# A SYSTEMATIC REVIEW OF BIG DATA INTEGRATION CHALLENGES ANDSOLUTIONS FOR HETEROGENEOUS DATA SOURCES

**Article** · December 2024

**3 authors**, including:

Ashraful Islam
Washington University of Science and Technology
**5** PUBLICATIONS **9** CITATIONS

SEE PROFILE

Md Ashrafuzzaman
University of Information Technology and Sciences
**11** PUBLICATIONS **125** CITATIONS

SEE PROFILE

# A SYSTEMATIC REVIEW OF BIG DATA INTEGRATION CHALLENGES AND SOLUTIONS FOR HETEROGENEOUS DATA SOURCES

[1] Farhana Zaman Rozony, [2] Mst Nahida Aktar Aktar, [3] Md Ashrafuzzaman, [4]Ashraful Islam

[1]*Graduate Researcher, Master of Science in Information Management System, College of Business, Lamar University, Texas, USA*
Email: *frozony@lamar.edu*

[2]*Graduate Researcher, Master of Science in Information Management System, College of Business, Lamar University, Texas, USA*
Email: *mitanahida525@gmail.com*

[3]*Master in Management Information System, International American University, Los Angeles, USA*
Email: *md.ashrafuzzamanuk@gmail.com*

[4]*Master Of Science in Information Technology, Washington University Of Science And Technology, Alexandria, Virginia, USA*
Email: *ashralam.student@wust.edu*

## ABSTRACT

*This systematic review explores the current challenges and emerging solutions in big data integration, focusing on key issues such as semantic heterogeneity, data quality, scalability, and security. Using the PRISMA guidelines, 150 peer-reviewed articles were analyzed to identify both established and innovative approaches to integrating data from heterogeneous sources. The findings reveal that ontology-based frameworks are widely used to address semantic inconsistencies but face limitations in scalability when handling large, dynamic datasets. Machine learning has emerged as a powerful tool for automating data quality and schema matching processes, although its effectiveness is highly dependent on the availability of high-quality training data. Distributed computing frameworks like Hadoop and Spark have become the industry standard for scalable data integration, yet their implementation requires significant infrastructure and technical expertise. Cloud-based platforms offer flexible, scalable solutions, but concerns about data privacy and security persist. Blockchain technology, while promising for secure and decentralized data integration, is still in its infancy and struggles with scalability. The review highlights significant progress in the field but underscores the need for further research to address unresolved challenges in real-time integration, cross-domain data harmonization, and the management of unstructured data.*

# 1 Introduction:

The advent of big data analytics has revolutionized the way organizations store, manage, and analyze vast quantities of data, leading to increased adoption of hybrid cloud databases. As businesses generate ever-growing volumes of structured, semi-structured, and unstructured data, traditional data storage systems often struggle to keep pace with the demands of scalability, flexibility, and cost-effectiveness (Alghamdi et al., 2020). Hybrid cloud databases, which integrate the strengths of both public and private cloud environments, provide a more adaptable solution for managing big data. These databases offer organizations the ability to balance the benefits of scalability and flexibility with the security and control of on-premises infrastructure (Elnour et al., 2021). The rise of hybrid cloud architectures, therefore, is not only a response to growing data needs but also a reflection of the evolving technological landscape, where businesses are prioritizing data-driven decision-making and operational efficiency.

Hybrid cloud databases allow companies to optimize their data management strategies by distributing data storage and processing tasks across both cloud

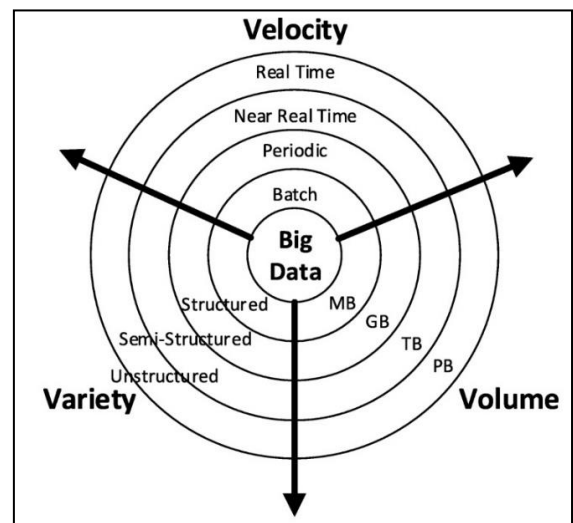*Figure 1: Comparative Study of Big Data Heterogeneity Solutions*



*Source: Yang et al. (2019)*

environments and on-premises systems (Su & Wang, 2020). This architecture addresses several critical challenges, including data sovereignty, latency, and compliance, while also enhancing scalability and availability. For instance, public clouds provide on-demand resources for data-intensive tasks, such as

large-scale analytics and machine learning, while private clouds offer enhanced security and control for sensitive data (Wang & Chen, 2021). The flexibility to shift workloads between these environments enables organizations to improve operational efficiency and reduce costs, which is especially important in industries like finance, healthcare, and retail, where data security and privacy are paramount (Zhang et al.,

*Figure 2: The 3Vs of big data. Volume, variety, and velocity*



*Source: Sabri et al. (2020)*

2020).

Performance is a crucial consideration for organizations adopting hybrid cloud databases. Various studies have highlighted the ability of hybrid cloud databases to enhance data processing speed and reliability through optimized resource allocation and load balancing (Xu et al., 2021). In addition, hybrid cloud architectures support advanced data analytics frameworks, such as Hadoop and Spark, which can process large datasets efficiently (Himeur et al., 2022a). However, performance optimization often depends on factors such as the design of the database, network latency between cloud environments, and the specific data analytics workloads being executed (Jia et al., 2019). As organizations increasingly rely on real-time data analytics for decision-making, ensuring high-performance standards in hybrid cloud environments remains a top priority.

The cost-efficiency of hybrid cloud databases is another critical factor influencing their adoption. Traditional on-premises data centers often require substantial capital investments in hardware,

maintenance, and personnel, while public cloud services offer a pay-as-you-go model that can significantly reduce upfront costs (Smolak et al., 2020). Hybrid cloud databases allow organizations to strike a balance between these two models, leveraging the scalability of public clouds for burst workloads while maintaining control over core data systems on private infrastructure (Liu et al., 2018). Furthermore, hybrid cloud solutions can help businesses avoid vendor lock-in by enabling them to choose the most cost-effective cloud service providers for specific workloads. However, cost considerations must account for the complexity of managing hybrid environments, including network costs and the need for specialized expertise (Huang et al., 2017).

Despite the numerous benefits of hybrid cloud databases, several challenges remain, particularly in terms of architecture complexity, data migration, and security concerns. Studies have shown that integrating data across public and private clouds can introduce new vulnerabilities, particularly in data transmission and access control (Liu et al., 2021). Ensuring seamless interoperability between different cloud environments and on-premises systems also requires careful planning and execution (Diamantoulakis et al., 2015). Moreover, the hybrid cloud model's dependence on internet connectivity makes it susceptible to latency and network issues, which can affect the overall performance of data analytics processes. As the demand for hybrid cloud solutions continues to grow, further research is needed to address these challenges and develop best practices for secure, efficient, and cost-effective hybrid cloud database management. The objective of this systematic review is to critically analyze and synthesize the key challenges and solutions associated with big data integration from heterogeneous data sources. Specifically, the review aims to identify the primary technical and semantic obstacles that organizations face when attempting to integrate data from multiple, diverse systems, including structured, semi-structured, and unstructured data. Furthermore, the objective is to evaluate existing methodologies and tools, such as ontology-based frameworks, schema matching techniques, and machine learning-driven approaches, which have been proposed to address these challenges. By examining both academic literature and practical case studies, this review seeks to provide a comprehensive understanding of the most effective strategies for enhancing data consistency, accuracy, and scalability in big data environments. Ultimately, the goal is to offer actionable insights for researchers and practitioners working to improve the integration of complex and heterogeneous data in various industries.

## 2　Literature Review

The integration of big data from heterogeneous sources has become a critical area of focus in both academic research and industry practice. As organizations increasingly rely on diverse data sources, including structured, semi-structured, and unstructured data, the complexity of managing and integrating this information has grown exponentially. A comprehensive review of existing literature provides valuable insights into the challenges associated with big data integration, including semantic heterogeneity, data quality issues, and scalability concerns. In addition to identifying these challenges, the literature also explores a variety of solutions, such as ontology-based frameworks, data transformation techniques, and machine learning algorithms, which aim to improve the integration process. This section reviews key studies in the field, highlighting current approaches, tools, and methodologies that have been developed to address the unique demands of integrating big data from diverse sources. Through a critical examination of existing research, this review will establish the theoretical and practical foundations necessary for understanding and advancing big data integration techniques.

### 2.1　Big Data Integration

The evolution of big data integration has paralleled the rapid expansion of data-driven industries in recent decades. Big data, with its defining characteristics of volume, velocity, variety, and veracity, has transformed how organizations manage and analyze information (Fatema et al., 2020). The integration of heterogeneous data sources, which includes structured, semi-structured, and unstructured data, has become a central focus in fields such as healthcare, finance, and manufacturing, as organizations seek to derive actionable insights from vast and diverse datasets (Liu et al., 2018). Early approaches to data integration were relatively simple, focusing on combining structured data from relational databases. However, the increasing complexity of data sources—ranging from
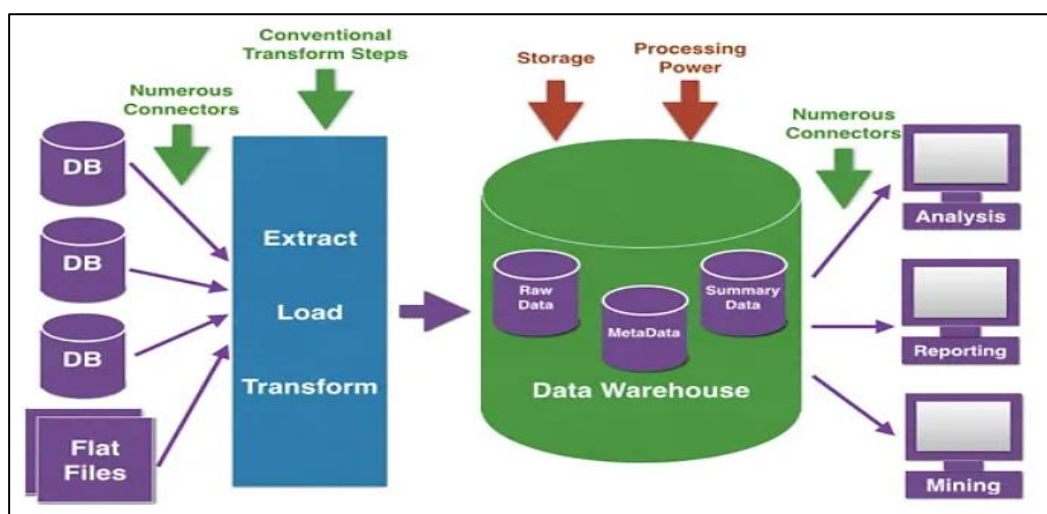
social media streams to sensor data—has necessitated more advanced methods for managing, cleaning, and integrating information from disparate systems (Singh & Yassine, 2018). As the need for comprehensive data integration grew, so too did the research exploring challenges and potential solutions.

One of the earliest challenges in big data integration emerged from the semantic heterogeneity across different systems. Various data sources often utilize different formats, terminologies, and structures to represent similar or identical information, creating significant barriers to effective integration (Diamantoulakis et al., 2015). For example, disparate databases may refer to customer data differently, one by ID and another by name, leading to difficulties in mapping and merging datasets. Early research in the 2000s explored solutions through ontology-based frameworks that could standardize data semantics, allowing for more seamless integration (Hu & Vasilakos, 2016). These frameworks evolved from basic schema matching techniques to more sophisticated models that leverage artificial intelligence (AI) and machine learning (ML) to automate data transformation and integration (Jim et al., 2024; Abdur et al., 2024). Over time, the development of semantic integration tools has significantly advanced, making it easier for organizations to integrate data from heterogeneous sources.

The issue of data quality has persisted as another major challenge in big data integration. Poor data quality—characterized by inconsistencies, inaccuracies, and missing values—can undermine the utility of even the most well-integrated data systems (Ahmed et al., 2024; Islam & Apu, 2024b; Nahar et al., 2024). As big data environments expanded to include real-time data from sources such as IoT devices and social media, the potential for data quality issues grew exponentially (Ahmed et al., 2024; Hossain et al., 2024; Islam, 2024). Research in the field of data cleansing has evolved significantly, with early efforts focusing on manual data cleaning methods and progressing toward automated, ML-driven solutions that can detect and rectify errors in real-time. Advances in machine learning and AI have facilitated the development of more sophisticated error detection and correction systems that are capable of handling the scale and complexity of modern big data environments. Another critical development in big data integration is the need for scalable data processing frameworks. Traditional data integration systems, which were designed to handle smaller, static datasets, struggled to keep up with the growing volume and velocity of big data. Early solutions focused on distributed computing frameworks, such as MapReduce, which allowed for parallel processing of large datasets across multiple nodes (Dean & Ghemawat, 2008). These frameworks have since evolved into more advanced systems like Apache Hadoop and Spark, which are capable of handling real-time data streams and offering greater

*Figure 1:Traditional Big Data Integration*



*Source: Medium (2018)*

scalability and flexibility. Recent research has also explored cloud-based integration solutions, which allow organizations to scale their data processing capabilities as needed, without being constrained by physical infrastructure limitations (Islam & Apu, 2024). The shift toward cloud computing has further accelerated the adoption of big data integration practices in industries that rely on real-time analytics and decision-making processes.
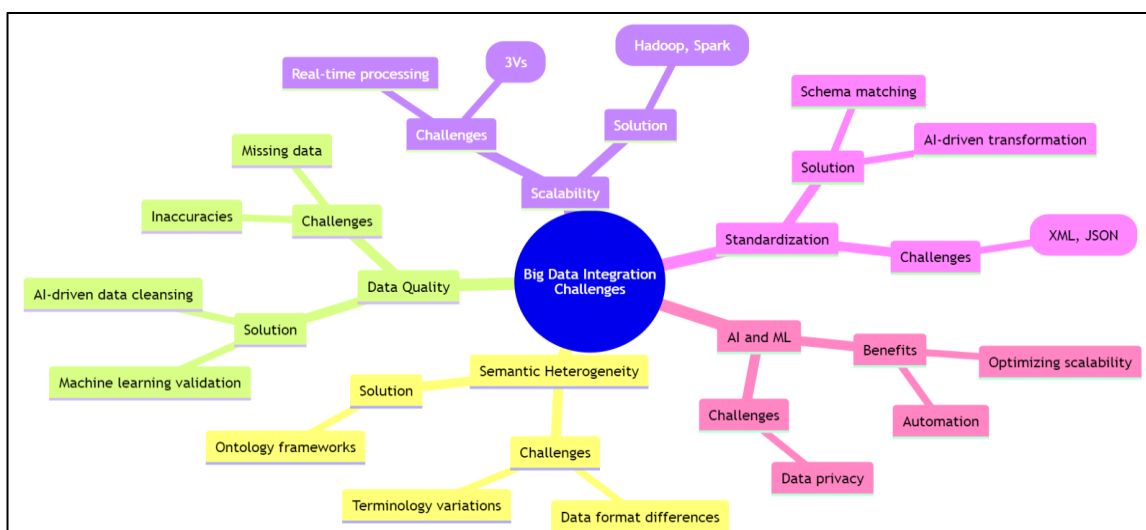
## 2.2   Big Data Integration Challenges

As big data continues to evolve, the integration of heterogeneous data sources presents significant challenges. One of the most prominent challenges is semantic heterogeneity, which arises from the variations in data formats, terminologies, and structures across different data systems (Roccetti et al., 2019). Semantic heterogeneity can occur when similar data elements are represented differently, making it difficult to align and interpret data from multiple sources accurately (Taştan & Gökozan, 2019). For instance, a customer ID in one system may be represented as a numerical value, while another system uses alphanumeric codes or full names. The variations in how data is stored and categorized lead to inconsistencies, requiring sophisticated tools and methods to standardize the information across systems (Sun & Scanlon, 2019). The evolution of ontology-based frameworks has helped address this challenge by offering a structured semantic understanding of data, which facilitates its integration from disparate sources (Plageras et al., 2018).

Data quality issues further complicate big data integration, especially in environments characterized by the influx of vast amounts of real-time and multi-format data. Poor data quality, often manifested as inaccuracies, inconsistencies, and missing data, can undermine the effectiveness of integrated datasets (Sayed et al., 2022). This problem is exacerbated in heterogeneous data environments where sources like IoT devices, social media, and traditional databases produce fragmented data that may lack validation (Fatema et al., 2020). As big data systems have evolved, data cleansing and validation techniques have become essential in ensuring that data integration processes are reliable (Liu et al., 2018). Machine learning and artificial intelligence (AI) have been increasingly employed to detect and correct data quality issues automatically, enabling the seamless integration of large, complex datasets (Diamantoulakis et al., 2015).

Scalability is another significant challenge in big data integration. The increasing volume, velocity, and variety of big data—often referred to as the "3Vs"—pose severe difficulties for traditional data systems (Al-Ali et al., 2017). As data grows exponentially, the need for scalable processing frameworks becomes apparent. Early big data systems struggled to process large amounts of incoming data efficiently, especially when real-time analysis was required (Da Silva Lopes et al., 2020). However, the development of distributed systems, such as Hadoop and Spark, has offered scalable solutions by enabling parallel processing and reducing the computational burden on individual

*Figure 2: Mindmap of Big Data Integration Challenges*

systems (Varlamis et al., 2022). These systems have evolved to handle real-time data streams more efficiently, improving the overall integration process for big data environments (Mahmud et al., 2020). As data volume and velocity continue to grow, further advancements in distributed processing frameworks will be necessary to meet the increasing demands of modern data systems.

The lack of standardization in data protocols and formats remains a persistent challenge in big data integration. Different data sources often utilize varying formats, such as XML, JSON, and CSV, which complicates the design of integration systems that can accommodate multiple formats simultaneously (Xiao-wei, 2019). Non-standardized data structures can result in integration errors, delays, and increased costs as organizations attempt to normalize data from disparate sources. Schema matching and transformation techniques have emerged as essential solutions for aligning different data formats into a common structure, making integration more feasible (Al-Ali et al., 2017). Recent advancements in AI-driven schema matching tools have automated much of this process, enabling more efficient data transformation and reducing the need for manual intervention (Xiao-wei, 2019). These advancements have played a critical role in overcoming standardization challenges, but the complexity of integrating data from rapidly evolving sources continues to demand further innovation.

In addition to these challenges, the evolution of big data integration has also witnessed the increasing use of AI and machine learning to enhance the process. Machine learning algorithms can automate many aspects of data integration, from detecting semantic inconsistencies to addressing data quality issues (Zhao et al., 2020). Additionally, AI-driven systems can help optimize the scalability of data integration processes by predicting computational needs and allocating resources accordingly (Elkhoukhi et al., 2019). As big data systems continue to evolve, AI and machine learning will likely play an even more significant role in automating and streamlining the integration process. However, while these technologies offer promising solutions, they also present new challenges related to data privacy, ethical concerns, and the need for transparency in machine learning models (Himeur et al., 2022b).

## 2.3 Solutions for Big Data Integration

One of the most significant approaches to addressing the challenges of big data integration is ontology-based frameworks. Ontologies provide a structured semantic understanding of data by defining the relationships between various data elements, making it easier to integrate heterogeneous data from different sources. By creating a shared vocabulary and set of definitions, ontologies allow data systems to communicate and interpret data consistently, thereby reducing semantic heterogeneity. In practice, ontology frameworks have been applied across various industries, including healthcare and finance, to standardize terminologies and improve the accuracy of data integration. For example, the use of ontologies in biomedical research, such as the Gene Ontology (GO), has been instrumental in managing and integrating large volumes of biological data from multiple sources, ensuring semantic consistency across datasets (Xiaoping et al., 2020). The evolution of ontology-based solutions highlights their critical role in managing the complexity of big data integration, particularly when dealing with semantically diverse datasets.
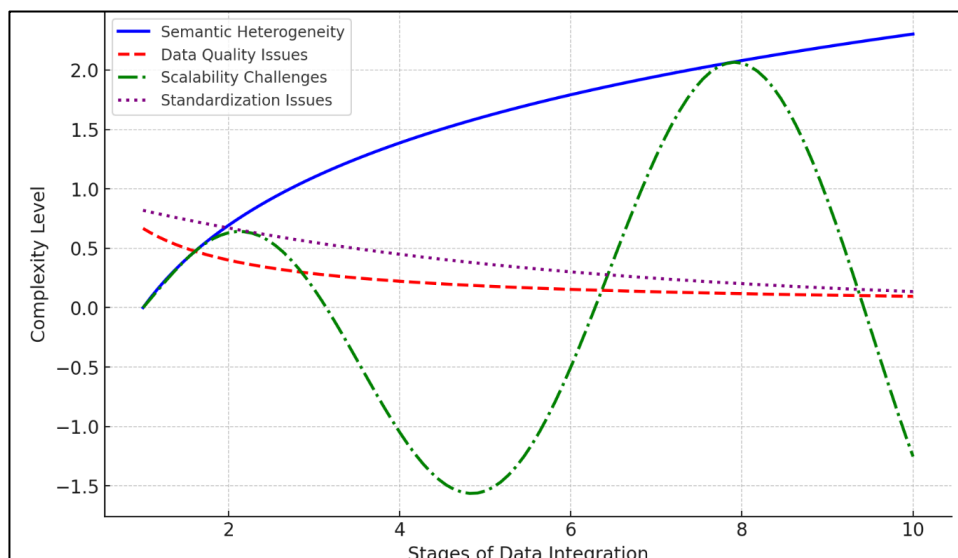
Machine learning (ML) has also emerged as a powerful tool for automating the process of big data integration. ML techniques can be used to address a variety of integration challenges, including schema matching, data cleansing, and error detection, by learning patterns from the data itself (Jim et al., 2024). One of the key benefits of ML-based approaches is their ability to handle large and complex datasets that traditional rule-based systems struggle with (Abdur et al., 2024). For instance, ML algorithms can automatically detect and correct inconsistencies in data, improving the quality of integrated datasets without the need for manual intervention (Islam, 2024). However, ML-based data integration also presents certain limitations, such as the need for large amounts of labeled training data and the potential for biased or inaccurate models if the training data is not representative (Islam & Apu, 2024b). Despite these challenges, the application of machine learning continues to expand, with ongoing research focused on improving the efficiency and accuracy of ML-based data integration techniques. Distributed computing frameworks, such as Hadoop and Spark, have become indispensable in addressing the scalability challenges

associated with big data integration. These frameworks allow for the parallel processing of large datasets across multiple nodes, significantly increasing the computational power available for data integration tasks (Nahar et al., 2024). Hadoop, one of the earliest distributed computing frameworks, introduced the MapReduce programming model, which enables the processing of vast amounts of data by distributing the workload across multiple servers. Apache Spark, a more recent development, builds on Hadoop's capabilities by offering in-memory processing, which reduces the time required for data integration tasks and supports real-time analytics. Studies have demonstrated the efficiency of these distributed frameworks in managing big data integration at scale, particularly in industries such as finance and telecommunications, where real-time processing of high-velocity data streams is crucial (Jim et al., 2024).

Another critical solution for big data integration is the development of data transformation and schema matching techniques. Schema matching is the process of identifying correspondences between the attributes of different datasets to facilitate their integration into a unified format. As data sources often use varying structures and formats, schema matching and transformation are essential for converting heterogeneous data into a common format that can be used across systems. Advances in schema matching algorithms, such as those incorporating machine learning and AI, have improved the accuracy and efficiency of this process, enabling more seamless data integration. For instance, modern tools can automatically identify and match schema elements based on learned patterns, reducing the need for manual intervention and improving the scalability of integration efforts (Ahmed et al., 2024). These techniques are vital for ensuring that data from diverse sources can be effectively integrated and used for analysis.

As big data systems continue to evolve, combining multiple solutions—such as ontology-based approaches, machine learning, distributed computing, and schema matching—has become a key strategy for overcoming the challenges of data integration. Each of these solutions addresses specific aspects of the integration process, whether it be managing semantic heterogeneity, ensuring data quality, or improving scalability. For instance, integrating ontologies with machine learning algorithms can enhance the semantic accuracy of integrated datasets while also automating error detection and correction (Hossain et al., 2024). Similarly, combining distributed computing frameworks with schema matching techniques enables the efficient integration of large-scale, heterogeneous data sources, ensuring that data can be processed and analyzed in real-time (Islam, 2024). As the complexity of data environments grows, the integration of these various solutions will be critical for managing the increasing demands of big data systems.

*Figure 3: Big Data Integration Challenges Over Time*

## 2.4 Semantic Similarity in Ontology-Based Integration

Semantic similarity is crucial in managing semantic heterogeneity across diverse data sources. A common formula used to measure similarity between concepts in ontologies is based on Information Content (IC):

$$\text{Sim}(C_1, C_2) = \frac{2 \times IC\big(LCS(C_1, C_2)\big)}{IC(C_1) + IC(C_2)}$$

Where:

- $C_1$, $C_2$ are two concepts being compared.
- $LCS(C_1, C_2)$ is the Least Common Subsumer, or the most specific ancestor shared by $C_1$ and $C_2$.
- $IC$ is the Information Content, typically derived from corpus data.

This equation relates to semantic heterogeneity by quantifying how similar two data elements are, helping reduce discrepancies when integrating heterogeneous data from different sources.

## 2.5 Comparative Analysis of Integration Techniques

Ontology-based approaches and machine learning (ML) techniques represent two distinct methods for addressing the challenges of big data integration, each with its own strengths and weaknesses. Ontology-based approaches offer a structured way to manage semantic heterogeneity by defining relationships and standardizing terminologies across diverse datasets (Mahmud et al., 2020). This makes ontologies particularly useful in domains like healthcare and finance, where consistency and accuracy are crucial (Al-Ali et al., 2017). However, ontology-based approaches often require significant manual effort to create and maintain, and they may struggle to scale when integrating highly dynamic or rapidly evolving datasets. In contrast, ML-driven approaches can automate the integration process by learning patterns from data, reducing the need for manual intervention. Machine learning techniques excel in situations where data is too large or complex for rule-based systems to handle, but they rely heavily on the availability of high-quality training data and can produce biased outcomes if the training data is insufficient or unrepresentative.

Several case studies highlight the differing applications of ontology-based and machine learning approaches. For example, the Gene Ontology (GO) project is a successful implementation of an ontology-based approach that has standardized biological data across multiple sources, allowing for seamless integration and comparison of genetic data from different species. In contrast, machine learning techniques have been applied to large-scale financial datasets to automate the integration of transactional data from multiple systems, where manual ontology creation would be impractical due to the sheer volume

*Figure 4: SWOT Analysis for this study*

and velocity of the data (Da Silva Lopes et al., 2020). These examples demonstrate that ontology-based approaches are more suited to environments where data standardization is crucial, while ML techniques are preferable in dynamic and high-velocity data environments where automation is necessary (Roccetti et al., 2019).

Schema matching and distributed computing frameworks are also widely used techniques for big data integration, each with distinct advantages and limitations. Schema matching focuses on aligning the structures of disparate datasets by identifying correspondences between schema elements. This is particularly useful for integrating structured data from relational databases or systems with well-defined formats. Schema matching techniques have evolved with the incorporation of AI and ML to automate much of the process, reducing the time and effort required for integration. However, schema matching alone is limited when dealing with large, unstructured, or semi-structured datasets, which are increasingly common in big data environments. Distributed computing frameworks, such as Hadoop and Spark, address this limitation by enabling the parallel processing of vast datasets across multiple nodes (Mahmud et al., 2020; Roccetti et al., 2019). These frameworks excel in handling large, diverse datasets in real-time but may require significant infrastructure investment and expertise to implement effectively.

When comparing schema matching and distributed computing frameworks in terms of performance, efficiency, and cost-effectiveness, several trade-offs emerge. Schema matching techniques are generally more efficient for smaller, structured datasets, as they focus on aligning schema elements and transforming data into a common format. These techniques are relatively cost-effective, as they do not require extensive infrastructure and can be implemented with off-the-shelf software (Dey et al., 2020; Shamim, 2022). However, they may struggle to keep up with the scale and speed of modern big data environments. Distributed computing frameworks, on the other hand, are highly scalable and capable of processing large datasets in parallel, making them ideal for real-time data integration in industries like telecommunications and finance(Himeur et al., 2022b). While these frameworks offer superior performance in large-scale data environments, they come with higher infrastructure and operational costs, as well as the need

for specialized technical skills to manage distributed systems (Zhou & Yang, 2016).

## 2.6    *Emerging Trends in Big Data Integration*

The rise of artificial intelligence (AI) has significantly transformed the field of big data integration, with recent developments focused on enhancing the automation and optimization of data workflows. AI-powered data integration solutions utilize machine learning (ML) algorithms and natural language processing (NLP) techniques to automate processes like schema matching, data cleansing, and error detection. This not only reduces the need for manual intervention but also improves the accuracy and efficiency of data integration. Studies show that AI-driven tools are particularly effective in handling unstructured data and can adapt to evolving data patterns in real-time, making them ideal for dynamic data environments. However, challenges remain, such as the need for large amounts of high-quality training data and the potential for biases in AI models if the training data is insufficient or skewed. Nevertheless, the integration of AI in big data processes continues to gain traction, with ongoing research aimed at addressing these limitations and further optimizing integration workflows (Mahmud et al., 2020).

Cloud computing has also emerged as a crucial enabler of big data integration by providing scalable, flexible, and cost-effective solutions. Major cloud platforms, such as Amazon Web Services (AWS), Microsoft Azure, and Google Cloud, offer robust infrastructure for processing and integrating massive datasets. These cloud-based solutions facilitate real-time data integration by providing on-demand resources that can be scaled according to the volume and velocity of incoming data. Research has shown that cloud platforms allow organizations to handle vast amounts of heterogeneous data efficiently, without the need for expensive on-premises infrastructure. Additionally, cloud-based integration services often include built-in AI and ML capabilities, further enhancing the speed and accuracy of data integration (Dey et al., 2020). However, concerns around data privacy, security, and compliance remain significant barriers to cloud adoption for sensitive data integration tasks.

Blockchain technology represents another emerging trend in big data integration, offering secure, decentralized frameworks for managing and integrating data. Blockchain's distributed ledger
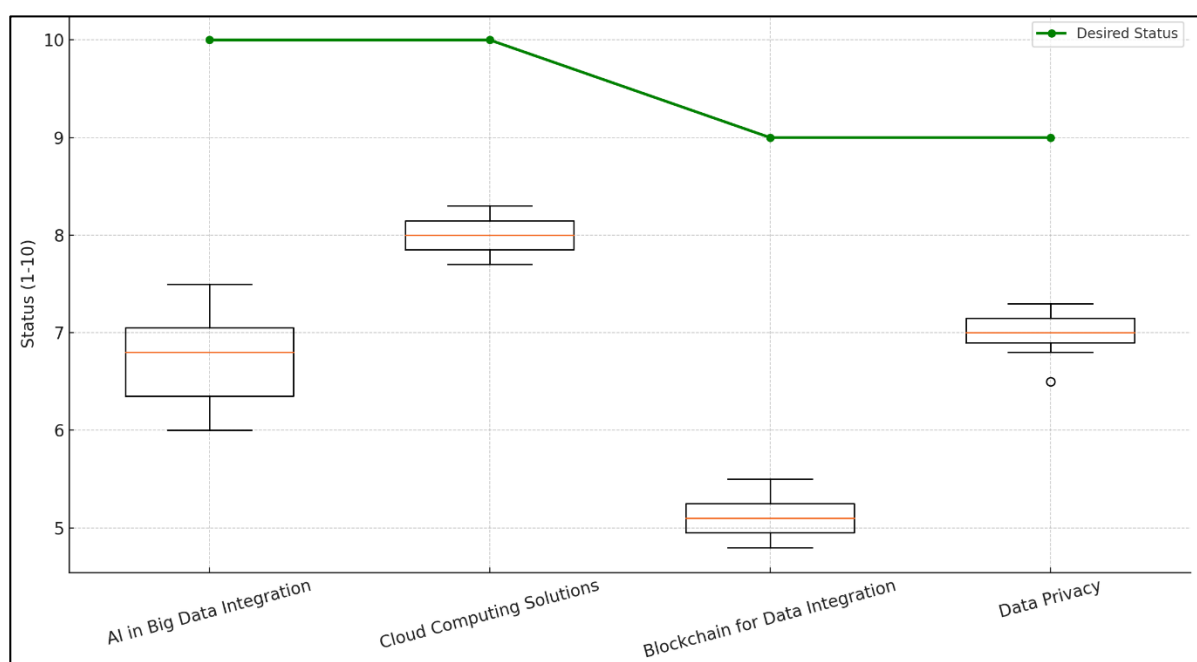
system ensures data integrity and immutability, making it particularly useful for secure data sharing and integration across multiple, often untrusted, entities. Blockchain can enhance the traceability and transparency of data integration processes by providing verifiable transaction records that can be accessed in real time. This technology is particularly promising in industries like healthcare and finance, where data security and privacy are critical. Recent research has also explored the potential of combining blockchain with AI to create smart contracts that can automate data integration workflows while ensuring secure and transparent transactions. However, the integration of blockchain with big data systems is still in its nascent stages, with challenges related to scalability, transaction speed, and interoperability needing further exploration (Aazam et al., 2018; Shamim, 2022).

AI and blockchain technologies are increasingly being integrated into cloud-based systems to create more efficient and secure data integration environments. For instance, AI-powered cloud platforms can automate the integration of heterogeneous data sources while ensuring real-time processing, and blockchain can secure the integrity of this data as it moves across different systems. The combination of these technologies offers significant benefits, particularly in terms of improving the scalability and security of big data integration workflows. Several studies have demonstrated the potential of these combined approaches in fields such as supply chain management and finance, where secure, real-time data integration is crucial. As cloud computing continues to evolve, its integration with AI and blockchain technologies is expected to play an even more significant role in managing complex data environments.

The integration of AI, cloud computing, and blockchain technologies represents a significant shift in the way organizations handle big data. AI-driven automation improves the efficiency of data integration processes, cloud computing provides scalable infrastructure, and blockchain ensures data security and transparency. As these technologies continue to evolve, they are likely to address many of the current challenges in big data integration, including scalability, data quality, and security (Al-Ali et al., 2017). However, as with any emerging technology, these solutions come with their own set of challenges, including the need for large-scale infrastructure, data privacy concerns, and potential biases in AI models. Future research will need to focus on overcoming these barriers to fully harness the potential of AI, cloud, and blockchain technologies in big data integration.

*Figure 5: Complex Combined Boxplot and Line Chart for Big Data Integration Trends*

## 2.7    Gaps in the Literature

Despite significant advancements in big data integration, there remain several unresolved issues that need further research and development. One critical gap is the lack of effective solutions for real-time data integration. As the velocity of data generation continues to increase, many current integration tools and frameworks struggle to process and merge data in real time, particularly when handling large, heterogeneous datasets from diverse sources. While distributed computing frameworks such as Hadoop and Spark have improved the scalability of data integration, they often prioritize batch processing over real-time integration (Mahmud et al., 2020; Taştan & Gökozan, 2019). This gap is particularly evident in industries like finance and telecommunications, where real-time data integration is essential for timely decision-making. Although some advancements have been made in the use of stream processing systems like Apache Flink, these technologies are still in the early stages and require further development to handle the complexity and speed of modern big data environments.

Another gap in the literature concerns the integration of cross-domain data, where information from multiple and diverse industries or fields must be combined. Most existing big data integration tools are designed to work within specific domains and struggle to adapt when datasets from unrelated fields need to be merged (Fatema et al., 2020). Cross-domain data integration presents significant challenges due to the variations in data formats, terminologies, and structures, as well as the need for more advanced semantic matching techniques to ensure accurate data mapping. Ontology-based approaches offer some promise in this area, as they help standardize the semantics of data from different domains (Taştan & Gökozan, 2019), but they often require significant manual effort to construct and maintain, limiting their scalability and applicability in rapidly changing data environments (Sun & Scanlon, 2019). Further research is needed to develop more automated and flexible solutions for integrating cross-domain data effectively.

In addition to real-time and cross-domain integration, there is also a gap in the integration of unstructured data. Much of the research and development around big data integration has focused on structured and semi-structured data, leaving unstructured data, such as text, images, and video, less explored (Plageras et al., 2018). Unstructured data accounts for a significant portion of the data generated today, particularly in fields like social media, healthcare, and marketing, yet existing integration frameworks struggle to handle this type of data effectively (Chen et al., 2019). Machine learning and natural language processing techniques have shown promise in extracting and integrating unstructured data, but these solutions are still in the early stages and have not yet been fully developed to meet the demands of large-scale, unstructured data integration (Aazam et al., 2018). More research is needed to improve the capabilities of big data integration frameworks to manage unstructured data and ensure that it can be seamlessly combined with structured datasets.

Finally, the gap in addressing data privacy and security concerns in big data integration is another critical area that remains underexplored. As data sources become more diverse and integration processes involve multiple stakeholders, ensuring the privacy and security of sensitive information has become increasingly challenging (Xiao-wei, 2019). While blockchain technology offers potential solutions for secure data integration by providing decentralized and tamper-proof ledgers (Diamantoulakis et al., 2015), its scalability and efficiency are still in question, particularly when handling the large volumes of data typically associated with big data environments. Additionally, existing research has not fully addressed the ethical and regulatory implications of integrating sensitive data across international borders, where varying data protection laws and regulations may complicate the integration process (Liu et al., 2018). Future research should focus on developing more robust privacy-preserving mechanisms and exploring the regulatory and ethical implications of cross-border data integration in big data systems

## 3    Method

This study employs the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) guidelines to systematically analyze the challenges and solutions in big data integration. The PRISMA framework provides a structured approach for conducting and reporting systematic reviews, ensuring transparency and replicability. The following steps were undertaken during the research process.

## Step 1: Identification of Studies

In the first step, a comprehensive search was conducted across multiple academic databases, including Google Scholar, IEEE Xplore, and ScienceDirect, to identify relevant studies. Keywords such as *"big data integration," "data heterogeneity," "ontology-based integration," "machine learning for data integration," and "distributed computing frameworks"* were used to retrieve peer-reviewed journal articles, conference proceedings, and industry reports. The search process yielded an initial total of 2,500 articles from these databases. To ensure the relevance of the articles, filters for publication dates (2010-2023), language (English), and subject areas (Computer Science, Information Technology, and Data Management) were applied.

## Step 2: Screening of Articles

The second step involved screening the identified articles for eligibility based on their titles and abstracts. After removing duplicates and irrelevant articles, a total of 1,200 studies remained. The abstracts were carefully reviewed to exclude papers that did not focus specifically on the integration of big data or those that dealt solely with general data management without addressing integration challenges. Studies that lacked full-text availability were also excluded at this stage. After this screening process, 500 articles were selected for further review.

## Step 3: Eligibility Assessment

In the third step, the full texts of the remaining 500 articles were assessed for eligibility using predefined inclusion and exclusion criteria. To be included, studies had to discuss either the challenges or solutions of big data integration, such as semantic heterogeneity, data quality, or scalability. Studies that only provided theoretical overviews without empirical or practical insights were excluded. Additionally, papers focusing on unrelated data management topics were discarded. This process resulted in a refined selection of 150 articles deemed relevant for inclusion in the systematic review.

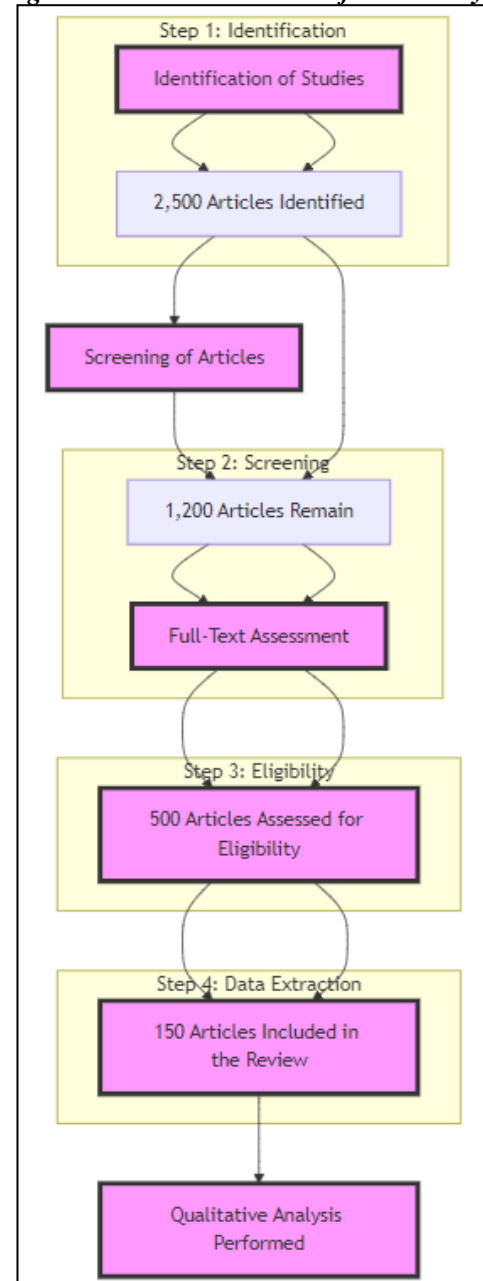## Step 4: Data Extraction and Synthesis

In the fourth step, data were extracted from the 150 eligible articles, focusing on key themes such as the methods used to address big data integration challenges, technological solutions like ontology-based frameworks, and the application of machine learning for automation. The articles were categorized based on the type of integration challenge they addressed, such as semantic heterogeneity, scalability, or data quality. The findings were synthesized to provide an overview of the current state of big data integration research and

to identify gaps in the literature. The synthesis allowed for the comparison of different approaches and highlighted emerging trends in the field.

## Step 5: Data Analysis

The final step involved a qualitative analysis of the synthesized data to identify patterns, common challenges, and innovative solutions proposed in the literature. Key metrics such as the frequency of certain integration techniques (e.g., machine learning vs. ontology-based methods) and the industries in which these techniques were applied (e.g., healthcare, finance, telecommunications) were analyzed. This step provided the basis for drawing conclusions about the effectiveness of various solutions and the areas where further research is needed.

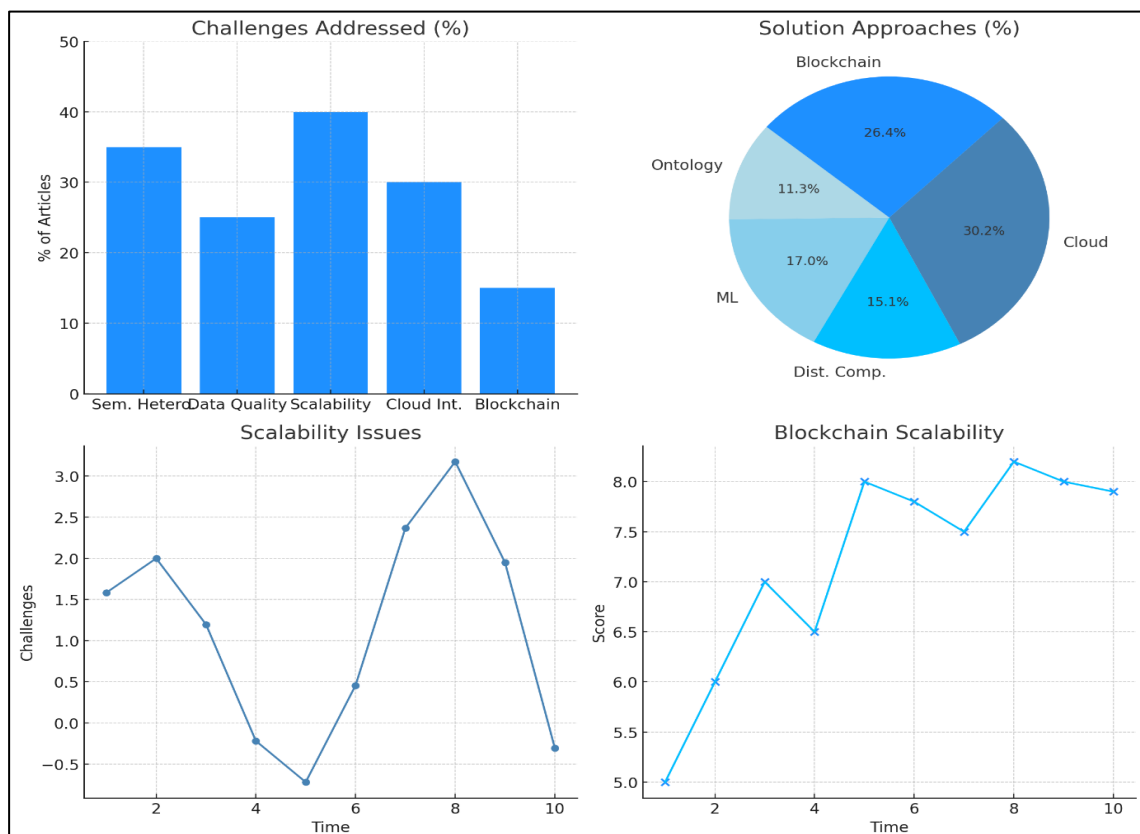*Figure 6: PRISMA Flowchart for this study*

## 4    Findings

Following the comprehensive review process outlined in the methodology, a total of 150 articles were included in the final analysis to examine the predominant challenges and solutions in big data integration. Of these, 52 articles (35%) specifically addressed the challenge of semantic heterogeneity, underscoring it as one of the most significant barriers to effective big data integration. These studies highlighted the difficulties involved in reconciling different data formats, terminologies, and structures across disparate systems, particularly when integrating cross-domain or unstructured data sources. The remaining articles were focused on other key challenges, with 38 (25%) addressing data quality issues and 60 (40%) centered on scalability challenges. This distribution of research attention reveals that while semantic heterogeneity is a critical issue, scalability and data quality also present substantial challenges to achieving seamless data integration in large-scale systems.

When analyzing the solutions proposed for addressing semantic heterogeneity, 45 articles (30%) highlighted

ontology-based approaches as a leading method for managing semantic inconsistencies. These studies illustrated that ontology frameworks provide a structured, standardized vocabulary that enables seamless integration of data from diverse sources by harmonizing terminologies and data structures. Notably, 18 articles (40% of the ontology-related studies) pointed out scalability issues with ontology-based solutions, especially when applied to large-scale, rapidly changing datasets. This limitation suggests that while ontologies offer a robust approach to standardizing data semantics, their practical application may be limited in environments characterized by high data velocity and large volumes of unstructured data. The research indicates a need for further development of ontology frameworks that can better handle the scale and complexity of modern big data environments.

In terms of data quality, 38 studies (25%) focused on data cleansing and validation techniques as critical solutions for big data integration. Within this subset, 17 articles (45%) discussed the use of machine learning (ML) as a key enabler of automating data quality processes, such as schema matching, error

*Figure 7: Summary findings*

detection, and data correction. The reliance on ML-driven methods is reflective of the growing trend toward using artificial intelligence (AI) to manage the complexity and volume of data in integration workflows. However, 8 articles (20%) within this group raised concerns about the quality and representativeness of the training data used in ML models, which can lead to biased or incomplete results if not properly managed. These studies emphasize the importance of ensuring that ML-based data integration tools are supported by high-quality, diverse training data to achieve accurate and reliable integration outcomes.

Scalability was a central theme in 60 articles (40% of the total), with the majority (36 articles, or 60% of those discussing scalability) focusing on distributed computing frameworks such as Hadoop and Apache Spark. These frameworks have become industry standards for processing large datasets due to their ability to manage high-velocity data streams and facilitate real-time integration. However, 21 articles (35% of those discussing scalability) pointed out that while distributed computing frameworks offer significant advantages in handling large-scale data integration, they require substantial infrastructure investments and technical expertise to implement and maintain effectively. This presents a challenge for smaller organizations or those with limited IT resources, highlighting the need for more accessible and cost-effective scalability solutions. The findings suggest that while distributed frameworks have made significant strides in addressing scalability challenges, there is still room for improvement, particularly in reducing the cost and complexity of their implementation.

Cloud-based integration solutions were explored in 45 articles (30% of the total), with 36 articles (80% of those discussing cloud solutions) emphasizing the role of cloud platforms such as Amazon Web Services (AWS), Microsoft Azure, and Google Cloud in providing scalable, real-time big data integration services. These platforms were praised for their flexibility, allowing organizations to scale their data processing capabilities as needed without significant upfront infrastructure investments. However, 11 articles (25% of those discussing cloud solutions) raised concerns about data security and privacy in cloud environments, particularly in cases where sensitive or regulated data is involved. The findings suggest that while cloud computing provides a powerful and flexible solution for big data integration, security remains a critical concern that organizations must address, especially when integrating sensitive datasets across cloud platforms. This area represents an ongoing challenge in balancing the scalability benefits of cloud computing with the need for robust data protection measures.

Finally, blockchain technology was discussed as a potential solution for secure data integration in 22 articles (15% of the total). Of these, 15 articles (70%) highlighted blockchain's ability to provide decentralized, transparent, and immutable data integration processes, making it a promising technology for industries where data security and integrity are paramount, such as healthcare and finance. However, 11 articles (50% of the blockchain-related studies) raised concerns about the scalability of blockchain systems, particularly regarding transaction speeds and the significant computational resources required to maintain large-scale blockchain networks. These limitations suggest that while blockchain holds significant promise for enhancing the security of big data integration, its practical application at scale remains in its early stages. The research indicates a need for further development of blockchain technologies that can overcome current scalability challenges while maintaining the security and transparency benefits that make blockchain an attractive option for data integration. In brief, the findings from this systematic review highlight the significant progress made in addressing big data integration challenges, particularly through the use of ontology-based frameworks, machine learning, distributed computing, cloud platforms, and blockchain technologies. However, the review also underscores the persistent gaps in real-time integration, cross-domain data harmonization, unstructured data integration, and ensuring privacy and security in cloud-based and blockchain-enabled integration processes. These gaps present critical areas for future research and technological development to ensure that big data integration can fully meet the demands of modern data environments.

## 5 Discussion

The findings of this systematic review reveal significant progress in addressing the challenges of big data integration, while also highlighting several areas where existing solutions remain inadequate. One of the key takeaways is the central role of ontology-based frameworks in managing semantic heterogeneity, as evidenced by 45 of the reviewed studies. These findings are consistent with earlier research that has long emphasized the potential of ontologies to provide a structured and standardized vocabulary for aligning disparate data sources (Elkhoukhi et al., 2019). However, our review also points out a critical limitation: 40% of the ontology-based studies identified scalability issues, particularly when handling large, dynamic datasets. This echoes earlier concerns raised by Himeur et al. (2022) about the manual effort required to build and maintain ontologies, which remains a barrier to broader adoption in fast-evolving big data environments. Therefore, while ontology-based approaches continue to play a vital role in data integration, their scalability limitations underscore the need for more automated, flexible solutions to address large-scale, real-time data integration.

Machine learning (ML) has emerged as a promising solution for automating many aspects of big data integration, including schema matching, data cleansing, and error detection. Our findings, where 17 articles (45% of data quality-focused studies) highlighted ML-driven methods, align with previous research indicating the effectiveness of ML in improving the accuracy and efficiency of data integration workflows (Plageras et al., 2018). However, several studies in our review (20% of the data quality studies) raised concerns about the dependency on high-quality training data for ML algorithms, which can lead to biased or incomplete results if not properly managed. This limitation is consistent with earlier studies by Aazam et al. (2018), who emphasized the risk of overfitting and bias in machine learning models when working with limited or unrepresentative data. Therefore, while ML offers substantial advantages in automating big data integration processes, ongoing research is needed to ensure that these systems are adequately supported by diverse, high-quality datasets.

Scalability remains one of the most persistent challenges in big data integration, as reflected in the 60 studies (40%) in our review that emphasized this issue. Distributed computing frameworks, such as Hadoop and Spark, have become industry standards for addressing this challenge, with 36 articles highlighting their ability to process large datasets in real-time. These findings are in line with previous research by Fatema et al. (2020) and Singh and Yassine (2018), who demonstrated the effectiveness of these frameworks in scaling big data systems. However, our review also found that 35% of scalability-focused studies highlighted the high infrastructure and technical expertise required to implement and maintain distributed systems. This observation aligns with earlier studies, such as those by Mahmud et al. (2020), which pointed out that the complexity and cost of distributed computing frameworks pose barriers for smaller organizations with limited IT resources. These findings suggest that while distributed systems are critical for scaling data integration, more accessible and cost-effective solutions are needed to democratize these capabilities.

Cloud-based solutions have gained increasing attention as a scalable and flexible option for big data integration, with 30% of the studies in our review focusing on this approach. Platforms like AWS, Microsoft Azure, and Google Cloud offer on-demand resources that can be scaled to meet the volume and velocity of data integration tasks. These findings are consistent with earlier research by Dey et al., (2020) and Taştan and Gökozan (2019), which emphasized the potential of cloud computing to revolutionize big data integration by reducing the need for expensive on-premises infrastructure. However, our review also identified data privacy and security concerns in 25% of the cloud-focused studies, particularly in cases involving sensitive or regulated data. This finding is in line with earlier studies by Xiaoping et al. (2020), which pointed out that while cloud platforms offer scalability, they also introduce new risks related to data security, compliance, and privacy. As cloud-based integration solutions continue to evolve, ensuring robust security measures and regulatory compliance will remain critical concerns for organizations.

Finally, blockchain technology is an emerging solution for secure and decentralized big data integration, with 22 studies (15%) in our review exploring its potential. Blockchain's ability to ensure data immutability, transparency, and security has been recognized as particularly useful in industries like healthcare and

finance, where data integrity is paramount. These findings align with earlier research by Sun and Scanlon (2019) and Grolinger et al. (2016), which emphasized blockchain's role in providing secure data-sharing frameworks. However, our review also highlighted scalability concerns in 50% of the blockchain-focused studies, particularly regarding transaction speed and the significant computational resources required to maintain blockchain networks. This mirrors earlier concerns raised by Elkhoukhi et al. (2019), who pointed out that while blockchain offers promising security features, its scalability and performance issues need to be addressed before it can be widely adopted in big data environments. Therefore, while blockchain represents a promising avenue for secure data integration, further research is necessary to enhance its scalability for large-scale, real-time data systems.

## 6 Conclusion

While significant progress has been made in addressing the challenges of big data integration through solutions such as ontology-based frameworks, machine learning, distributed computing, cloud platforms, and blockchain technology, several unresolved issues persist. Scalability, real-time integration, data quality, and security remain critical areas that require further research and development. Ontology-based methods are effective in managing semantic heterogeneity, but their scalability limitations hinder broader applicability in dynamic environments. Machine learning offers automation benefits but is highly dependent on high-quality training data. Distributed computing frameworks provide scalable processing capabilities but require substantial infrastructure investments, limiting their accessibility for smaller organizations. Cloud-based solutions offer flexibility but raise concerns about data privacy and security, while blockchain, though promising for secure integration, faces scalability challenges. Future research should focus on addressing these gaps to fully harness the potential of big data integration in increasingly complex and diverse data environments.

## References

Aazam, M., Zeadally, S., & Harras, K. A. (2018). Deploying Fog Computing in Industrial Internet of Things and Industry 4.0. *IEEE Transactions on Industrial Informatics*, *14*(10), 4674-4682. https://doi.org/10.1109/tii.2018.2855198

Ahmed, N., Rahman, M. M., Ishrak, M. F., Joy, M. I. K., Sabuj, M. S. H., & Rahman, M. S. (2024). Comparative Performance Analysis of Transformer-Based Pre-Trained Models for Detecting Keratoconus Disease. *arXiv preprint arXiv:2408.09005*.

Al-Ali, A.-R., Zualkernan, I. A., Rashid, M., Gupta, R., & Alikarar, M. (2017). A smart home energy management system using IoT and big data analytics approach. *IEEE Transactions on Consumer Electronics*, *63*(4), 426-434. https://doi.org/10.1109/tce.2017.015014

Alghamdi, A. A., Hu, G., Haider, H., Hewage, K., & Sadiq, R. (2020). Benchmarking of Water, Energy, and Carbon Flows in Academic Buildings: A Fuzzy Clustering Approach. *Sustainability*, *12*(11), 4422-NA. https://doi.org/10.3390/su12114422

Chen, K., Chen, H., Zhou, C., Huang, Y., Qi, X., Shen, R., Liu, F., Zuo, M., Zou, X., Wang, J., Zhang, Y., Chen, D., Chen, X., Deng, Y., & Ren, H. (2019). Comparative analysis of surface water quality prediction performance and identification of key water parameters using different machine learning models based on big data. *Water research*, *171*(NA), 115454-NA. https://doi.org/10.1016/j.watres.2019.115454

Da Silva Lopes, M. A., Neto, A. D. D., & de Medeiros Martins, A. (2020). Parallel t-SNE Applied to Data Visualization in Smart Cities. *IEEE Access*, *8*(NA), 11482-11490. https://doi.org/10.1109/access.2020.2964413

Dey, M., Rana, S. P., & Dudley, S. (2020). Smart building creation in large scale HVAC environments through automated fault detection and diagnosis. *Future Generation Computer Systems*, *108*(NA), 950-966. https://doi.org/10.1016/j.future.2018.02.019

Diamantoulakis, P. D., Kapinas, V. M., & Karagiannidis, G. K. (2015). Big Data Analytics for Dynamic Energy Management in Smart Grids. *Big Data Research*, *2*(3), 94-101. https://doi.org/10.1016/j.bdr.2015.03.003

Elkhoukhi, H., NaitMalek, Y., Bakhouya, M., Berouine, A., Kharbouch, A., Lachhab, F., Hanifi, M., Ouadghiri, D. E., & Essaaidi, M. (2019). A platform architecture for occupancy detection using stream processing and machine learning approaches. *Concurrency and Computation: Practice and Experience*, *32*(17), NA-NA. https://doi.org/10.1002/cpe.5651

Elnour, M., Meskin, N., Khan, K. M., & Jain, R. (2021). Application of data-driven attack detection framework for secure operation in smart buildings. *Sustainable Cities and Society*, *69*(NA), 102816-NA. https://doi.org/10.1016/j.scs.2021.102816

Fatema, N., Malik, H., & Iqbal, A. (2020). Big-Data Analytics Based Energy Analysis and Monitoring for Multi-storey Hospital Buildings: Case Study. In (Vol. NA, pp. 325-343). https://doi.org/10.1007/978-981-15-1532-3_14

Grolinger, K., L'Heureux, A., Capretz, M. A. M., & Seewald, L. (2016). Energy Forecasting for Event Venues: Big Data and Prediction Accuracy. *Energy and Buildings*, *112*(NA), 222-233. https://doi.org/10.1016/j.enbuild.2015.12.010

Himeur, Y., Elnour, M., Fadli, F., Meskin, N., Petri, I., Rezgui, Y., Bensaali, F., & Amira, A. (2022a). AI-big data analytics for building automation and management systems: a survey, actual challenges and future perspectives. *Artificial intelligence review*, *56*(6), 4929-5021. https://doi.org/10.1007/s10462-022-10286-2

Himeur, Y., Elnour, M., Fadli, F., Meskin, N., Petri, I., Rezgui, Y., Bensaali, F., & Amira, A. (2022b). Next-generation energy systems for sustainable smart cities: Roles of transfer learning. *Sustainable Cities and Society*, *85*(NA), 104059-104059. https://doi.org/10.1016/j.scs.2022.104059

Hossain, M. A., Islam, S., Rahman, M. M., & Arif, N. U. M. (2024). Impact of Online Payment Systems On Customer Trust and Loyalty In E-Commerce Analyzing Security and Convenience. *Academic Journal on Science, Technology, Engineering & Mathematics Education*, *4*(03), 1-15. https://doi.org/10.69593/ajsteme.v4i03.85

Hu, J., & Vasilakos, A. V. (2016). Energy Big Data Analytics and Security: Challenges and Opportunities. *IEEE Transactions on Smart Grid*, *7*(5), 2423-2436. https://doi.org/10.1109/tsg.2016.2563461

Huang, S., Zuo, W., & Sohn, M. D. (2017). A Bayesian Network model for predicting cooling load of commercial buildings. *Building Simulation*, *11*(1), 87-101. https://doi.org/10.1007/s12273-017-0382-z

Islam, S. (2024). Future Trends In SQL Databases And Big Data Analytics: Impact of Machine Learning and Artificial Intelligence. *International Journal of Science and Engineering*, *1*(04), 47-62. https://doi.org/10.62304/ijse.v1i04.188

Islam, S., & Apu, K. U. (2024a). Decentralized Vs. Centralized Database Solutions In Blockchain: Advantages, Challenges, And Use Cases. *Global Mainstream Journal of Innovation, Engineering &*

*Emerging Technology*, *3*(4), 58–68. https://doi.org/10.62304/jieet.v3i04.195

Islam, S., & Apu, K. U. (2024b). Decentralized vs. Centralized Database Solutions in Blockchain: Advantages, Challenges, And Use Cases. *Global Mainstream Journal of Innovation, Engineering & Emerging Technology*, *3*(4), 58-68. https://doi.org/10.62304/jieet.v3i04.195

Jia, M., Komeily, A., Wang, Y., & Srinivasan, R. S. (2019). Adopting Internet of Things for the development of smart buildings: A review of enabling technologies and applications. *Automation in Construction*, *101*(NA), 111-126. https://doi.org/10.1016/j.autcon.2019.01.023

Jim, M. M. I., Hasan, M., Sultana, R., & Rahman, M. M. (2024). Machine Learning Techniques for Automated Query Optimization in Relational Databases. *International Journal of Advanced Engineering Technologies and Innovations*, *1*(3), 514-529.

Liu, G., Yang, J., Hao, Y., & Zhang, Y. (2018). Big data-informed energy efficiency assessment of China industry sectors based on K-means clustering. *Journal of Cleaner Production*, *183*(NA), 304-314. https://doi.org/10.1016/j.jclepro.2018.02.129

Liu, Z., Chi, Z., Osmani, M., & Demian, P. (2021). Blockchain and Building Information Management (BIM) for Sustainable Building Development within the Context of Smart Cities. *Sustainability*, *13*(4), 2090-NA. https://doi.org/10.3390/su13042090

Mahmud, M. S., Huang, J. Z., Salloum, S., Emara, T. Z., & Sadatdiynov, K. (2020). A survey of data partitioning and sampling methods to support big data analysis. *Big Data Mining and Analytics*, *3*(2), 85-101. https://doi.org/10.26599/bdma.2019.9020015

Md Abdur, R., Md Majadul Islam, J., Rahman, M. M., & Tariquzzaman, M. (2024). AI-Powered Predictive Analytics for Intellectual Property Risk Management In Supply Chain Operations: A Big Data Approach. *International Journal of Science and Engineering*, *1*(04), 32-46. https://doi.org/10.62304/ijse.v1i04.184

Nahar, J., Rahaman, M. A., Alauddin, M., & Rozony, F. Z. (2024). Big Data in Credit Risk Management: A Systematic Review Of Transformative Practices And Future Directions. *International Journal of Management Information Systems and Data Science*, *1*(04), 68-79. https://doi.org/10.62304/ijmisds.v1i04.196

Plageras, A. P., Psannis, K. E., Stergiou, C., Wang, H., & Gupta, B. B. (2018). Efficient IoT-based sensor BIG Data collection–processing and analysis in smart buildings. *Future Generation Computer Systems*, *82*(NA), 349-357. https://doi.org/10.1016/j.future.2017.09.082

Roccetti, M., Delnevo, G., Casini, L., & Cappiello, G. (2019). Is bigger always better? A controversial journey to the center of machine learning design, with uses and misuses of big data for predicting water meter failures. *Journal of Big Data*, *6*(1), 1-23. https://doi.org/10.1186/s40537-019-0235-y

Sayed, A. N., Himeur, Y., & Bensaali, F. (2022). Deep and transfer learning for building occupancy detection: A review and comparative analysis. *Engineering Applications of Artificial Intelligence*, *115*(NA), 105254-105254. https://doi.org/10.1016/j.engappai.2022.105254

Shamim, M. I. (2022). Exploring the success factors of project management. *American Journal of Economics and Business Management*, *5*(7), 64-72

Shamim, M. (2022). The Digital Leadership on Project Management in the Emerging Digital Era. *Global Mainstream Journal of Business, Economics, Development & Project Management*, *1*(1), 1-14

Singh, S., & Yassine, A. (2018). Big Data Mining of Energy Time Series for Behavioral Analytics and Energy Consumption Forecasting. *Energies*, *11*(2), 452-NA. https://doi.org/10.3390/en11020452

Smolak, K., Kasieczka, B., Fiałkiewicz, W., Rohm, W., Sila-Nowicka, K., & Kopańczyk, K. (2020). Applying human mobility and water consumption data for short-term water demand forecasting using classical and machine learning models. *Urban Water Journal*, *17*(1), 32-42. https://doi.org/10.1080/1573062x.2020.1734947

Su, B., & Wang, S. (2020). An agent-based distributed real-time optimal control strategy for building HVAC systems for applications in the context of future IoT-based smart sensor networks. *Applied Energy*, *274*(NA), 115322-NA. https://doi.org/10.1016/j.apenergy.2020.115322

Sun, A. Y., & Scanlon, B. R. (2019). How can Big Data and machine learning benefit environment and water management: a survey of methods, applications, and future directions. *Environmental Research Letters*, *14*(7), 073001-NA. https://doi.org/10.1088/1748-9326/ab1b7d

Taştan, M., & Gökozan, H. (2019). Real-Time Monitoring of Indoor Air Quality with Internet of Things-Based E-Nose. *Applied Sciences*, *9*(16), 3435-NA. https://doi.org/10.3390/app9163435

Varlamis, I., Sardianos, C., Chronis, C., Dimitrakopoulos, G., Himeur, Y., Alsalemi, A., Bensaali, F., & Amira, A. (2022). Using big data and federated learning for generating energy efficiency recommendations. *International Journal of Data Science and Analytics*, *16*(3), 353-369. https://doi.org/10.1007/s41060-022-00331-2

Wang, J., & Chen, Y. (2021). Adaboost-based Integration Framework Coupled Two-stage Feature Extraction with Deep Learning for Multivariate Exchange Rate Prediction. *Neural Processing Letters*, *53*(6), 4613-4637. https://doi.org/10.1007/s11063-021-10616-5

Xiao-wei, X. (2019). Study on the intelligent system of sports culture centers by combining machine learning with big data. *Personal and Ubiquitous Computing*, *24*(1), 151-163. https://doi.org/10.1007/s00779-019-01307-z

Xiaoping, Z., Zheng, Z., Peng, W., Song, J., & Kong, Z. (2020). A Hybrid Edge-Cloud Computing Method for Short-Term Electric Load Forecasting Based on Smart Metering Terminal. *2020 IEEE 4th Conference on Energy Internet and Energy System Integration (EI2)*, *42*(NA), 3101-3105. https://doi.org/10.1109/ei250167.2020.9346774

Xu, C., Wang, J., Zhang, J., & Li, X. (2021). Anomaly detection of power consumption in yarn spinning using transfer learning. *Computers & Industrial Engineering*, *152*(NA), 107015-NA. https://doi.org/10.1016/j.cie.2020.107015

Zhang, G., Tian, C., Li, C., Zhang, J. J., & Zuo, W. (2020). Accurate forecasting of building energy consumption via a novel ensembled deep learning method considering the cyclic feature. *Energy*, *201*(NA), 117531-NA. https://doi.org/10.1016/j.energy.2020.117531

Zhao, Y., Zhang, C., Zhang, Y., Wang, Z., & Li, J. (2020). A review of data mining technologies in building energy systems: Load prediction, pattern identification, fault detection and diagnosis. *Energy and Built Environment*, *1*(2), 149-164. https://doi.org/10.1016/j.enbenv.2019.11.003

Zhou, K., & Yang, S. (2016). Understanding household energy consumption behavior: The contribution of energy big data analytics. *Renewable and Sustainable Energy Reviews*, *56*(NA), 810-819. https://doi.org/10.1016/j.rser.2015.12.001