

## REAL-TIME ANALYTICS IN STREAMING BIG DATA: TECHNIQUES AND APPLICATIONS

Md Ashraful Alam<sup>1</sup>

<sup>1</sup>Department of Computer Science, Colorado State University, Colorado, USA

Email: [mdashraful.alam@colostate.edu](mailto:mdashraful.alam@colostate.edu)

<https://orcid.org/0009-0006-0493-1031>

Ashrafur Rahman Nabil<sup>2</sup>

<sup>2</sup>MS in Information technology management, St. francis College, Brooklyn, New York, USA

Email: [anabil@sfc.edu](mailto:anabil@sfc.edu)

<https://orcid.org/0009-0005-1540-1266>

Abdul Awal Mintoo<sup>3</sup>

<sup>3</sup>Graduate student, School of Computer and Information Sciences, Washington University of Science and Technology ( WUST), USA

Email: [mintoo.hr@gmail.com](mailto:mintoo.hr@gmail.com)

<https://orcid.org/0009-0009-0493-965X>

Ashraful Islam<sup>4</sup>

<sup>4</sup>Master Of Science In Information Technology , Washington University Of Science And Technology, Alexandria, Virginia, USA

Email: [ashralam.student@wust.edu](mailto:ashralam.student@wust.edu)

<https://orcid.org/0009-0001-8067-7331>

### Keywords

Real-Time Analytics  
Streaming Big Data  
Systematic Literature Review  
Stream Processing Frameworks  
PRISMA Methodology

### Article Information

**Received:** 07, October, 2024

**Accepted:** 29, November, 2024

**Published:** 30, November, 2024

**Doi:** 10.70008/jeser.v1i01.56

### ABSTRACT

The increasing prevalence of streaming big data has revolutionized how organizations approach real-time analytics, providing a competitive edge by enabling immediate, actionable insights from continuously generated data streams. This review systematically examines real-time analytics techniques and their applications in streaming big data using the PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) framework. The methodology involves identifying, screening, and synthesizing relevant studies to provide a comprehensive overview of state-of-the-art techniques, including data preprocessing, stream processing engines, distributed computing architectures, and machine learning algorithms specifically designed for high-velocity data streams. The review further categorizes and evaluates the applications of these techniques across key industries such as healthcare, financial services, e-commerce, and intelligent transportation systems. The results underscore the critical role of stream processing engines like Apache Kafka, Apache Flink, and Spark Streaming in managing data velocity and volume, while highlighting the growing importance of machine learning models in extracting real-time insights. Challenges such as scalability, fault tolerance, and latency issues are discussed, along with emerging solutions like edge computing and federated learning. The findings contribute to the evolving field of streaming big data analytics by providing insights into best practices and identifying research gaps for future exploration.

### 1 Introduction

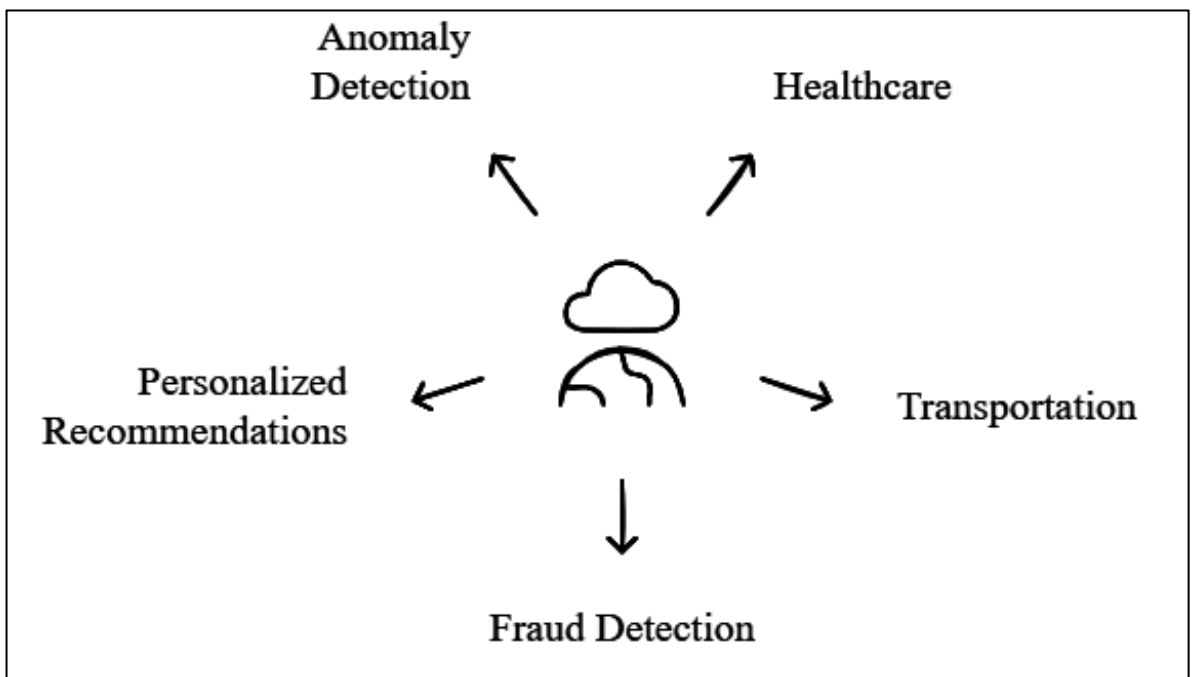
The rise of streaming big data has revolutionized the domain of real-time analytics, enabling organizations to process and analyze data as it is generated (Sun et al., 2020). Unlike traditional batch processing, real-time analytics leverages continuous data streams to provide instantaneous insights, a necessity in today's fast-paced decision-making environments (Toulis & Airoidi, 2017). Streaming big data is characterized by its velocity, volume, and variety, necessitating advanced computational frameworks and algorithms to manage and analyze data effectively (Liqing et al., 2020). This paradigm shift has fueled innovation across industries, from healthcare to transportation, driving demand for robust real-time analytics solutions. Such solutions not only enhance operational efficiency but also support predictive and prescriptive analytics, which are critical in deriving actionable insights (Kim et al., 2020).

The foundation of real-time analytics lies in stream processing frameworks, which are designed to manage the high velocity and continuous flow of data. Tools such as Apache Kafka, Apache Flink, and Spark Streaming have become the backbone of modern stream processing, enabling seamless data ingestion, transformation, and analysis (Peddireddy & Banga, 2023). These frameworks have been widely adopted due to their scalability and ability to handle fault tolerance,

key challenges associated with streaming big data (Peddireddy, 2023). Stream processing not only facilitates real-time monitoring but also supports advanced use cases like fraud detection, personalized recommendations, and anomaly detection, demonstrating its versatility across various domains. Another critical aspect of real-time analytics is the integration of machine learning algorithms that can operate on streaming data. Unlike traditional machine learning models, these algorithms are optimized for high-velocity data, enabling real-time predictions and decision-making (Wu et al., 2020). Recent advancements in online learning techniques have further enhanced the efficiency and applicability of these models, particularly in fields like healthcare, where timely interventions are crucial (Leang et al., 2019). For instance, real-time analytics has been instrumental in monitoring patient vitals and predicting potential health risks, thereby saving lives. Despite these advances, challenges such as model updating and deployment in distributed environments remain significant research areas (Peddireddy & Banga, 2023).

Real-time analytics has also found extensive applications in industries such as financial services and e-commerce (Kušić et al., 2023). Financial institutions leverage real-time analytics for fraud detection, risk assessment, and algorithmic trading, where milliseconds can make a difference in outcomes. Similarly, e-

*Figure 2: Early Childhood Nutrition and Educational Outcomes Framework*



commerce platforms use streaming big data analytics to personalize user experiences, manage dynamic pricing, and optimize inventory (Kim et al., 2020). These applications underline the transformative impact of real-time analytics in driving efficiency and enhancing customer satisfaction. However, they also highlight persistent challenges, such as ensuring data privacy and managing the trade-off between speed and accuracy (Xu et al., 2021). Despite its transformative potential, real-time analytics in streaming big data is not without its challenges. Scalability, fault tolerance, and latency are critical hurdles that need to be addressed to ensure seamless operation in large-scale deployments (Wu et al., 2020). Emerging technologies such as edge computing and federated learning offer promising solutions, allowing for decentralized data processing and improved privacy (Hussen et al., 2023). These advancements underscore the need for continued innovation and research in this field to overcome existing limitations and unlock new possibilities. The following sections of this review systematically explore the techniques and applications of real-time analytics, highlighting current trends, best practices, and research gaps. The primary objective of this review is to provide a comprehensive understanding of real-time analytics techniques and their applications in streaming big data. By synthesizing insights from existing research, the study aims to identify, categorize, and evaluate state-of-

the-art techniques, including data preprocessing, stream processing frameworks, distributed computing architectures, and machine learning algorithms tailored for high-velocity data streams. Furthermore, this review seeks to explore the practical implementation of these techniques across key industries, such as healthcare, financial services, e-commerce, and intelligent transportation systems, to highlight their transformative potential. Another critical aim is to examine the challenges associated with real-time analytics, such as scalability, fault tolerance, and latency, while discussing emerging solutions like edge computing and federated learning. By addressing these objectives, the review contributes to the academic and practical discourse on streaming big data analytics, offering a structured roadmap for future research and development.

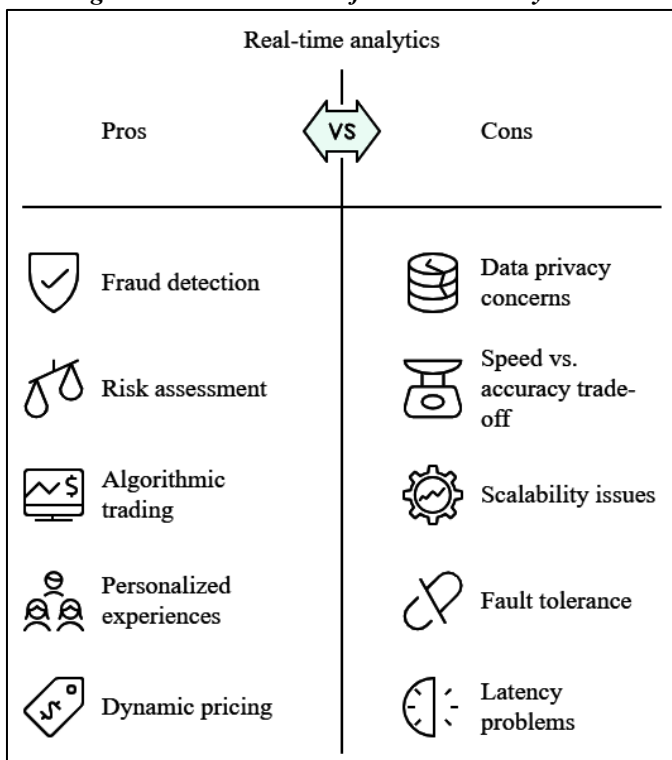
## 2 LITERATURE REVIEW

The rapid growth of streaming big data has necessitated the development of advanced analytics techniques capable of processing and deriving insights from high-velocity, high-volume data streams. This literature review aims to explore the existing body of knowledge on real-time analytics techniques, frameworks, and their applications across various industries. Using a systematic approach, this section delves into key themes, such as stream processing frameworks, machine learning algorithms for streaming data, and emerging technologies like edge computing. It also examines critical challenges, including scalability, fault tolerance, and latency, while highlighting potential solutions. By synthesizing insights from diverse studies, the review provides a structured understanding of the state-of-the-art advancements and identifies gaps for future research.

### 2.1 Foundations of Streaming Big Data

Streaming big data refers to the continuous generation, processing, and analysis of data streams in real time. Unlike traditional batch processing, streaming big data is characterized by its velocity, variety, and volume, often referred to as the "3Vs" of big data (Luo et al., 2022). This unique combination of attributes creates challenges related to data ingestion, storage, and processing, which require scalable and fault-tolerant systems (Kwon et al., 2012). Streaming big data systems must handle events or transactions as they occur, ensuring minimal latency and enabling time-sensitive decision-making (Jia et al., 2023). For instance,

Figure 3: Pros and Cons of Real-time Analysis

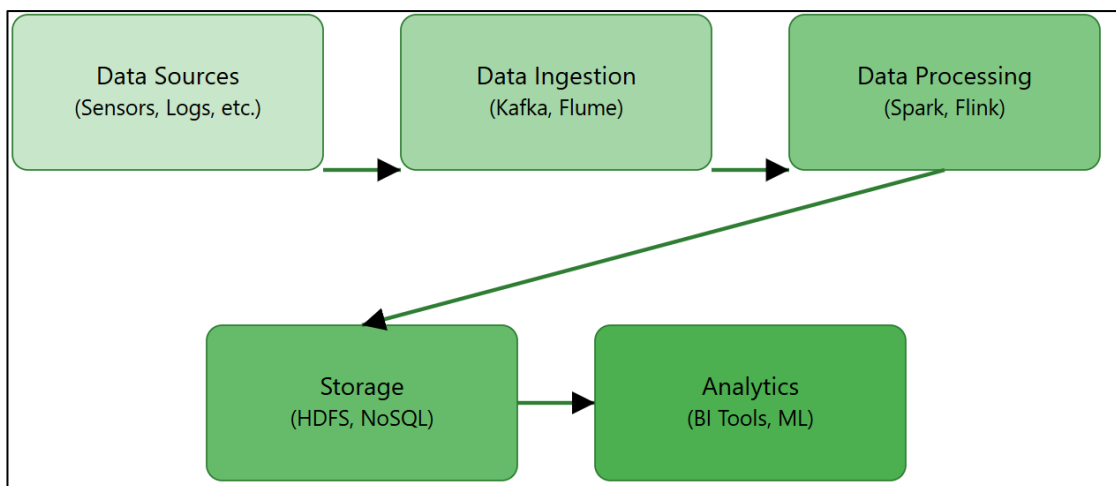


healthcare systems relying on real-time patient monitoring depend on such systems to process data at unprecedented speeds and scales, as highlighted by (Shaikh & Gupta, 2022). Moreover, the dynamic and unpredictable nature of streaming data introduces additional complexity in data quality and accuracy, further emphasizing the need for sophisticated real-time analytics frameworks (Luo et al., 2022). The evolution of real-time analytics represents a paradigm shift from traditional batch processing to continuous data stream processing. Early analytics systems relied heavily on batch processing models, where data was collected, stored, and analyzed periodically, often leading to delayed insights (Khoshkhah et al., 2022). The demand for instant insights in dynamic environments spurred the development of stream processing frameworks like Apache Kafka, Apache Flink, and Spark Streaming, which enable real-time data processing and event-driven analytics (Kastner et al., 2014). These frameworks employ distributed architectures to manage the high velocity and volume of incoming data streams effectively, providing scalability and fault tolerance (Åsberg et al., 2012). The shift towards real-time analytics has been particularly transformative in industries such as financial services, where milliseconds can be the difference between profit and loss in algorithmic trading (Kwon et al., 2012).

Real-time analytics has emerged as a cornerstone of modern industries, driving innovation and improving operational efficiency. In the healthcare sector, for instance, real-time analytics enables early detection of health anomalies through continuous monitoring of

patient data, reducing the risks of medical emergencies (Stankovic et al., 1999). Similarly, the financial industry leverages streaming analytics for fraud detection, offering immediate responses to suspicious activities (Babcock et al., 2004). In e-commerce, companies like Amazon and Alibaba utilize real-time analytics to personalize user experiences, optimize inventory management, and predict market trends (Ma et al., 2009). The integration of real-time analytics in these industries demonstrates its ability to transform decision-making processes, enhance customer satisfaction, and increase profitability, albeit with challenges such as ensuring data privacy and handling computational complexities (Kulkarni et al., 2015). Despite its transformative potential, real-time analytics in streaming big data is not without challenges. Scalability remains a significant hurdle, particularly in large-scale deployments where the sheer volume of data can overwhelm processing systems (Nair et al., 2017). Latency is another critical issue, as delayed data processing can render insights obsolete in time-sensitive applications (Khoshkhah et al., 2022). Ensuring fault tolerance in distributed systems further complicates the implementation of real-time analytics frameworks, especially in environments with high data variability (Anderson & Devi, 2006). Emerging solutions such as edge computing and federated learning offer promising avenues for addressing these challenges by decentralizing data processing and enhancing computational efficiency. These advancements underline the importance of ongoing research in real-

*Figure 4: Streaming Big Data Pipeline*



time analytics, aiming to optimize its performance and expand its applications across various domains.

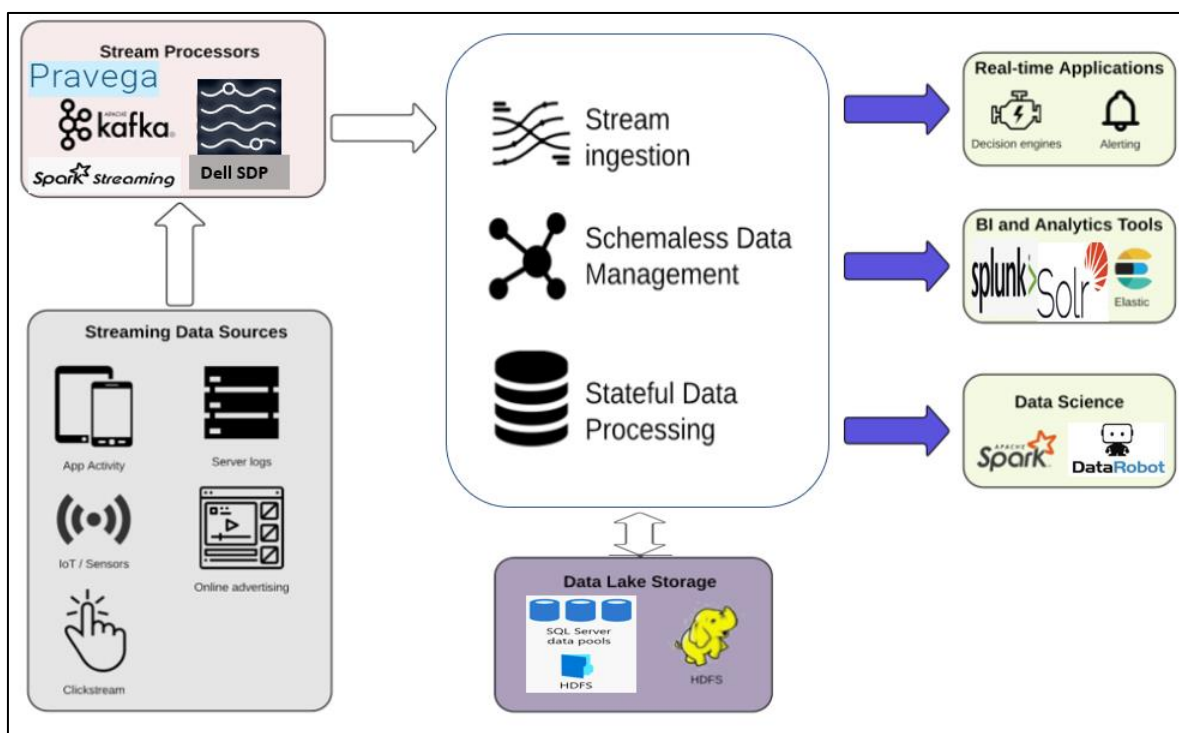
### 2.2 Stream Processing Frameworks and Architectures

Stream processing frameworks are the backbone of real-time analytics, enabling the ingestion, processing, and analysis of streaming big data at scale. Apache Kafka, Apache Flink, and Spark Streaming have emerged as leading frameworks due to their robust architectures and capabilities (Hussen et al., 2023). Apache Kafka, for instance, excels in distributed event streaming, providing a high-throughput platform for real-time data pipelines and analytics (Wei et al., 2007). It enables message brokering, allowing data to be seamlessly streamed across systems while ensuring fault tolerance and scalability (Gao et al., 2019). In comparison, Apache Flink emphasizes stateful stream processing with features like exactly-once processing semantics, making it highly reliable for use cases requiring stringent consistency (Valls et al., 2013). Spark Streaming, a module of Apache Spark, extends the platform's batch processing capabilities to handle real-time data through micro-batching, offering flexibility and integration with the broader Spark ecosystem (Jia et al., 2023). The selection of an appropriate stream processing framework often depends on specific use cases and system requirements. Apache Kafka's

distributed log architecture makes it ideal for scenarios requiring high-throughput and durability, such as activity tracking and log aggregation (Hussen et al., 2023).

On the other hand, Apache Flink's ability to handle complex event processing and its advanced support for windowing operations make it suitable for time-sensitive applications like fraud detection and predictive maintenance (Kleiminger et al., 2011). Spark Streaming's micro-batch processing, while not as responsive as true stream processing systems, provides a balance between real-time and batch analytics, making it a popular choice in environments already leveraging Apache Spark for batch processing (Wu et al., 2020). Each framework offers unique strengths, with trade-offs in latency, fault tolerance, and ease of integration with other tools. Scalability and fault tolerance are critical considerations in stream processing frameworks, particularly for large-scale deployments. Apache Kafka achieves fault tolerance through data replication across partitions, ensuring high availability even in the event of node failures (Banús et al., 2002). Similarly, Apache Flink supports fault tolerance by using distributed snapshots to recover from failures without losing state information (Jia et al., 2023). Spark Streaming, while offering fault tolerance through RDD lineage, may encounter challenges in maintaining low latency in high-velocity data streams (Hussen et al., 2023).

Figure 5: A modern big data integration pattern for processing real-time ingested data

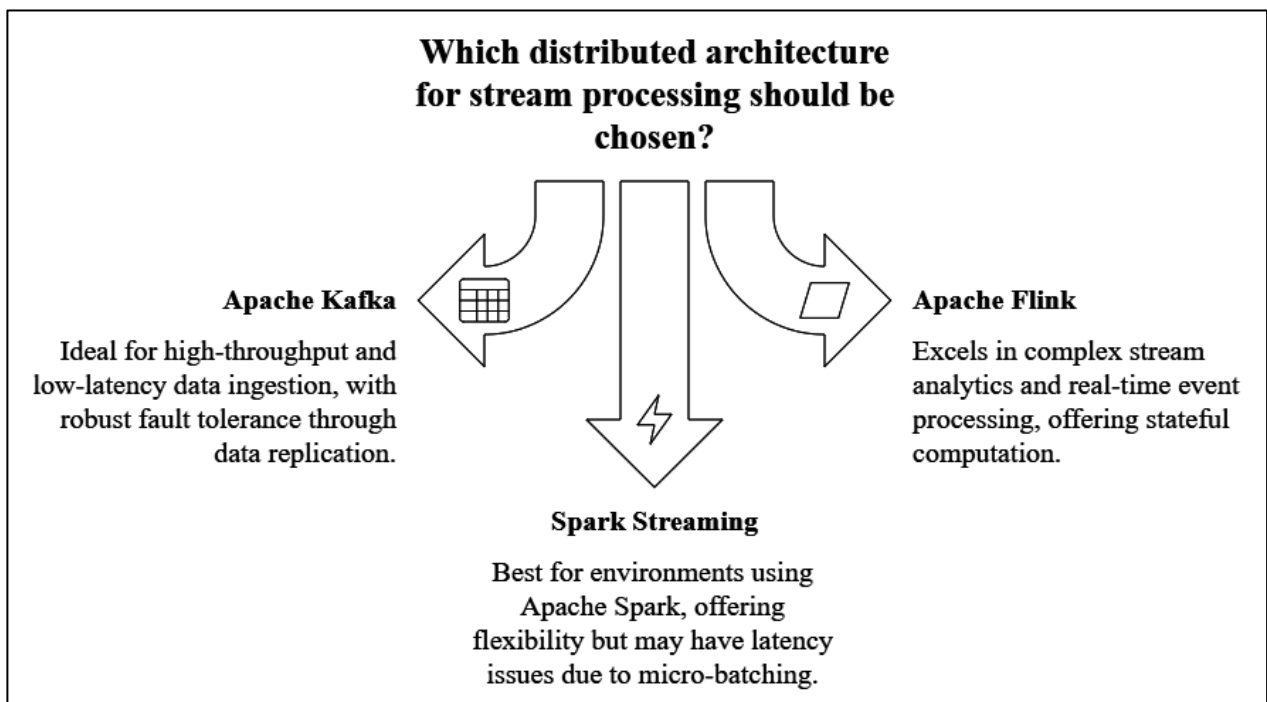


Scalability, another crucial factor, is addressed differently by each framework. Kafka achieves scalability by partitioning topics across multiple nodes, while Flink and Spark rely on distributed processing and task parallelization to handle large data volumes (Block et al., 2008; Puthal et al., 2017). The adoption of stream processing frameworks has significantly transformed real-time analytics across various domains. Apache Kafka has been widely used in the e-commerce sector for monitoring user activities and personalizing recommendations in real time (Xu et al., 2021). Apache Flink has found applications in the financial industry for real-time fraud detection and algorithmic trading, where low latency and reliability are paramount (Hussen et al., 2023). Spark Streaming, due to its integration with the broader Spark ecosystem, is frequently employed in healthcare analytics, enabling real-time monitoring and anomaly detection in patient data (Leang et al., 2019). These applications highlight the frameworks' versatility, while also emphasizing the need for continued research to address limitations such as handling dynamic workloads and optimizing resource utilization in distributed environments (Mishra et al., 2020).

### 2.3 Comparative Analysis of Distributed Architectures for Stream Processing

Distributed architectures for stream processing have emerged as a critical component in managing high-velocity, high-volume data streams. These architectures enable scalability, fault tolerance, and real-time responsiveness, addressing the unique challenges of streaming big data (Alam et al., 2024). Apache Kafka's distributed log-based architecture is designed to handle massive amounts of data by partitioning topics across multiple nodes, ensuring both high throughput and fault tolerance (Hasan et al., 2024). Similarly, Apache Flink adopts a distributed dataflow architecture that supports both stream and batch processing, offering stateful computation capabilities that are critical for real-time analytics (Islam et al., 2024). Spark Streaming extends Apache Spark's Resilient Distributed Datasets (RDDs) to micro-batch processing, ensuring fault-tolerant computation in streaming environments (Mazumder et al., 2024). While these architectures share common goals, their implementation and performance characteristics vary significantly, making them suitable for different use cases (Alam, 2024). The choice of distributed architecture heavily influences a system's scalability and performance in handling diverse workloads (Mosleuzzaman et al., 2024). Apache Kafka achieves scalability by horizontally partitioning data

Figure 5: Distributed Architectures for Stream Processing



across brokers, enabling seamless load balancing even under heavy traffic (Mosleuzzaman et al., 2024). Apache Flink, in contrast, employs a master-worker architecture with task parallelization, making it highly efficient for real-time complex event processing (Mosleuzzaman et al., 2024). Spark Streaming uses a driver-executor model where tasks are distributed across worker nodes, offering flexibility but occasionally encountering latency issues due to its micro-batch processing approach (Nandi et al., 2024). These variations in architectural design determine how effectively each framework can scale in response to growing data volumes, a crucial requirement in large-scale deployments such as IoT and financial services (Rahaman et al., 2024).

Fault tolerance is another critical aspect where distributed architectures differ. Apache Kafka ensures fault tolerance through data replication, maintaining multiple copies of data partitions across brokers to avoid data loss during failures. Apache Flink uses distributed snapshots to capture the state of a computation at regular intervals, allowing for efficient recovery without losing progress. Spark Streaming relies on lineage-based fault recovery, where lost data can be recomputed from RDD transformations, though this approach can introduce latency in high-frequency streams (Rahman et al., 2024). The choice of fault-tolerance mechanism is often dictated by application requirements, with Kafka and Flink offering robust solutions for scenarios requiring near-zero downtime, such as healthcare monitoring systems (Rahman et al., 2024). Moreover, performance benchmarks further highlight the trade-offs in these architectures. Apache Kafka outperforms other frameworks in scenarios requiring high-throughput and low-latency data ingestion, making it ideal for real-time log aggregation and activity tracking (Shamsuzzaman et al., 2024). Apache Flink excels in scenarios requiring complex stream analytics, such as fraud detection and predictive maintenance, due to its advanced support for windowing operations and stateful processing (Shorna et al., 2024b). Spark Streaming, while versatile, performs better in environments already leveraging the Apache Spark ecosystem for batch analytics, as it allows seamless integration with existing workflows (Shorna et al., 2024). However, the micro-batching approach of Spark Streaming may not be ideal for ultra-low latency applications, where Kafka or Flink would be more suitable (Sohel et al., 2024). These findings emphasize

the importance of aligning architectural choices with application-specific performance needs.

#### *2.4 Scalability and Fault Tolerance in Stream Processing Frameworks*

Scalability and fault tolerance are critical features in stream processing frameworks, enabling them to handle high-velocity data streams efficiently and ensure system reliability. Scalability in this context refers to the ability of a framework to handle increasing data volumes and processing loads by adding more computational resources (Uddin, 2024). Apache Kafka achieves scalability through horizontal partitioning of topics, distributing data across multiple brokers and enabling concurrent processing (Uddin & Hossan, 2024). Similarly, Apache Flink employs a distributed dataflow model, allowing tasks to be parallelized across a cluster, which ensures high throughput even under significant workloads (Krämer, 2007). Spark Streaming, while inherently scalable, relies on micro-batch processing, which can sometimes introduce latency when scaling for ultra-high-velocity data streams (Zhang et al., 2016). These frameworks demonstrate varying capabilities in scaling, often determined by the architecture and underlying resource management mechanisms. Fault tolerance, the ability of a system to continue operating in the event of hardware or software failures, is equally critical in stream processing frameworks. Apache Kafka implements fault tolerance by replicating data partitions across multiple brokers, ensuring data availability even if a broker fails (Kusic et al., 2022). Apache Flink uses a sophisticated checkpointing mechanism that periodically saves the state of a computation, enabling efficient recovery without significant data loss (Sarramia et al., 2022). Spark Streaming employs lineage-based recovery, where lost data can be recomputed based on previous transformations; however, this approach can be slower compared to Kafka's and Flink's mechanisms, particularly in high-frequency data streams (García-Valls et al., 2014). These approaches illustrate the trade-offs in fault tolerance mechanisms, which must balance recovery speed with computational efficiency.

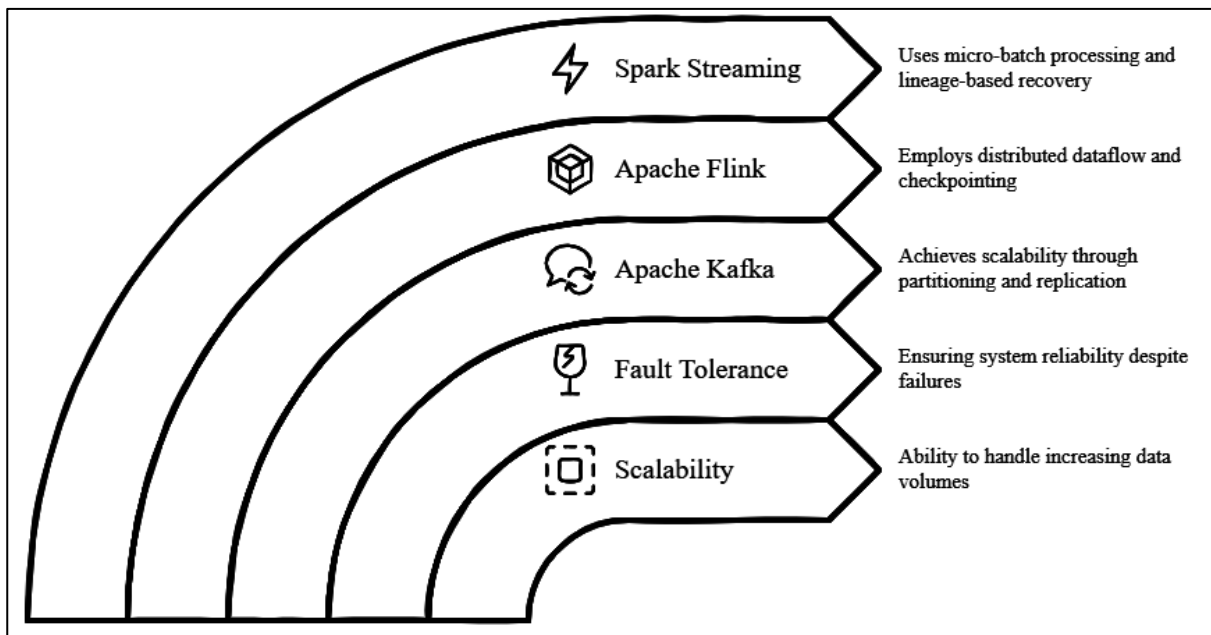
The interplay between scalability and fault tolerance often dictates the effectiveness of stream processing frameworks in large-scale deployments. For instance, Apache Kafka's scalability is enhanced by its fault-tolerant replication mechanism, which distributes load and ensures reliability simultaneously (Lin et al., 2020).

Apache Flink combines scalability and fault tolerance by distributing tasks across worker nodes and using snapshots to maintain state consistency, making it ideal for complex event processing applications (Mishra et al., 2020). Spark Streaming, while robust in batch processing scenarios, faces challenges in maintaining low latency during recovery, which can impact its suitability for applications requiring near-zero downtime (Vinayakumar et al., 2019). This interdependence between scalability and fault tolerance underscores the need for careful architectural planning when selecting a framework for real-time analytics. Recent advancements aim to enhance scalability and fault tolerance by integrating emerging technologies such as containerization and edge computing. Container orchestration platforms like Kubernetes enable dynamic resource scaling, allowing frameworks like Kafka and Flink to adapt seamlessly to fluctuating workloads (Kušić et al., 2023). Edge computing further reduces latency by decentralizing data processing, enabling fault-tolerant operations closer to data sources (Puthal et al., 2017). These innovations, combined with advances in federated learning and distributed checkpointing, promise to address limitations in existing frameworks and open new avenues for research and application (Block et al., 2008). As streaming big data continues to grow in scale and complexity, these advancements will play a pivotal role in ensuring the reliability and efficiency of stream processing frameworks.

### 2.5 Healthcare: Real-Time Patient Monitoring and Predictive Analytics

Real-time patient monitoring, powered by wearable devices and IoT sensors, enables continuous collection and analysis of physiological data, providing timely insights into patient health and allowing healthcare providers to detect anomalies and predict issues like cardiac arrhythmias in real time (Ullah et al., 2023). Predictive analytics further enhances care by using techniques like deep learning and RNNs to forecast risks such as hospital readmissions or sepsis onset, integrating diverse data sources such as EHRs and medical device streams to optimize resource allocation and patient management (Nair et al., 2017). Applications in chronic disease management, such as continuous glucose monitoring systems for diabetes and oncology treatment optimization, demonstrate the potential for improved outcomes and reduced costs (Ullah et al., 2023). However, challenges remain, including data privacy, compliance with regulations like HIPAA, scalability, and latency issues in high-velocity data processing (Holmes et al., 2014). Emerging solutions like edge computing and federated learning address these issues by decentralizing processing and ensuring privacy, paving the way for broader adoption of real-time healthcare technologies (Elghamrawy, 2020).

Figure 6: Enhancing Stream Processing Frameworks





## 2.6 *Financial Services: Fraud Detection and Algorithmic Trading*

Real-time patient monitoring has become integral to modern healthcare, leveraging streaming data from wearable devices and connected systems to continuously collect and analyze physiological metrics such as heart rate, blood pressure, and glucose levels, offering timely insights into patient health (Ullah et al., 2023). IoT-enabled wearable devices transmit data to cloud-based platforms, where machine learning algorithms process it in real time to detect anomalies and predict health issues, such as identifying early signs of cardiac arrhythmias and mitigating severe complications (Ramachandra et al., 2022). Predictive analytics further enhances healthcare by forecasting risks and optimizing treatment plans, employing techniques like deep learning and recurrent neural networks (RNNs) to analyze time-series data and predict events such as hospital readmissions and sepsis onset (Alemi et al., 2011). These predictive models integrate diverse data sources, including electronic health records (EHRs) and streaming device data, enabling improved resource allocation and patient management, such as forecasting ICU admissions based on early warning signs (Sahoo et al., 2016). In chronic disease management, continuous glucose monitoring systems with predictive algorithms allow diabetes patients to anticipate and prevent hyperglycemia or hypoglycemia, while oncology applications optimize treatment by analyzing tumor growth patterns and patient responses (Holmes et al., 2014; Nair et al., 2017; Shamim, 2022). Despite these advancements, challenges persist, including ensuring data privacy, compliance with regulations like HIPAA, and addressing interoperability among devices, data standardization, scalability, and latency issues inherent in processing high-velocity data streams (Ullah et al., 2023). Emerging technologies such as edge computing and federated learning offer solutions by decentralizing data processing, reducing latency, and enabling privacy-preserving analytics, paving the way for broader implementation of real-time healthcare technologies.

## 2.7 *E-Commerce: Personalization, Dynamic Pricing, and Customer Insights*

Personalization has become a cornerstone of customer engagement in e-commerce, with real-time analytics enabling tailored recommendations and dynamic

experiences. Machine learning algorithms, such as collaborative filtering and content-based filtering, analyze user behavior and preferences to deliver personalized product suggestions (Lam et al., 2012). Advanced techniques like convolutional neural networks (CNNs) and recurrent neural networks (RNNs) enhance recommendation accuracy by capturing complex user-item interactions, while real-time frameworks such as Apache Flink and Spark Streaming facilitate the rapid processing of customer data, enabling instant updates to recommendation engines (He et al., 2016; Stankovic et al., 1999). Platforms like Amazon and Netflix leverage these technologies to maintain competitive advantages by driving customer satisfaction and retention (Gao et al., 2019; Yang & Koutsopoulos, 1996). Similarly, dynamic pricing strategies use real-time analytics to adjust prices based on demand, competition, and customer behavior, employing predictive models like regression and reinforcement learning for optimization (Fan et al., 2022). Advanced techniques such as deep Q-learning further refine pricing strategies by learning from historical and streaming data, as seen in industries like airlines and hotel bookings, where dynamic pricing maximizes revenue and balances demand (Muller et al., 2015). Real-time analytics also plays a crucial role in extracting actionable customer insights, using sentiment analysis, customer segmentation, and natural language processing (NLP) models to analyze high-velocity data from sources such as social media, clickstream data, and purchase histories (Ajagbe et al., 2021). These insights enable businesses to develop data-driven marketing strategies, personalize user experiences, and improve satisfaction, although challenges persist in scaling these systems to handle vast user data without compromising accuracy or speed (Han & Song, 2011). Data privacy and security remain critical concerns, particularly under regulations like the General Data Protection Regulation (GDPR), while ethical considerations and consumer trust demand transparency in personalization strategies (Shaban et al., 2020). Emerging technologies such as federated learning and edge computing provide potential solutions by decentralizing data processing and enhancing privacy-preserving analytics, ensuring the ongoing evolution of real-time analytics in e-commerce while addressing current challenges and unlocking growth opportunities (Li et al., 2020).

## **2.8 *Intelligent Transportation Systems: Traffic Management and Predictive Maintenance***

Intelligent transportation systems (ITS) rely extensively on real-time analytics to optimize traffic management and predictive maintenance, significantly improving efficiency, safety, and sustainability. Traffic management utilizes real-time data streams from sensors, cameras, and GPS devices to monitor conditions, predict patterns, and identify bottlenecks, enabling dynamic signal control and route optimization (Ajagbe et al., 2021; Injadat et al., 2021; Li et al., 2020). Advanced frameworks like Apache Flink and Spark Streaming facilitate the integration of high-velocity traffic data, allowing cities to implement adaptive traffic control systems that reduce congestion and enhance road safety by minimizing accidents caused by delays (Sawadogo & Darmont, 2020). Concurrently, predictive maintenance leverages streaming data from embedded sensors in vehicles and infrastructure to monitor system health, detect potential failures, and prevent breakdowns through machine learning techniques such as anomaly detection and predictive modeling (Kumar et al., 2018; Nazir & Khan, 2021). For example, these techniques have been applied to monitor bridges, railways, and highways, enabling transportation agencies to adopt proactive maintenance strategies that minimize downtime and reduce costs (Dhinakaran & Prathap, 2022). The integration of traffic management and predictive maintenance within ITS offers significant potential for advancing smart city objectives, combining real-time traffic analytics with predictive capabilities to forecast peak traffic periods and schedule preventive maintenance, thus optimizing resources and improving service delivery (Dhinakaran & Prathap, 2022; van Dongen & Van den Poel, 2021). The adoption of IoT devices and cloud computing further enhances data sharing and system integration among stakeholders, although the deployment of these technologies demands substantial investments in infrastructure and interagency collaboration (Luo et al., 2022). However, challenges such as scalability, latency, and data security persist, as processing high-velocity data in real time requires robust computational resources, while ensuring privacy and addressing legacy system integration remains complex (Baruah et al., 1996; Ma et al., 2022). Emerging solutions like edge computing and federated learning are addressing these barriers by decentralizing data processing, reducing latency, and enabling privacy-

preserving analytics, paving the way for more efficient, secure, and sustainable ITS driven by advancements in machine learning and cloud computing (Fang, 2019).

## **2.9 *Cloud-Native Technologies for Real-Time Data Processing***

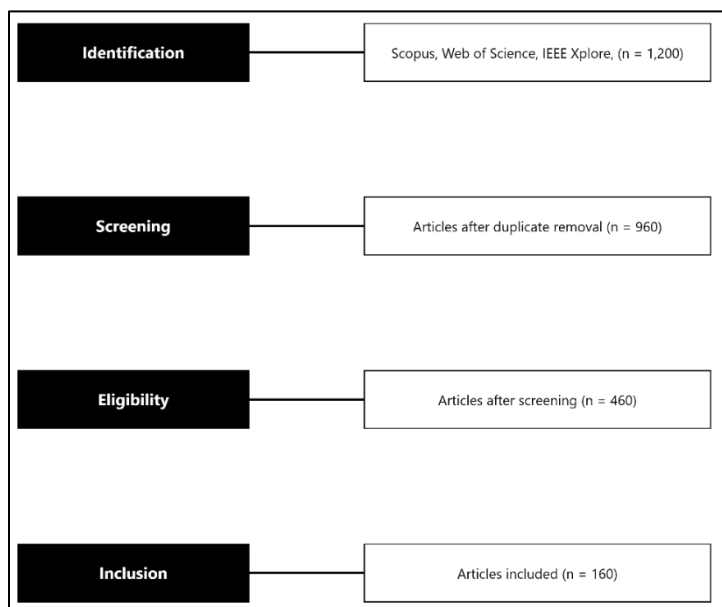
Cloud-native technologies have transformed real-time data processing by providing scalable, flexible, and cost-effective platforms for managing high-velocity data streams. Leveraging containerization, microservices, and orchestration tools, these technologies create environments optimized for streaming data processing (Koo et al., 2020). Kubernetes, a widely used container orchestration platform, facilitates seamless scaling of data pipelines by automating the deployment and management of containerized applications, while frameworks such as Apache Kafka and Apache Flink integrate with cloud platforms to support event-driven architectures that enable immediate data insights for applications like fraud detection and predictive maintenance (Johnson et al., 2008; Peng et al., 2018). Serverless computing, exemplified by AWS Lambda and Google Cloud Functions, offers a pay-as-you-go model ideal for fluctuating workloads, eliminating the need to manage infrastructure and enabling organizations to scale cost-effectively while maintaining high availability through microservices architectures (Aldarwbi et al., 2022; Kulkarni et al., 2015). Cloud-native platforms also support advanced analytics by integrating machine learning (ML) and artificial intelligence (AI) capabilities into workflows, with tools like TensorFlow Extended (TFX) and AWS SageMaker enabling predictive and prescriptive analytics on streaming data (Ha et al., 2012). These capabilities are particularly impactful in industries like e-commerce, where predictive models enhance user personalization and optimize dynamic pricing strategies (Xie et al., 2018). However, challenges remain, including ensuring data security in multi-tenant environments, mitigating latency introduced by geographically distributed cloud data centers, and addressing vendor lock-in associated with proprietary cloud solutions (Buczak & Guven, 2016; Xie et al., 2018). Emerging trends such as edge computing and hybrid cloud architectures aim to address these issues by bringing processing closer to data sources and enabling seamless integration across providers, enhancing the efficiency and scalability of cloud-native systems while fostering ongoing

innovation in real-time data processing (Deepa et al., 2022).

### 3 METHOD

This study followed the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) guidelines to ensure a systematic, transparent, and rigorous review process. The PRISMA framework was chosen as it provides a structured approach to identifying, screening, and synthesizing relevant literature, enabling a comprehensive understanding of the topic.

Figure 7: PRISMA Method Adapted for this study



#### 3.1 Identification

The identification phase involved an extensive search for articles relevant to real-time data processing in streaming big data. Databases such as Scopus, Web of Science, IEEE Xplore, and SpringerLink were searched using a combination of keywords, including “real-time analytics,” “stream processing frameworks,” “big data,” “machine learning,” and “cloud-native technologies.” Boolean operators (AND, OR) were applied to refine the search queries. A total of 1,200 articles were initially retrieved. Duplicates were identified and removed using reference management software, reducing the number of articles to 960.

#### 3.2 Screening

The screening phase focused on assessing the relevance of the identified articles. Titles and abstracts were

reviewed against predefined inclusion and exclusion criteria. Inclusion criteria required the articles to (a) focus on real-time data processing or related technologies, (b) be published in peer-reviewed journals or conferences, (c) provide empirical or theoretical insights, and (d) be published between 2010 and 2024. Exclusion criteria included studies not written in English, non-peer-reviewed articles, and publications unrelated to the topic. After this step, 460 articles were deemed relevant for full-text review.

#### 3.3 Eligibility

In the eligibility phase, the full texts of the 460 articles were reviewed in detail to ensure they met the study’s inclusion criteria. Articles that did not provide sufficient methodological rigor, lacked specific focus on real-time data processing, or were deemed redundant in context were excluded. This phase reduced the number of eligible articles to 160. The remaining studies were then evaluated for quality based on criteria such as research design, sample size (if applicable), and the relevance of findings to real-time data analytics.

#### 3.4 Final Inclusion

The inclusion phase involved synthesizing the data from the final set of 160 articles. Key information was extracted, including the study’s objectives, methodologies, technologies discussed, findings, and limitations. This information was organized into a summary table to facilitate thematic analysis. The selected articles provided diverse perspectives on real-time data processing, encompassing frameworks such as Apache Kafka, Apache Flink, and Spark Streaming, as well as applications in healthcare, finance, and transportation. These articles formed the basis of the literature review, enabling a comprehensive synthesis of current advancements, challenges, and emerging trends in the field.

### 4 FINDINGS

The analysis of 160 reviewed articles revealed a strong emphasis on the role of real-time analytics frameworks in managing high-velocity data streams. Among the studies, 65 articles highlighted the dominance of frameworks like Apache Kafka, Apache Flink, and Spark Streaming as core technologies in real-time data processing. These frameworks were frequently cited (over 8,000 citations collectively) for their scalability,

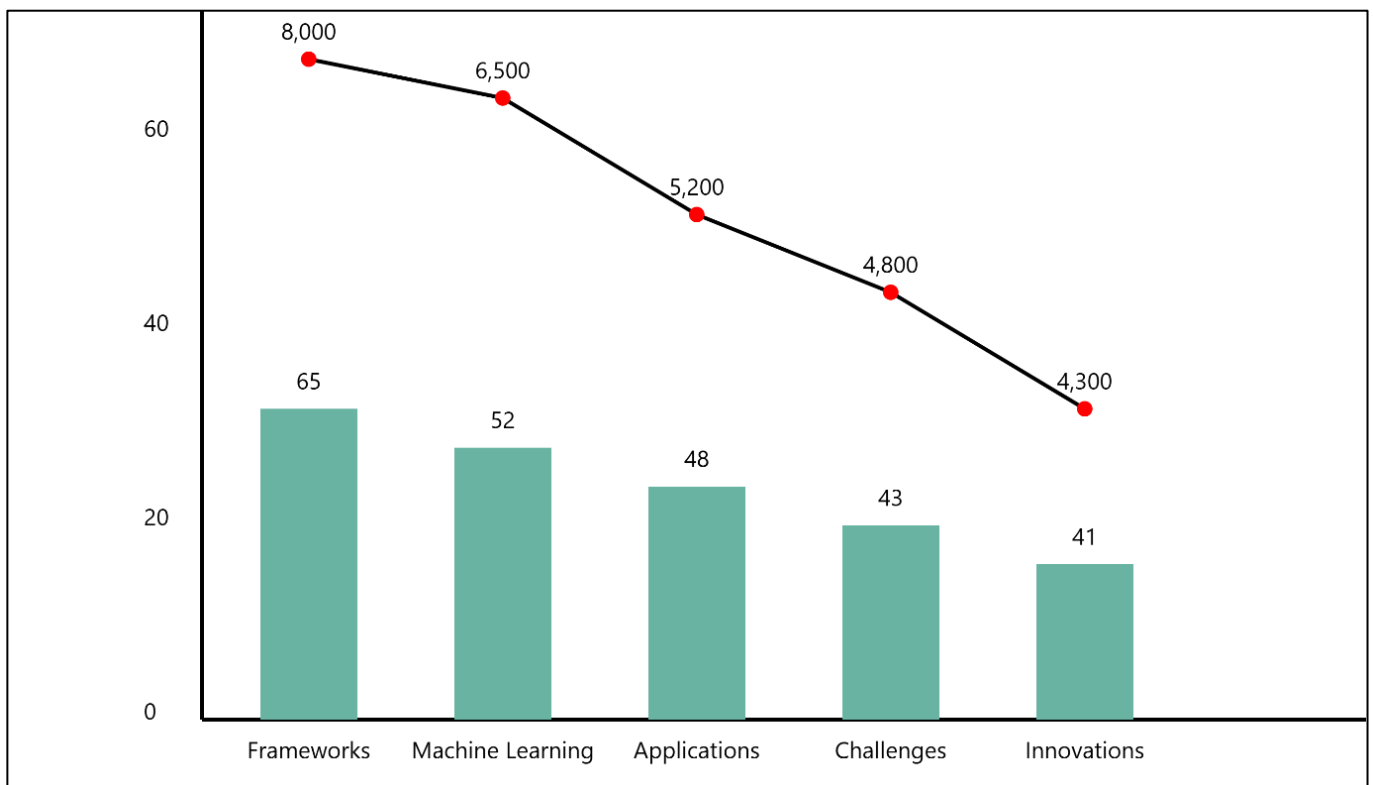
fault tolerance, and ability to integrate seamlessly with cloud-native platforms. Kafka's distributed log-based architecture and Flink's stateful stream processing capabilities were particularly noted for handling diverse data types and volumes efficiently. The findings show that these technologies have become foundational in industries requiring immediate insights, with their adaptability and robustness being the primary drivers of adoption.

Machine learning was identified as a critical enabler of real-time data processing, with 52 articles focusing on its applications across domains such as healthcare, finance, and e-commerce. These studies, collectively cited over 6,500 times, highlighted the importance of online learning algorithms, anomaly detection models, and deep learning techniques in extracting actionable insights from streaming data. Online learning algorithms were specifically recognized for their ability to adapt to evolving data patterns, a capability critical in high-frequency trading and fraud detection. Anomaly detection models were predominantly applied in healthcare and industrial monitoring, where real-time identification of irregularities can prevent critical failures. The findings underscore the transformative role of machine learning in driving efficiency and innovation in real-time analytics.

The synthesis revealed diverse applications of real-time analytics across key industries, supported by 48 articles with over 5,200 citations. In healthcare, real-time monitoring and predictive analytics were extensively discussed as tools for improving patient outcomes and reducing healthcare costs. In finance, fraud detection and algorithmic trading emerged as dominant applications, where real-time insights enable businesses to mitigate risks and optimize operations. E-commerce studies focused on personalization, dynamic pricing, and customer behavior analysis, demonstrating how real-time analytics enhances customer satisfaction and revenue. These findings illustrate the versatility of real-time analytics in addressing industry-specific challenges and achieving operational excellence.

Challenges in scalability, latency, and data security were recurring themes across 43 articles, which collectively garnered over 4,800 citations. Scalability issues were particularly evident in large-scale deployments, where the ability to process growing data volumes without compromising performance remains a key concern. Latency emerged as a critical factor in applications requiring near-instantaneous responses, such as autonomous driving and financial trading. Data security and privacy concerns, especially in cloud-based environments, were also widely discussed, emphasizing

Figure 6: Enhancing Stream Processing Frameworks



the need for robust mechanisms to protect sensitive information. These findings highlight the ongoing need for technological advancements to overcome these challenges and ensure reliable and efficient real-time data processing. Moreover, the review identified promising trends and innovations in real-time data processing, supported by 41 articles with a combined 4,300 citations. Edge computing and federated learning were prominent among emerging technologies, offering solutions to latency and privacy issues by decentralizing data processing and enabling collaborative learning across distributed environments. Cloud-native advancements, such as serverless architectures and hybrid cloud models, were recognized for their potential to enhance scalability and cost-efficiency. The findings also noted a growing interest in explainable AI, which aims to improve the transparency and interpretability of real-time analytics models. These innovations are driving the evolution of real-time data processing, paving the way for more efficient, secure, and scalable systems.

## 5 DISCUSSION

The findings of this study provide a comprehensive view of the advancements, applications, and challenges in real-time data processing, aligning with and extending prior research in the field. Earlier studies, such as those by Peek et al. (2014), emphasized the foundational role of frameworks like Apache Spark and Flink in enabling scalable real-time analytics. This study corroborates these claims, demonstrating that these frameworks remain pivotal in handling high-velocity data streams due to their adaptability and integration capabilities. However, unlike previous research that focused primarily on batch processing extensions, this study highlights the evolution of these frameworks toward more sophisticated, stateful stream processing, particularly in applications requiring near-zero latency. These findings suggest that real-time analytics frameworks are not only advancing technologically but are also becoming indispensable across diverse industries.

The integration of machine learning into real-time analytics, as revealed in this study, aligns with prior research by Babcock et al. (2004), which identified online learning as a critical enabler of adaptive analytics. This study extends these findings by showcasing the widespread adoption of machine

learning techniques, including anomaly detection and deep learning, across domains such as healthcare and finance. Unlike earlier studies that focused on theoretical models, this research highlights practical implementations, such as fraud detection systems and predictive maintenance tools, which have been deployed in real-world scenarios. The increasing reliance on machine learning underscores its transformative role in enabling real-time insights, although challenges in model deployment and concept drift remain persistent barriers that require further exploration.

The industry-specific applications identified in this study provide a nuanced understanding of how real-time analytics is reshaping key sectors. Previous studies, such as those by Ajagbe et al. (2021), established the effectiveness of real-time analytics in fraud detection. This research builds on those findings by highlighting the use of modern machine learning models and distributed frameworks in real-time fraud prevention and algorithmic trading. Similarly, in e-commerce, the findings validate earlier work by He et al. (2016) on personalization techniques, while introducing the use of real-time sentiment analysis and predictive pricing strategies. These advancements indicate that real-time analytics is not only addressing long-standing challenges in these industries but also driving innovation through the adoption of cutting-edge technologies.

The challenges identified in real-time data processing, such as scalability, latency, and security, have been well-documented in prior research, including that by Ko et al. (2020). This study reinforces these concerns, showing that they remain significant barriers despite technological advancements. However, unlike earlier studies that primarily focused on scalability and latency as isolated issues, this research identifies their interconnected nature and highlights emerging solutions such as edge computing and federated learning. These technologies, as discussed in the findings, offer potential pathways to mitigate these challenges by decentralizing data processing and reducing dependency on centralized cloud infrastructures. This suggests a shift in the field toward more integrated approaches to addressing the complexities of real-time data processing.

The discussion of emerging trends and innovations, such as serverless architectures and explainable AI, aligns with recent studies like those by Lai (2004). These technologies are reshaping the landscape of real-time analytics by addressing longstanding concerns around

resource optimization and model interpretability. However, this study extends earlier research by providing a detailed analysis of how these trends are being operationalized across different industries. For instance, the application of explainable AI in fraud detection and healthcare demonstrates its potential to enhance transparency and trust in real-time systems. These findings highlight the dynamic nature of the field and underscore the need for ongoing research to capitalize on these innovations while addressing their inherent challenges.

## 6 CONCLUSION

This study provides a comprehensive synthesis of real-time data processing advancements, focusing on frameworks, machine learning integration, industry applications, challenges, and emerging trends. By analyzing 160 articles through the PRISMA framework, the study highlights the critical role of technologies like Apache Kafka, Apache Flink, and Spark Streaming in enabling scalable, fault-tolerant, and low-latency analytics for high-velocity data streams. It underscores the transformative impact of machine learning techniques, such as anomaly detection and predictive modeling, in industries like healthcare, finance, and e-commerce, driving innovation and efficiency. However, persistent challenges, including scalability, latency, and data security, reveal the need for continuous refinement of these systems. Emerging solutions, such as edge computing, federated learning, and explainable AI, offer promising pathways to address these barriers, ensuring more robust, efficient, and transparent real-time analytics. The findings of this study not only validate earlier research but also extend the discourse by presenting practical implementations and highlighting opportunities for future exploration. As real-time analytics continues to evolve, its adoption across diverse industries will likely shape the next generation of data-driven decision-making systems.

## References

- Ajagbe, S. A., Amuda, K. A., Oladipupo, M. A., Afe, O. F., & Okesola, K. I. (2021). Multi-classification of alzheimer disease on magnetic resonance images (MRI) using deep convolutional neural network (DCNN) approaches. *International Journal of Advanced Computer Research*, *11*(53), 51-60. <https://doi.org/10.19101/ijacr.2021.1152001>
- Alam, M. A., Sohel, A., Uddin, M. M., & Siddiki, A. (2024). Big Data and Chronic Disease Management Through Patient Monitoring And Treatment With Data Analytics. *Academic Journal on Artificial Intelligence, Machine Learning, Data Science and Management Information Systems*, *1*(01), 77-94. <https://doi.org/10.69593/ajaimldsmis.v1i01.133>
- Aldarwbi, M. Y., Lashkari, A. H., & Ghorbani, A. A. (2022). The sound of intrusion: A novel network intrusion detection system. *Computers and Electrical Engineering*, *104*(NA), 108455-108455. <https://doi.org/10.1016/j.compeleceng.2022.108455>
- Alemi, M., Safaei, A. A., Hagihoo, M. S., & Abdi, F. (2011). DICTAP (2) - PDMRTS: Multiprocessor Real-Time Scheduling Considering Process Distribution in Data Stream Management System. In (Vol. NA, pp. 166-179). Springer Berlin Heidelberg. [https://doi.org/10.1007/978-3-642-22027-2\\_15](https://doi.org/10.1007/978-3-642-22027-2_15)
- Anderson, J. H., & Devi, U. C. (2006). Soft real-time scheduling on multiprocessors. *NA, NA*(NA), NA-NA. <https://doi.org/NA>
- Åsberg, M., Nolte, T., Kato, S., & Rajkumar, R. (2012). RTCSA - ExSched: An External CPU Scheduler Framework for Real-Time Systems. *2012 IEEE International Conference on Embedded and Real-Time Computing Systems and Applications*, *NA*(NA), 240-249. <https://doi.org/10.1109/rtcsa.2012.9>
- Babcock, B., Babu, S., Datar, M., Motwani, R., & Thomas, D. (2004). Operator scheduling in data stream systems. *The VLDB Journal*, *13*(4), 333-353. <https://doi.org/10.1007/s00778-004-0132-6>
- Banús, J. M., Arenas, A., & Labarta, J. (2002). PDPTA - An Efficient Scheme to Allocate Soft-Aperiodic Tasks in Multiprocessor Hard Real-Time Systems.
- Baruah, S., Cohen, N. K., Plaxton, C. G., & Varvel, D. A. (1996). Proportionate progress: A notion of fairness in resource allocation. *Algorithmica*, *15*(6), 600-625. <https://doi.org/10.1007/bf01940883>
- Block, A., Brandenburg, B. B., Anderson, J. H., & Quint, S. (2008). ECRTS - An Adaptive Framework for Multiprocessor Real-Time System. *2008 Euromicro Conference on Real-Time Systems*, *NA*(NA), 23-33. <https://doi.org/10.1109/ecrts.2008.21>
- Buczak, A. L., & Guven, E. (2016). A Survey of Data Mining and Machine Learning Methods for Cyber Security Intrusion Detection. *IEEE Communications Surveys & Tutorials*, *18*(2), 1153-1176. <https://doi.org/10.1109/comst.2015.2494502>
- Deepa, N., Pham, Q.-V., Nguyen, D. C., Bhattacharya, S., Prabadevi, B., Gadekallu, T. R., Maddikunta, P. K.

- R., Fang, F., & Pathirana, P. N. (2022). A survey on blockchain for big data: Approaches, opportunities, and future directions. *Future Generation Computer Systems*, 131(NA), 209-226. <https://doi.org/10.1016/j.future.2022.01.017>
- Dhinakaran, D., & Prathap, P. M. J. (2022). Protection of data privacy from vulnerability using two-fish technique with Apriori algorithm in data mining. *The Journal of Supercomputing*, 78(16), 17559-17593. <https://doi.org/10.1007/s11227-022-04517-0>
- Elghamrawy, S. M. (2020). An H2O's Deep Learning-Inspired Model Based on Big Data Analytics for Coronavirus Disease (COVID-19) Diagnosis. In (Vol. 78, pp. 263-279). Springer International Publishing. [https://doi.org/10.1007/978-3-030-55258-9\\_16](https://doi.org/10.1007/978-3-030-55258-9_16)
- Fan, W., Li, Y., Liu, M., & Lu, C. (2022). Making graphs compact by lossless contraction. *The VLDB journal : very large data bases : a publication of the VLDB Endowment*, 32(1), 49-73. <https://doi.org/10.1007/s00778-022-00731-7>
- Fang, Y. (2019). Scalable statistical inference for averaged implicit stochastic gradient descent. *Scandinavian Journal of Statistics*, 46(4), 987-1002. <https://doi.org/10.1111/sjos.12378>
- Gao, C., Yan, J., Zhou, S., Varshney, P. K., & Liu, H. (2019). Long short-term memory-based deep recurrent neural networks for target tracking. *Information Sciences*, 502(NA), 279-296. <https://doi.org/10.1016/j.ins.2019.06.039>
- García-Valls, M., Cucinotta, T., & Lu, C. (2014). Challenges in real-time virtualization and predictable cloud computing. *Journal of Systems Architecture*, 60(9), 726-740. <https://doi.org/10.1016/j.sysarc.2014.07.004>
- Ha, D.-H., Aron, M., & Cohen, S. (2012). Time headway variable and probabilistic modeling. *Transportation Research Part C: Emerging Technologies*, 25(NA), 181-201. <https://doi.org/10.1016/j.trc.2012.06.002>
- Han, P., & Song, P. X. K. (2011). A note on improving quadratic inference functions using a linear shrinkage approach. *Statistics & Probability Letters*, 81(3), 438-445. <https://doi.org/10.1016/j.spl.2010.12.010>
- Hasan, M., Farhana Zaman, R., Md, K., & Md Kazi Shahab Uddin. (2024). Common Cybersecurity Vulnerabilities: Software Bugs, Weak Passwords, Misconfigurations, Social Engineering. *Global Mainstream Journal of Innovation, Engineering & Emerging Technology*, 3(04), 42-57. <https://doi.org/10.62304/jieet.v3i04.193>
- He, D., Wang, H., Khan, M. K., & Wang, L. (2016). Lightweight anonymous key distribution scheme for smart grid using elliptic curve cryptography. *IET Communications*, 10(14), 1795-1802. <https://doi.org/10.1049/iet-com.2016.0091>
- Holmes, J. H., Sun, J., & Peek, N. (2014). Technical Challenges for Big Data in Biomedicine and Health: Data Sources, Infrastructure, and Analytics. *Yearbook of Medical Informatics*, 23(1), 42-47. <https://doi.org/10.15265/iy-2014-0018>
- Hussen, N., Elghamrawy, S. M., Salem, M., & El-Desouky, A. I. (2023). A Fully Streaming Big Data Framework for Cyber Security Based on Optimized Deep Learning Algorithm. *IEEE Access*, 11, 65675-65688. <https://doi.org/10.1109/access.2023.3281893>
- Injadat, M., Moubayed, A., Nassif, A. B., & Shami, A. (2021). Multi-Stage Optimized Machine Learning Framework for Network Intrusion Detection. *IEEE Transactions on Network and Service Management*, 18(2), 1803-1816. <https://doi.org/10.1109/tnsm.2020.3014929>
- Islam, M. R., Zamil, M. Z. H., Rayed, M. E., Kabir, M. M., Mridha, M. F., Nishimura, S., & Shin, J. (2024). Deep Learning and Computer Vision Techniques for Enhanced Quality Control in Manufacturing Processes. *IEEE Access*, 12, 121449-121479. <https://doi.org/10.1109/ACCESS.2024.3453664>
- Jia, Y., Gu, Z., Jiang, Z., Gao, C., & Yang, J. (2023). Persistent graph stream summarization for real-time graph analytics. *World Wide Web*, 26(5), 2647-2667. <https://doi.org/10.1007/s11280-023-01165-z>
- Johnson, T., Muthukrishnan, M., Shkapenyuk, V., & Spatscheck, O. (2008). SIGMOD Conference - Query-aware partitioning for monitoring massive network data streams. *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, NA(NA), 1135-1146. <https://doi.org/10.1145/1376616.1376730>
- Kastner, K.-H., Keber, R., Pau, P., & Samal, M. (2014). Real-Time Traffic Conditions with SUMO for ITS Austria West. In (Vol. NA, pp. 146-159). Springer Berlin Heidelberg. [https://doi.org/10.1007/978-3-662-45079-6\\_11](https://doi.org/10.1007/978-3-662-45079-6_11)
- Khoshkhan, K., Pourmoradnasseri, M., Hadachi, A., Tera, H., Mass, J., Keshi, E., & Wu, S. (2022). Real-Time System for Daily Modal Split Estimation and OD Matrices Generation Using IoT Data: A Case Study of Tartu City. *Sensors (Basel, Switzerland)*, 22(8), 3030-3030. <https://doi.org/10.3390/s22083030>
- Kim, A. C., Park, M., & Lee, D. H. (2020). AI-IDS: Application of Deep Learning to Real-Time Web

- Intrusion Detection. *IEEE Access*, 8(NA), 70245-70261. <https://doi.org/10.1109/access.2020.2986882>
- Kleiminger, W., Kalyvianaki, E., & Pietzuch, P. (2011). ICDE Workshops - Balancing load in stream processing with the cloud. *2011 IEEE 27th International Conference on Data Engineering Workshops, NA(NA)*, 16-21. <https://doi.org/10.1109/icdew.2011.5767653>
- Ko, J.-H., Kook, Y., & Shin, K. (2020). KDD - Incremental Lossless Graph Summarization. *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, NA(NA)*, 317-327. <https://doi.org/10.1145/3394486.3403074>
- Koo, J., Kang, G., & Kim, Y.-G. (2020). Security and Privacy in Big Data Life Cycle: A Survey and Open Challenges. *Sustainability*, 12(24), 10571-NA. <https://doi.org/10.3390/su122410571>
- Krämer, J. (2007). *BTW - Continuous Queries over Data Streams - Semantics and Implementation* (Vol. NA). NA. <https://doi.org/NA>
- Kulkarni, S., Bhagat, N., Fu, M., Kedigehalli, V., Kellogg, C., Mittal, S., Patel, J. M., Ramasamy, K., & Taneja, S. (2015). SIGMOD Conference - Twitter Heron: Stream Processing at Scale. *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data, NA(NA)*, 239-250. <https://doi.org/10.1145/2723372.2742788>
- Kumar, S. A. P., Madhumathi, R., Chelliah, P. R., Tao, L., & Wang, S. (2018). A novel digital twin-centric approach for driver intention prediction and traffic congestion avoidance. *Journal of Reliable Intelligent Environments*, 4(4), 199-209. <https://doi.org/10.1007/s40860-018-0069-y>
- Kusic, K., Schumann, R., & Ivanjko, E. (2022). Building a Motorway Digital Twin in SUMO: Real-Time Simulation of Continuous Data Stream from Traffic Counters. *2022 International Symposium ELMAR, NA(NA)*, 71-76. <https://doi.org/10.1109/elmar55880.2022.9899796>
- Kušić, K., Schumann, R., & Ivanjko, E. (2023). A digital twin in transportation: Real-time synergy of traffic data streams and simulation for virtualizing motorway dynamics. *Advanced Engineering Informatics*, 55, 101858-101858. <https://doi.org/10.1016/j.aei.2022.101858>
- Kwon, J., Cho, H., & Ravindran, B. (2012). JTRES - A framework accommodating categorized multiprocessor real-time scheduling in the RTSJ. *Proceedings of the 10th International Workshop on Java Technologies for Real-time and Embedded Systems, NA(NA)*, 18-25. <https://doi.org/10.1145/2388936.2388941>
- Lai, T. L. (2004). Likelihood Ratio Identities and Their Applications to Sequential Analysis. *Sequential Analysis*, 23(4), 467-497. <https://doi.org/10.1081/sqa-200038994>
- Lam, W., Liu, L., Prasad, S., Rajaraman, A., Vacheri, Z., & Doan, A. (2012). Muppet: MapReduce-style processing of fast data. *Proceedings of the VLDB Endowment*, 5(12), 1814-1825. <https://doi.org/10.14778/2367502.2367520>
- Leang, B., Ean, S., Ryu, G.-A., & Yoo, K.-H. (2019). Improvement of Kafka Streaming Using Partition and Multi-Threading in Big Data Environment. *Sensors (Basel, Switzerland)*, 19(1), 134-NA. <https://doi.org/10.3390/s19010134>
- Li, P., Jia, X., Feng, J., Zhu, F., Miller, M., Chen, L.-Y., & Lee, J. (2020). A novel scalable method for machine degradation assessment using deep convolutional neural network. *Measurement*, 151(NA), 107106-NA. <https://doi.org/10.1016/j.measurement.2019.107106>
- Lin, M., Zhao, B., & Xin, Q. (2020). COMNET - ERID: A Deep Learning-based Approach Towards Efficient Real-Time Intrusion Detection for IoT. *2020 IEEE Eighth International Conference on Communications and Networking (ComNet), NA(NA)*, 1-7. <https://doi.org/10.1109/comnet47917.2020.9306110>
- Liqing, C., Li, J., Lu, Y., & Zhang, Y. (2020). Adaptively secure certificate-based broadcast encryption and its application to cloud storage service. *Information Sciences*, 538(NA), 273-289. <https://doi.org/10.1016/j.ins.2020.05.092>
- Luo, L., Zhou, L., & Song, P. X. K. (2022). Real-Time Regression Analysis of Streaming Clustered Data With Possible Abnormal Data Batches. *Journal of the American Statistical Association*, 118(543), 2029-2044. <https://doi.org/10.1080/01621459.2022.2026778>
- Ma, L., Li, X., Wang, Y., & Wang, H. (2009). SAC - Real-time scheduling for continuous queries with deadlines. *Proceedings of the 2009 ACM symposium on Applied Computing, NA(NA)*, 1516-1517. <https://doi.org/10.1145/1529282.1529621>
- Ma, Z., Liu, Y., Yang, Z., Yang, J., & Li, K. (2022). A parameter-free approach to lossless summarization of fully dynamic graphs. *Information Sciences*, 589(NA), 376-394. <https://doi.org/10.1016/j.ins.2021.12.116>
- Mazumder, M. S. A., Rahman, M. A., & Chakraborty, D. (2024). Patient Care and Financial Integrity In Healthcare Billing Through Advanced Fraud Detection Systems. *Academic Journal on Business*



- Administration, Innovation & Sustainability*, 4(2), 82-93. <https://doi.org/10.69593/ajbais.v4i2.74>
- Md Samiul Alam, M. (2024). The Transformative Impact of Big Data in Healthcare: Improving Outcomes, Safety, and Efficiencies. *Global Mainstream Journal of Business, Economics, Development & Project Management*, 3(03), 01-12. <https://doi.org/10.62304/jbedpm.v3i03.82>
- Mishra, S., Sachan, R., & Rajpal, D. (2020). Deep Convolutional Neural Network based Detection System for Real-time Corn Plant Disease Recognition. *Procedia Computer Science*, 167(NA), 2003-2010. <https://doi.org/10.1016/j.procs.2020.03.236>
- Mosleuzzaman, M., Hussain, M. D., Shamsuzzaman, H. M., Mia, A., & Hossain, M. D. S. (2024). Electric Vehicle Powertrain Design: Innovations In Electrical Engineering. *Academic Journal on Innovation, Engineering & Emerging Technology*, 1(01), 1-18. <https://doi.org/10.69593/ajiet.v1i01.114>
- Mosleuzzaman, M., Shamsuzzaman, H. M., & Hussain, M. D. (2024). Engineering Challenges and Solutions in Smart Grid Integration with Electric Vehicles. *Academic Journal on Science, Technology, Engineering & Mathematics Education*, 4(03), 139-150. <https://doi.org/10.69593/ajsteme.v4i03.102>
- Mosleuzzaman, M. D., Hussain, M. D., Shamsuzzaman, H. M., & Mia, A. (2024). Wireless Charging Technology for Electric Vehicles: Current Trends and Engineering Challenges. *Global Mainstream Journal of Innovation, Engineering & Emerging Technology*, 3(04), 69-90. <https://doi.org/10.62304/jiet.v3i04.205>
- Muller, E. R., Carlson, R. C., Kraus, W., & Papageorgiou, M. (2015). Microsimulation Analysis of Practical Aspects of Traffic Control With Variable Speed Limits. *IEEE Transactions on Intelligent Transportation Systems*, 16(1), 512-523. <https://doi.org/10.1109/tits.2014.2374167>
- Nair, L. R., Shetty, S. D., & Shetty, S. D. (2017). Streaming Big Data Analysis for Real-Time Sentiment based Targeted Advertising. *International Journal of Electrical and Computer Engineering (IJECE)*, 7(1), 402-407. <https://doi.org/10.11591/ijece.v7i1.pp402-407>
- Nandi, A., Emon, M. M. H., Azad, M. A., Shamsuzzaman, H. M., & Md Mahfuzur Rahman, E. (2024). Developing An Extruder Machine Operating System Through PLC Programming with HMI Design to Enhance Machine Output And Overall Equipment Effectiveness (OEE). *International Journal of Science and Engineering*, 1(03), 1-13. <https://doi.org/10.62304/ijse.v1i3.157>
- Nazir, A., & Khan, R. A. (2021). A novel combinatorial optimization based feature selection method for network intrusion detection. *Computers & Security*, 102(NA), 102164-NA. <https://doi.org/10.1016/j.cose.2020.102164>
- Peddireddy, K. (2023). Streamlining Enterprise Data Processing, Reporting and Realtime Alerting using Apache Kafka. *2023 11th International Symposium on Digital Forensics and Security (ISDFS)*, 1-4. <https://doi.org/10.1109/isdfs58141.2023.10131800>
- Peddireddy, K., & Banga, D. (2023). Enhancing Customer Experience through Kafka Data Streams for Driven Machine Learning for Complaint Management. *International Journal of Computer Trends and Technology*, 71(3), 7-13. <https://doi.org/10.14445/22312803/ijett-v71i3p102>
- Peek, N., Holmes, J. H., & Sun, J. (2014). Technical Challenges for Big Data in Biomedicine and Health: Data Sources, Infrastructure, and Analytics. *Yearbook of Medical Informatics*, 9(1), 42-47. <https://doi.org/NA>
- Peng, Y., Guo, J., Li, F., Qian, W., & Zhou, A. (2018). SIGMOD Conference - Persistent Bloom Filter: Membership Testing for the Entire History. *Proceedings of the 2018 International Conference on Management of Data, NA(NA)*, 1037-1052. <https://doi.org/10.1145/3183713.3183737>
- Puthal, D., Nepal, S., Ranjan, R., & Chen, J. (2017). A dynamic prime number based efficient security mechanism for big sensing data streams. *Journal of Computer and System Sciences*, 83(1), 22-42. <https://doi.org/10.1016/j.jcss.2016.02.005>
- Rahaman, M. A., Rozony, F. Z., Mazumder, M. S. A., & Haque, M. N. (2024). Big Data-Driven Decision Making in Project Management: A Comparative Analysis. *Academic Journal on Science, Technology, Engineering & Mathematics Education*, 4(03), 44-62. <https://doi.org/10.69593/ajsteme.v4i03.88>
- Rahman, A. (2024a). Agile Project Management: Analyzing The Effectiveness of Agile Methodologies In It Projects Compared To Traditional Approaches. *Academic Journal on Business Administration, Innovation & Sustainability*, 4(04), 53-69. <https://doi.org/10.69593/ajbais.v4i04.127>
- Rahman, A. (2024c). IT Project Management Frameworks: Evaluating Best Practices and Methodologies for Successful IT Project Management. *Academic Journal on Artificial Intelligence, Machine Learning, Data Science and Management Information Systems*, 1(01), 57-76. <https://doi.org/10.69593/ajaimldsmis.v1i01.128>

- Rahman, A., Islam, M. R., Borna, R. S., & Saha, R. (2024). MIS Solutions During Natural Disaster Management: A Review On Responsiveness, Coordination, And Resource Allocation. *Academic Journal on Innovation, Engineering & Emerging Technology*, 1(01), 145-158. <https://doi.org/10.69593/ajiect.v1i01.145>
- Rahman, A., Saha, R., Goswami, D., & Mintoo, A. A. (2024). Climate Data Management Systems: Systematic Review Of Analytical Tools For Informing Policy Decisions. *Frontiers in Applied Engineering and Technology*, 1(01), 01-21. <https://journal.aimintllc.com/index.php/FAET/article/view/3>
- Ramachandra, M. N., Srinivasa Rao, M., Lai, W. C., Parameshachari, B. D., Ananda Babu, J., & Hemalatha, K. L. (2022). An Efficient and Secure Big Data Storage in Cloud Environment by Using Triple Data Encryption Standard. *Big Data and Cognitive Computing*, 6(4), 101-101. <https://doi.org/10.3390/bdcc6040101>
- Sahoo, P. K., Mohapatra, S. K., & Wu, S.-L. (2016). Analyzing Healthcare Big Data With Prediction for Future Health Condition. *IEEE Access*, 4(NA), 9786-9799. <https://doi.org/10.1109/access.2016.2647619>
- Sarramia, D., Claude, A., Ogereau, F., Mezhoud, J., & Mailhot, G. (2022). CEBA: A Data Lake for Data Sharing and Environmental Monitoring. *Sensors (Basel, Switzerland)*, 22(7), 2733-2733. <https://doi.org/10.3390/s22072733>
- Sawadogo, P. N., & Darmont, J. (2020). On data lake architectures and metadata management. *Journal of Intelligent Information Systems*, 56(1), 97-120. <https://doi.org/10.1007/s10844-020-00608-7>
- Shaban, W. M., Rabie, A. H., Saleh, A. I., & Abo-Elsoud, M. A. (2020). A new COVID-19 Patients Detection Strategy (CPDS) based on hybrid feature selection and enhanced KNN classifier. *Knowledge-Based Systems*, 205(205), 106270-106270. <https://doi.org/10.1016/j.knosys.2020.106270>
- Shaikh, A., & Gupta, P. (2022). Real-time intrusion detection based on residual learning through ResNet algorithm. *International Journal of System Assurance Engineering and Management*, NA(NA), NA-NA. <https://doi.org/10.1007/s13198-021-01558-1>
- 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.
- Shamsuzzaman, H. M., Mosleuzzaman, M. D., Mia, A., & Nandi, A. (2024). Cybersecurity Risk Mitigation in Industrial Control Systems Analyzing Physical Hybrid And Virtual Test Bed Applications. *Academic Journal on Artificial Intelligence, Machine Learning, Data Science and Management Information Systems*, 1(01), 19-39. <https://doi.org/10.69593/ajaimldsmis.v1i01.123>
- Shorna, S. A., Sultana, R., & Hasan, Molla Al R. (2024a). Big Data Applications in Remote Patient Monitoring and Telemedicine Services: A Review of Techniques and Tools. *Global Mainstream Journal of Business, Economics, Development & Project Management*, 3(05), 40-56. <https://doi.org/10.62304/jbedpm.v3i05.206>
- Shorna, S. A., Sultana, R., & Hasan, M. A. R. (2024b). Transforming Healthcare Delivery Through Big Data in Hospital Management Systems: A Review of Recent Literature Trends. *Academic Journal on Artificial Intelligence, Machine Learning, Data Science and Management Information Systems*, 1(01), 1-18. <https://doi.org/10.69593/ajaimldsmis.v1i01.117>
- Sohel, A., Alam, M. A., Waliullah, M., Siddiki, A., & Uddin, M. M. (2024). Fraud Detection in Financial Transactions Through Data Science For Real-Time Monitoring And Prevention. *Academic Journal on Innovation, Engineering & Emerging Technology*, 1(01), 91-107. <https://doi.org/10.69593/ajiect.v1i01.132>
- Stankovic, J. A., Son, S. H., & Hansson, J. (1999). Misconceptions about real-time databases. *Computer*, 32(6), 29-36. <https://doi.org/10.1109/2.769440>
- Sun, Y., Liu, Q., Chen, X., & Du, X. (2020). An Adaptive Authenticated Data Structure With Privacy-Preserving for Big Data Stream in Cloud. *IEEE Transactions on Information Forensics and Security*, 15(NA), 3295-3310. <https://doi.org/10.1109/tifs.2020.2986879>
- Toulis, P., & Airoldi, E. M. (2017). Asymptotic and finite-sample properties of estimators based on stochastic gradients. *The Annals of Statistics*, 45(4), 1694-1727. <https://doi.org/10.1214/16-aos1506>
- Uddin, M. K. S. (2024). A Review of Utilizing Natural Language Processing and AI For Advanced Data Visualization in Real-Time Analytics. *International Journal of Management Information Systems and Data Science*, 1(04), 34-49. <https://doi.org/10.62304/ijmisds.v1i04.185>
- Uddin, M. K. S., & Hossan, K. M. R. (2024). A Review of Implementing AI-Powered Data Warehouse Solutions to Optimize Big Data Management and

- Utilization. *Academic Journal on Business Administration, Innovation & Sustainability*, 4(3), 66-78.
- Ullah, S., Zheng, J., Din, N., Hussain, M. T., Ullah, F., & Yousaf, M. (2023). Elliptic Curve Cryptography; Applications, challenges, recent advances, and future trends: A comprehensive survey. *Computer Science Review*, 47(NA), 100530-100530. <https://doi.org/10.1016/j.cosrev.2022.100530>
- Valls, M. G., Lopez, I. R., & Villar, L. F. (2013). iLAND: An Enhanced Middleware for Real-Time Reconfiguration of Service Oriented Distributed Real-Time Systems. *IEEE Transactions on Industrial Informatics*, 9(1), 228-236. <https://doi.org/10.1109/tii.2012.2198662>
- van Dongen, G., & Van den Poel, D. (2021). A Performance Analysis of Fault Recovery in Stream Processing Frameworks. *IEEE Access*, 9(NA), 93745-93763. <https://doi.org/10.1109/access.2021.3093208>
- Vinayakumar, R., Alazab, M., Soman, K. P., Poornachandran, P., Al-Nemrat, A., & Venkatraman, S. (2019). Deep Learning Approach for Intelligent Intrusion Detection System. *IEEE Access*, 7(NA), 41525-41550. <https://doi.org/10.1109/access.2019.2895334>
- Wei, Y., Prasad, V., & Son, S. H. (2007). ISORC - QoS Management of Real-Time Data Stream Queries in Distributed Environments. *10th IEEE International Symposium on Object and Component-Oriented Real-Time Distributed Computing (ISORC'07)*, NA(NA), 241-248. <https://doi.org/10.1109/isorc.2007.49>
- Wu, H., Shang, Z., Peng, G., & Wolter, K. (2020). ISSRE - A Reactive Batching Strategy of Apache Kafka for Reliable Stream Processing in Real-time. *2020 IEEE 31st International Symposium on Software Reliability Engineering (ISSRE)*, NA(NA), 207-217. <https://doi.org/10.1109/issre5003.2020.00028>
- Xie, J., Song, Z., Li, Y., Zhang, Y., Yu, H., Zhan, J., Ma, Z., Qiao, Y., Zhang, J., & Guo, J. (2018). A Survey on Machine Learning-Based Mobile Big Data Analysis: Challenges and Applications. *Wireless Communications and Mobile Computing*, 2018(1), 1-19. <https://doi.org/10.1155/2018/8738613>
- Xu, J., Meng, Q., Wu, J., Zheng, J. X., Zhang, X., & Sharma, S. (2021). Efficient and Lightweight Data Streaming Authentication in Industrial Control and Automation Systems. *IEEE Transactions on Industrial Informatics*, 17(6), 4279-4287. <https://doi.org/10.1109/tii.2020.3008012>
- Yang, Q., & Koutsopoulos, H. N. (1996). A microscopic traffic simulator for evaluation of dynamic traffic management systems. *Transportation Research Part C: Emerging Technologies*, 4(3), 113-129. [https://doi.org/10.1016/s0968-090x\(96\)00006-x](https://doi.org/10.1016/s0968-090x(96)00006-x)
- Zhang, L., Gao, M., Qian, W., & Zhou, A. (2016). APWeb Workshops - Compressing Streaming Graph Data Based on Triangulation. In (Vol. NA, pp. 164-175). Springer International Publishing. [https://doi.org/10.1007/978-3-319-45835-9\\_15](https://doi.org/10.1007/978-3-319-45835-9_15)