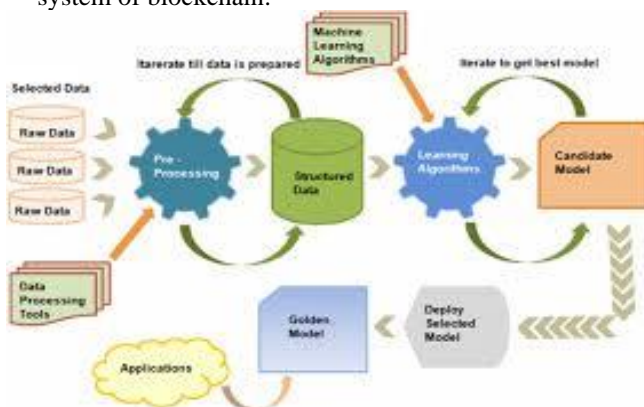


**Figure 8:** Blockchain meets machine learning

## 3. Challenges

Despite the numerous benefits, integrating Big Data, Blockchain, and Machine Learning (ML) presents several challenges: Despite the numerous benefits, integrating Big Data Blockchain, and Machine Learning (ML) presents several challenges:

- Complexity: These advanced technologies can be integrated into the system by applying data science techniques, cryptography, and Blockchain development skills. However, their considerable level of integration and coordination can be a major challenge in terms of their creation and sustenance.
- Performance: Running ML algorithms on Blockchain may cause computational overhead and latency due to training. Some consensus mechanisms, such as Proof of Work (PoW), are computationally expensive and might slow down the system's performance.
- Privacy: Although it brings the transparency feature, Blockchain also raises privacy issues since anyone in the

network can access information located in the block. Such threats are well understood, and methods such as zero - knowledge proofs and secure multi - party computation must be used.

- Scalability: The consensus mechanisms proposed by Blockchain, such as PoW, can be energy - consuming and somewhat restrictive regarding scalability. To address these problems, it is essential to work on improving the consensus mechanisms and applications outside the chain.

## Comparison with Existing Systems

Some examples of such systems include the traditional databases that store data centrally and are prone to being hacked or contain single points of failure. Centralized computing systems are characterized by centralizing all data storage and processing on a single server or a cluster of servers owned by a central authority. This architecture poses a potential risk of data loss if the central server collapses or is hacked by a cyber attacker (Grolinger et al., 2013). Furthermore, centralized systems need to be improved where they act as a bottleneck and cannot accommodate large volumes of data. However, decentralized data storage systems store data across different nodes in a network to reduce vulnerability and improve reliability. For example, the Blockchain technology used in a distributed database causes each node to store the whole data set, making it redundant and highly available (Nakamoto, 2008). Proposed structures for distributed storage, such as IPFS (InterPlanetary et al.), also add a layer of security by fragmenting data into chunks and using a P2P network to distribute them and minimize the exposure to security threats or single points of failure (Benet, 2014). This distributed approach means that data is still available and secure even if specific nodes are compromised, making decentralized systems less vulnerable than their centralized counterparts.

## Manual vs. Automated Processes

Hearths the idea of traditional systems, there are usually some interventions to perform data processing and decision making. For instance, data cleaning, conversion, and loading on analytical models often require the intervention of people to verify the correctness of the data being fed into the models and the adherence of the data to business rules. This approach can be tiresome, prone to errors, and more efficient when size increases (Sadalage & Fowler, 2012). Moreover, manual processes slow down response times to real - time data, making the decision - making process far less efficient. The combined systems apply smart contracts supported by Blockchain and Machine Learning (ML) to enhance data processing and decision - making. Smart contracts are automated digital contracts whereby the conditions of a contract are coded into a program and enforced on Blockchain by the computer system (Szabo, 1997). Such contracts can be designed to execute specific actions where the ML model has made confident predictions, such as starting a supply chain procedure where the inventory level is estimated to reach a particular low. They include purchasing, inventory control, production planning, and scheduling, among others; by automating these processes, organizations can reduce the time taken to complete those processes and be able to respond to changes in the environment more promptly. Blockchain and ML are symbiotic in that they allow for automated decision - making processes based on data collected in real time, thus

improving efficiency and minimizing the amount of human interference.

## Transparency

Automated systems in traditional systems often need to be more transparent; thus, data verification and auditing are complicated. In centralized systems, all the data is handled by a single authority, which can be a problem because the entire system is not transparent, and there is a high probability of data being manipulated (Underwood, 2016). Such practices compromise trust and confidence amongst the stakeholders since it is hard to confirm the validity and accuracy of the data. In addition, conventional databases do not offer an excellent record of accountability, which is a problem since it can be challenging to track data usage and change history. While centralized systems may be fast and efficient, blockchain systems provide transparency and accountability in data by decentralization and a distributed ledger. Every transaction made on a Blockchain is date - stamped and tied to the prior transaction, making it possible to track all the records and ascertain the sequence of transactions (Nakamoto, 2008). This brings efficiency to the network since every network member can have the data checked and trusted without central power. In the same way, the public ledger of Blockchain makes it easy for the network participant to have a clear check of the audit trail, thus increasing accountability. This transparency is beneficial in fields like supply chain management, where it is essential to monitor the origin and other activities of the products to guarantee their quality.

**Table 2:** Comparison with Existing Systems

| Traditional Systems | Integrated Systems |
|---|---|
| Centralized Data Storage: Traditional systems often rely on centralized databases, which can be vulnerable to single points of failure and data breaches. | Decentralized Data Storage: Blockchain provides a decentralized and secure way to store data, reducing the risk of data breaches. |
| Manual Processes: Many traditional systems require manual data processing and decision - making intervention. | Automated Processes: Smart contracts automate processes based on ML predictions, reducing the need for manual intervention. |
| Limited Transparency: Centralized systems may need more transparency, making verifying data integrity and provenance difficult. | Enhanced Transparency: Blockchain ensures transparency and traceability of data, enhancing trust among stakeholders |

## Case Study: Proof of Concept

### Scenario Description

Let a system employ IoT sensors to gather extensive data about industrial equipment and store them safely in a Blockchain. ML models process this information to determine the likelihood of equipment failures, and subsequent smart contracts perform maintenance tasks accordingly. This integrated approach helps coordinators reduce operational time, minimize downtime, and guarantee data accuracy and openness. As far as the parameters of equipment are concerned, the system should be able to monitor the temperature, vibrations, and pressure, among others, so that in case of any variance, it can easily pick on them, meaning that it may be easily able to predict the following action to take in case of failure. *The personal data collected are stored in a Blockchain, thus enabling an accurate and non - reversible trail.* The above information is then fed to the ML models, which utilize the patterns of failures from historical data to predict future failures efficiently. In the event of an identified concern, the ML model initiates a smart contract. This contract can automatically schedule maintenance, order required spares, and notify the maintenance workforce without the operator's intervention. Over and above the above mentioned, this proactive approach to maintenance has the following benefits: it not only helps avoid equipment breakdowns but also helps in the proper timing of the maintenance schedules, thus lowering the maintenance costs and extending the machines' useful life. Further, the use of Blockchain increases transparency and security, and the data and actions taken thereon can be trusted by all stakeholders. This system shows how the combination of IoT, Blockchain, and ML is an effective way of developing an optimal industrial maintenance system.

## Components Involved

- IoT Sensors: These devices capture real - time values of different parameters, including temperature, vibration, and pressure from industrial products. The information is then sent to a database, processed, and saved for further use.
- Blockchain: The sensor data collected is saved to guarantee that no one tampers with the data during processing. Every entry made is captured in a transaction file, creating a timeline of the equipment usage (Yaga et al., 2018).
- Machine Learning: The Blockchain contains historical data, which the ML models scrutinize to look for regularities and calculate future equipment failures. They allow preventive maintenance to minimize the chances of a machine or equipment developing a fault (Murphy, 2012).
- Smart Contracts: By relying on the ML results, smart contracts take the necessary actions, including initiating maintenance work orders, ordering new parts, or notifying technicians. This automation helps in regular and prompt maintenance to avoid more extended time loss for equipment and general interruption of operations (Christidis & Devetsikiotis, 2016).

## Pseudocode Examples
### Data Collection and Storage

```python
Copy code
import random
import time
import json
import hashlib

class IoTDevice:
 def __init__ (self, device_id):
 self. device_id = device_id

 def collect_data (self):
 data = {
 'device_id': self. device_id,
 'temperature': random. uniform (20, 100),
 'vibration': random. uniform (0, 10),
 'timestamp': time. time ()
```

```
}
return data

# Initialize IoT devices and collect data
devices = [IoTDevice (device_id=i) for i in range (10) ]
collected_data = [device. collect_data () for device in devices]

class Blockchain:
 def __init__ (self):
 self. chain = []
 def add_block (self, data):
 block = {
 'index': len (self. chain) + 1,
 'timestamp': time. time (),
 'data': data,
 'previous_hash': self. chain [ - 1] ['hash'] if self. chain else '0'
 }
 block ['hash'] = self. hash_block (block)
 self. chain. append (block)

 def hash_block (self, block):
 block_string = json. dumps (block, sort_keys=True). encode
()
 return hashlib. sha256 (block_string). hexdigest ()

# Initialize blockchain and add data
blockchain = Blockchain ()
for data in collected_data:
 blockchain. add_block (data)
```

**Data Analysis with ML**
python
Copy code

```
from sklearn. ensemble import RandomForestClassifier
import numpy as np

class MLModel:
 def __init__ (self):
 self. model = RandomForestClassifier ()

 def train (self, X, y):
 self. model. fit (X, y)

 def predict (self, X):
 return self. model. predict (X)

# Train ML model (assuming we have historical data)
ml_model = MLModel ()
historical_data = np. random. rand (100, 2) # 100 samples, 2
features
labels = np. random. randint (2, size=100) # Binary labels
ml_model. train (historical_data, labels)

# Predict using new data
collected_data = [
 {'temperature': 25.6, 'vibration': 0.3},
 {'temperature': 27.3, 'vibration': 0.1}
 # Add more data points as needed
]
new_data = np. array ([[data ['temperature'], data ['vibration']]
for data in collected_data])
predictions = ml_model. predict (new_data)
print (predictions)
```

**Smart Contract Integration**
solidity
Copy code

```
pragma solidity ^0.5.3;

contract PredictiveMaintenance {
 struct MaintenanceRecord {
 uint256 deviceId;
 bool needsMaintenance;
 uint256 timestamp;
 }
mapping (uint256 => MaintenanceRecord) public
maintenanceRecords;
 address public owner;

 event MaintenanceRequired (uint256 deviceId, uint256
timestamp);

 constructor () public {
 owner = msg. sender;
 }

 function updateMaintenanceStatus (uint256 deviceId, bool
needsMaintenance) public {
 require (msg. sender == owner, "Only owner can update
maintenance status");
 maintenanceRecords [deviceId] = MaintenanceRecord
(deviceId, needsMaintenance, now);

 if (needsMaintenance) {
 emit MaintenanceRequired (deviceId, now);
 }
 }

 function getMaintenanceStatus (uint256 deviceId) public
view returns (bool) {
 return maintenanceRecords [deviceId]. needsMaintenance;
 }
}
```

This example shows how IoT sensor data can be gathered and stored on a Blockchain, analyzed with Machine Learning models, and pointed to smart contracts that perform maintenance actions. This integrated approach brings more reliability, efficiency, and security to overall industrial processes.

**Blockchain Indexing: Paving the Way for Decentralized Big Data Storage**

Blockchain indexing plays a pivotal role in enhancing data storage and management within decentralized systems. By providing a structured approach to accessing and organizing data, blockchain indexing ensures efficient searchability, quick query responses, and improved system performance. Here's how blockchain indexing can contribute to the decentralization of big data:

1) **Enhanced Data Retrieval:** Indexing allows for faster and more efficient data retrieval by creating a systematic method to locate data within the blockchain. This is essential for managing large datasets, where traditional sequential searches would be impractical.

2) **Scalability**: Proper indexing techniques enable blockchain systems to scale effectively, handling vast amounts of data without compromising performance.

This is crucial for big data applications that require processing and storing enormous volumes of information.

3) **Improved Query Performance**: Indexing optimizes query operations, reducing the time it takes to access and retrieve specific data points. This leads to better system response times, which is vital for real - time applications and analytics.

4) **Data Integrity and Security**: By decentralizing data storage, blockchain ensures that data is not controlled by a single entity, enhancing security and reducing the risk of data breaches. Indexing supports this by maintaining a transparent and immutable record of all data transactions.

5) **Cost Efficiency**: Decentralized storage solutions, enhanced by effective indexing, can reduce costs associated with data management by eliminating the need for centralized data centers and reducing redundancy.

6) **Enhanced Transparency and Trust**: Blockchain's inherent transparency is amplified by efficient indexing, as users can easily track and verify data records. These fosters trust in the system and ensures the integrity of stored data.

7) **Interoperability**: Indexing can facilitate interoperability between different blockchain networks and traditional databases, allowing for seamless data integration and management across diverse systems.

By integrating robust indexing methods, blockchain technology can effectively manage big data in a decentralized manner, providing a secure, scalable, and efficient solution for modern data storage challenges.

## 4. Future Work

Further studies should focus on the improvement of the integrated systems that incorporate Blockchain, ML, and Big Data technologies in terms of performance and effectiveness. Another significant opportunity to enhance is the utilization of the consensus algorithm. Algorithms currently being deployed, such as the Proof of Work (PoW) and Proof of Stake (PoS), have inefficiencies in terms of energy and scalability (Dinh et al., 2018). Further investigation into the other consensus algorithms that are more advanced than the traditional Paxos and Raft consensus algorithms, like the Byzantine Fault Tolerance (BFT) or the Directed Acyclic Graphs (DAGs), could provide more efficient solutions (Castro & Liskov, 2002). Another crucial area is the creation of methods that would protect the individual's personal information. Since integrating these technologies may require processing sensitive information, particular emphasis should be placed on data protection. More specifically, techniques like secure multi - party computation (MPC), homomorphic encryption, and zero - knowledge proofs require further advancement and combination to offer good privacy assurances (Benet, 2014; Gentry, 2009).

It is also essential to discover additional applications of this integrated approach more broadly. Some possibilities are smart cities where information received from different sensors in real - time can be applied to manage different facilities and healthcare, where patients' data can be processed and analyzed for diagnostics without breaching their privacy (Zhang & Jacobsen, 2018). Predicted issues concern the issues of how to manage these combined technologies in terms of the higher level of difficulty and computing requirements they present. Blockchain, ML, and Big Data are complex methods that demand considerable computational power and knowledge in numerous fields to be combined (Hassani et al., 2018). As such, developing more user - friendly frameworks and tools that can facilitate the integration of such components will be crucial. It is also crucial to point out that the use of Blockchain imposes a certain amount of performance overhead, which has to be addressed to avoid impacting the overall performance of the systems. At the same time, ML algorithms require significant computational power to process extensive data.
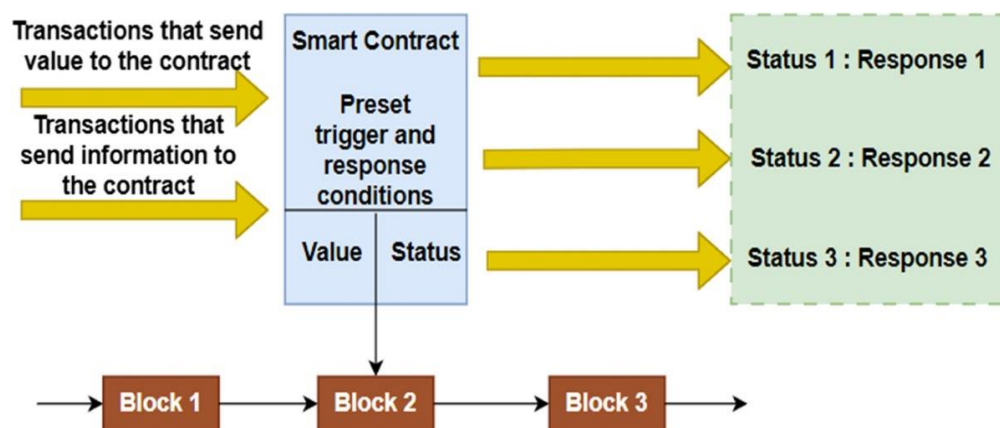


**Figure 9:** Blockchain and machine learning

## 5. Conclusion

Blockchain, Big Data, and Machine Learning can be effectively combined to provide a revolutionary approach to data management by offering increased security, transparency, and efficiency. When integrated into an organization, Blockchain removes the need for a central authority while simultaneously ensuring data reliability and confidentiality, two critical issues in centralized systems. Machine learning algorithms add to this framework by having more powerful data analysis and forecasting capacities integrated with Blockchain, which results in the capability of making decisions independently with the help of smart contracts. This integration increases productivity, saves time, and protects data and information. The proposed framework has a wide range of applicability across different domains,

including the industry 4.0 use case of predicting the need for maintenance in industrial settings, real - time data processing as part of smart cities, and secure management of patient data in the healthcare sector. However, some problems still need to be addressed, such as issues of scale, computational costs, and difficulties in integrating such sophisticated technologies. Therefore, Future studies should aim to improve the consensus algorithms, propose privacy preservation techniques, and identify new applications to fulfill the potential of this combined strategy. Integrating Blockchain, Big Data, and Machine Learning has challenges and issues. However, given its benefits, it is a direction that is likely to yield significant improvements in future data management and analysis, considering it can provide a robust, scalable, and secure solution to modern - day data issues.

## References

[1] Aggarwal, C. C. (2015). Data Mining: The Textbook. Springer.

[2] Angraal, S., Krumholz, H. M., & Schulz, W. L. (2017). Blockchain technology: Applications in health care. Circulation: Cardiovascular Quality and Outcomes, 10 (9), e003800. https: //doi. org/10.1161/CIRCOUTCOMES.117.003800

[3] Azaria, A., Ekblaw, A., Vieira, T., & Lippman, A. (2016). MedRec: Using Blockchain for Medical Data Access and Permission Management.2016 2nd International Conference on Open and Big Data (OBD), 25 - 30. https: //doi. org/10.1109/OBD.2016.11

[4] Benet, J. (2014). IPFS - Content Addressed, Versioned, P2P File System. Retrieved from

[5] Bogdanov, D., Laur, S., & Willemson, J. (2008). Sharemind: A Framework for Fast Privacy - Preserving Computations. Proceedings of the 13th European Symposium on Research in Computer Security (ESORICS), 192 - 206.

[6] Buterin, V. (2017). A Next - Generation Smart Contract and Decentralized Application Platform. Ethereum White Paper. Retrieved from

[7] Castro, M., & Liskov, B. (2002). Practical Byzantine Fault Tolerance and Proactive Recovery. ACM Transactions on Computer Systems (TOCS), 20 (4), 398 - 461.

[8] Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly Detection: A Survey. ACM Computing Surveys (CSUR), 41 (3), 1 - 58.

[9] Chang, F., Dean, J., Ghemawat, S., Hsieh, W. C., Wallach, D. A., Burrows, M.,. . . & Gruber, R. E. (2008). Bigtable: A distributed storage system for structured data. ACM Transactions on Computer Systems (TOCS), 26 (2), 1 - 26.

[10] Christidis, K., & Devetsikiotis, M. (2016). Blockchains and Smart Contracts for the Internet of Things. IEEE Access, 4, 2292 - 2303.

[11] Dinh, T. T. A., Liu, R., Zhang, M., Chen, G., Ooi, B. C., & Wang, J. (2018). Untangling Blockchain: A Data Processing View of Blockchain Systems. IEEE Transactions on Knowledge and Data Engineering, 30 (7), 1366 - 1385.

[12] Gandomi, A., & Haider, M. (2015). Beyond the hype: Big data concepts, methods, and analytics. International Journal of Information Management, 35 (2), 137 - 144.

[13] Gartner. (2018). IT Glossary: Big Data. Retrieved from https: //www.gartner. com/en/information - technology/glossary/big - data

[14] Gentry, C. (2009). Fully Homomorphic Encryption Using Ideal Lattices. Proceedings of the 41st Annual ACM Symposium on Theory of Computing, 169 - 178.

[15] Grolinger, K., Hayes, M., Higashino, W. A., L'Heureux, A., & Capretz, M. A. M. (2013). Challenges for MapReduce in Big Data.2013 IEEE World Congress on Services, 182 - 189.

[16] Grolinger, K., Hayes, M., Higashino, W. A., L'Heureux, A., & Capretz, M. A. M. (2013). Challenges for MapReduce in Big Data.2013 IEEE World Congress on Services, 182 - 189.

[17] Hassani, H., Huang, X., & Silva, E. (2018). Big - crypto: Big data, blockchain and cryptocurrency. *Big Data and Cognitive Computing*, *2* (4), 34.

[18] Kamath, R. (2018). Food Traceability on Blockchain: Walmart's Pork and Mango Pilots with IBM. The Journal of the British Blockchain Association, 1 (1), 3712.

[19] Khan, A., Othman, M., Madani, S. A., & Khan, S. U. (2014). A Survey of Mobile Cloud Computing Application Models. IEEE Communications Surveys & Tutorials, 16 (1), 393 - 413.

[20] Kitchin, R. (2014). The Data Revolution: Big Data, Open Data, Data Infrastructures and Their Consequences. SAGE Publications Ltd.

[21] Kreps, J., Narkhede, N., & Rao, J. (2011). Kafka: A Distributed Messaging System for Log Processing. Proceedings of the NetDB, 1 - 7.

[22] McMahan, B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication - Efficient Learning of Deep Networks from Decentralized Data. Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS), 30, 1273 - 1282. Retrieved from http: //proceedings. mlr. press/v54/mcmahan17a/mcmahan17a. pdf

[23] Murphy, K. P. (2012). Machine Learning: A Probabilistic Perspective. MIT Press.

[24] Nakamoto, S. (2008). Bitcoin: A Peer - to - Peer Electronic Cash System. Retrieved from

[25] Narayanan, A., Bonneau, J., Felten, E., Miller, A., & Goldfeder, S. (2016). Bitcoin and Cryptocurrency Technologies. Princeton University Press.

[26] Nawari, N. O., & Ravindran, S. (2019). Blockchain technology and BIM process: review and potential applications. *Journal of Information Technology in Construction*, *24*.

[27] Oussous, A., Benjelloun, F. Z., Ait Lahcen, A., & Belfkih, S. (2018). Big Data technologies: A survey. Journal of King Saud University - Computer and Information Sciences, 30 (4), 431 - 448.

[28] Peters, G. W., & Panayi, E. (2016). Understanding Modern Banking Ledgers through Blockchain Technologies: Future of Transaction Processing and Smart Contracts on the Internet of Money. In P. Tasca, T. Aste, L. Pelizzon, & N. Perony (Eds.), Banking Beyond Banks and Money (pp.239 - 278). Springer.

[29] Ohm Patel. (2020). Blockchain Indexing Strategies: Enhancing Performance in Software Applications. Journal of Engineering and Applied Sciences Technology.

[30] Puschmann, T. (2017). Fintech. Business & Information Systems Engineering, 59 (1), 69 - 76.

[31] Sadalage, P. J., & Fowler, M. (2012). NoSQL Distilled: A Brief Guide to the Emerging World of Polyglot Persistence. Addison - Wesley.

[32] Satyanarayanan, M. (2017). The emergence of edge computing. Computer, 50 (1), 30 - 39.

[33] Szabo, N. (1997). Formalizing and Securing Relationships on Public Networks. First Monday, 2 (9).

[34] Underwood, S. (2016). Blockchain Beyond Bitcoin. Communications of the ACM, 59 (11), 15 - 17.

[35] Vavilapalli, V. K., Murthy, A. C., Douglas, C., Agarwal, S., Konar, M., Evans, R.,. . . & Baldeschwieler, E. (2013). Apache Hadoop YARN: Yet Another Resource Negotiator. Proceedings of the 4th Annual Symposium on Cloud Computing, 1 - 16.

[36] White, T. (2012). Hadoop: The Definitive Guide. O'Reilly Media, Inc.

[37] Wilkinson, S., Boshevski, T., Brandoff, J., & Butterfield, J. (2014). Storj: A Peer - to - Peer Cloud Storage Network. Retrieved from https: //storj. io/storj. pdf

[38] Yaga, D., Mell, P., Roby, N., & Scarfone, K. (2018). Blockchain Technology Overview. National Institute of Standards and Technology.

[39] Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., & Stoica, I. (2016). Spark: Cluster Computing with Working Sets. Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud Computing. Retrieved from

[40] Zhang, P., & Jacobsen, H. A. (2018). Towards Dependable, Scalable, and Pervasive Distributed Ledgers with Blockchains. Proceedings of the VLDB Endowment, 11 (12), 1627 - 1630.

[41] Zheng, Z., Xie, S., Dai, H., Chen, X., & Wang, H. (2017). An Overview of Blockchain Technology: Architecture, Consensus, and Future Trends. 2017 IEEE International Congress on Big Data, 557 - 564.