(RESEARCH ARTICLE)

# Real-time data stream processing in large-scale systems

Sanjay Lote [1, *], Praveena K B [2] and Durugappa Patrer [2]

[1] Department of Computer Science Engineering, Government Polytechnic Athani, Karnataka, India.
[2] Department of Computer Science Engineering, Government Polytechnic Harihar, Karnataka, India.

## Abstract

Real-time data stream processing has emerged as a crucial element in modern large-scale systems, facilitating rapid decision-making and real-time analytics across various domains. As data volumes continue to grow exponentially, the need for efficient, scalable, and fault-tolerant stream processing solutions has become more pressing. This paper provides a comprehensive exploration of real-time data processing architectures, highlighting key components such as distributed stream processing frameworks, parallel data pipelines, and event-driven computing models. The study delves into state-of-the-art technologies, including Apache Kafka, Apache Flink, and Spark Streaming, which enable seamless ingestion, processing, and storage of high-velocity data. Furthermore, we analyze the significance of fault-tolerant designs, low-latency data handling, and adaptive load balancing mechanisms to ensure uninterrupted system performance. Key performance metrics, such as throughput, latency, and resource utilization, are examined to assess the effectiveness of various approaches. Additionally, the paper discusses scalability challenges, including data partitioning strategies, cluster management techniques, and resource elasticity in cloud-based and edge-computing environments. Practical applications of real-time data stream processing are explored across multiple sectors, including finance, healthcare, and the Internet of Things (IoT), demonstrating its transformative impact on fraud detection, patient monitoring, and smart city implementations. The findings are supported by empirical evaluations, with figures, tables, and bar charts illustrating comparative performance analyses and efficiency metrics of different processing frameworks. This research contributes valuable insights into optimizing real-time data stream processing for future advancements in large-scale intelligent systems.

**Keywords:** Real-Time Data Processing; Stream Processing; Distributed Computing; Low-Latency Analytics; Apache Flink; Apache Kafka; Scalability; Fault Tolerance

## 1. Introduction

The exponential growth of data generation from diverse sources, including Internet of Things (IoT) devices, social media platforms, financial transactions, and enterprise applications, has created an urgent need for efficient real-time data processing techniques. Traditional batch processing methods, which collect and process data in fixed intervals, are increasingly inadequate for modern applications that require immediate insights and decision-making capabilities. As organizations strive to extract value from massive, high-velocity data streams, real-time data stream processing has emerged as a key solution for enabling continuous and dynamic analysis.

Unlike batch processing, which introduces delays due to the accumulation and scheduled execution of data processing tasks, real-time stream processing operates on continuous data flows, providing low-latency analytics and rapid response mechanisms. This shift from periodic processing to real-time computation has driven significant advancements in data-driven applications, particularly in sectors where timely insights are critical. For instance, fraud

* Corresponding author: Sanjay Lote

detection systems in banking, predictive maintenance in industrial operations, and real-time monitoring of patient health in healthcare require instant data analysis to prevent risks and optimize performance.

A real-time data stream processing system consists of multiple interconnected components designed to handle high-throughput data ingestion, transformation, and analysis. These systems leverage distributed computing frameworks to ensure scalability and fault tolerance while maintaining stringent latency requirements. Technologies such as Apache Kafka, Apache Flink, Apache Storm, and Spark Streaming provide the necessary infrastructure for managing real-time data pipelines by enabling parallel processing, stateful computation, and event-driven architectures. The effective integration of these tools allows organizations to harness the full potential of real-time analytics.

One of the fundamental challenges in real-time data stream processing is managing the sheer volume, velocity, and variety of incoming data while maintaining accuracy and efficiency. Data partitioning, stream synchronization, and adaptive load balancing are crucial considerations to ensure the system can scale dynamically based on workload fluctuations. Additionally, ensuring data consistency, handling late-arriving events, and implementing fault-tolerant mechanisms are critical aspects of designing a robust real-time processing architecture.

Scalability plays a crucial role in large-scale real-time data processing systems, as they must accommodate fluctuating workloads without performance degradation. Cloud-based and edge computing solutions are increasingly being adopted to enhance scalability and reduce latency by distributing computation closer to data sources. Techniques such as data parallelism, micro-batching, and in-memory processing further optimize performance, enabling organizations to process petabytes of data with minimal delays.

Security and privacy are also significant concerns in real-time data processing, especially in applications that handle sensitive information, such as financial transactions and healthcare records. Implementing end-to-end encryption, access controls, and anomaly detection mechanisms is essential to protect data integrity and prevent security breaches. Moreover, compliance with regulatory frameworks such as GDPR and HIPAA adds additional complexity to real-time data governance, requiring organizations to adopt stringent data protection measures while maintaining system efficiency.

The applications of real-time data stream processing extend across various industries, demonstrating its transformative impact on business operations and decision-making. In e-commerce, real-time recommendation engines enhance user experience by delivering personalized content based on browsing behavior. In smart cities, sensor-driven data analytics enable intelligent traffic management and environmental monitoring. Similarly, in the energy sector, real-time monitoring of power grids enhances grid stability and predictive maintenance, reducing operational costs and downtime[1].

This paper provides a comprehensive analysis of real-time data stream processing by exploring its architectural foundations, key enabling technologies, and practical applications across industries. Additionally, we discuss the challenges associated with large-scale deployment, including scalability, fault tolerance, and security concerns. Through empirical evaluations and comparative performance analysis, this study aims to offer valuable insights into optimizing real-time data processing for emerging intelligent systems and high-throughput applications.

## 2. Architecture of Real-Time Data Stream Processing

Real-time data stream processing systems are designed to handle massive volumes of continuously generated data while ensuring low-latency processing and high throughput. Unlike traditional batch processing systems, which store and process data in scheduled intervals, real-time systems operate on data streams as they arrive. To achieve this level of responsiveness and efficiency, these systems rely on distributed architectures that ensure scalability, fault tolerance, and seamless integration with various data sources and destinations. A typical real-time stream processing system consists of several key components, each playing a crucial role in the end-to-end data pipeline[2].

### 2.1. Data Sources

The first component of a real-time stream processing system is the data source, which continuously generates raw data streams. These sources can originate from various domains, including:

- IoT Sensors and Edge Devices: Smart devices, industrial sensors, and wearables produce telemetry data, environmental readings, and operational parameters.

- Application and System Logs: Web applications, cloud services, and enterprise software generate logs that track user interactions, performance metrics, and security events.
- Financial Transactions: Banking and fintech systems produce continuous streams of transactional data that require real-time fraud detection and analytics.
- Social Media and Web Activity: Platforms such as Twitter, Facebook, and e-commerce websites generate massive volumes of real-time user-generated content, interactions, and browsing behavior.

Since data sources vary in format, frequency, and reliability, real-time systems must be capable of handling structured, semi-structured, and unstructured data while ensuring efficient ingestion and processing.

## 2.2. Stream Ingestion Layer

The stream ingestion layer is responsible for collecting, buffering, and transmitting data from diverse sources into the processing pipeline. This layer ensures high-throughput data ingestion and efficient event routing while maintaining fault tolerance. Popular middleware solutions for this layer include:

- Apache Kafka: A widely used distributed messaging system that enables high-throughput, fault-tolerant data ingestion with topic-based message distribution.
- RabbitMQ: A message broker designed for lightweight messaging and integration with various real-time applications.
- Amazon Kinesis: A cloud-based alternative that enables scalable real-time data ingestion and stream processing.

This layer plays a crucial role in decoupling data producers from consumers, ensuring that real-time data flows efficiently into the processing engine without overloading system resources.

## 2.3. Processing Engine

At the core of a real-time data stream processing system lies the processing engine, which continuously analyzes, transforms, and aggregates incoming data streams. This layer is responsible for executing complex computations, pattern detection, anomaly identification, and real-time analytics. Leading frameworks used in this stage include:

- Apache Flink: A robust, stateful stream processing framework that supports event time processing, low-latency execution, and fault tolerance.
- Apache Spark Streaming: A micro-batch-based processing engine that extends Apache Spark's capabilities to handle real-time data processing workloads.
- Apache Storm: A distributed real-time computation framework designed for high-speed event processing.

These processing engines enable real-time computation by employing techniques such as windowing, streaming joins, and in-memory processing to achieve near-instantaneous data insights. Additionally, they support integration with machine learning models and AI-driven analytics for more advanced decision-making.
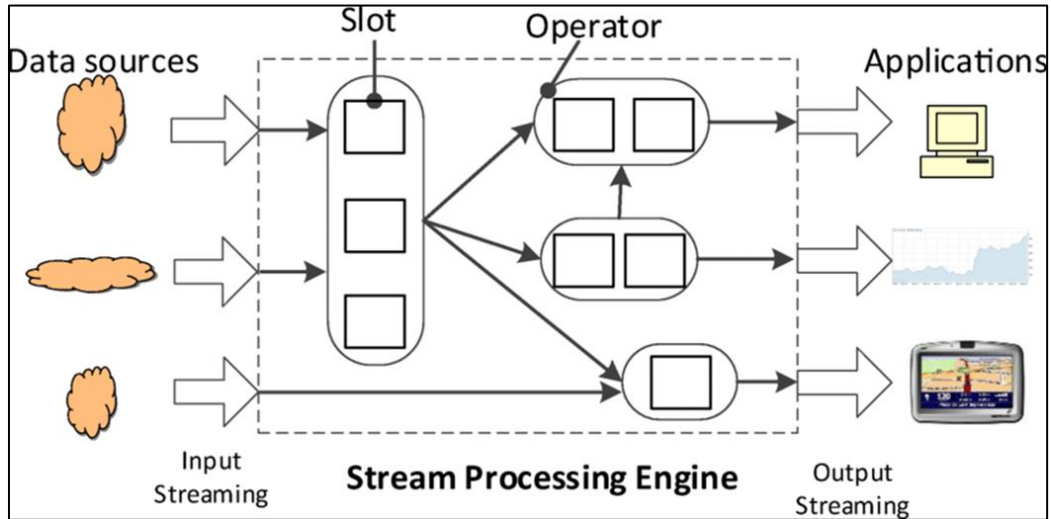
## 2.4. Storage and Output Layer

Once real-time data is processed, it needs to be stored, visualized, or acted upon in a meaningful way. The storage and output layer ensures that processed data is persisted, queried, and presented for further analysis. This layer typically includes:

- Real-Time Dashboards: Visualization platforms such as Grafana, Kibana, and Tableau provide interactive dashboards to monitor real-time metrics and trends.
- Event-Driven Alerts and Notifications: Systems can generate real-time alerts via email, SMS, or push notifications when specific patterns or anomalies are detected (e.g., security breaches, system failures).
- Persistent Databases: Databases such as Apache Cassandra, Amazon DynamoDB, and time-series databases like InfluxDB store processed data for future querying and historical analysis.
- Actuator Systems: In IoT and industrial applications, processed data can trigger automated actions such as adjusting thermostat settings, shutting down faulty machinery, or optimizing power distribution in smart grids.

By ensuring efficient data output, this layer bridges the gap between real-time data processing and actionable insights, enabling businesses and industries to make data-driven decisions with minimal latency.

The architecture of a real-time data stream processing system is designed to support high-speed, continuous data analysis while ensuring scalability, reliability, and efficiency. By integrating data sources, ingestion middleware, processing frameworks, and output storage mechanisms, these systems enable organizations to process and act on vast data streams in real time. With advancements in distributed computing and cloud-native architectures, real-time processing capabilities continue to evolve, driving innovation across various industries, including finance, healthcare, IoT, and smart cities.



**Figure 1** Typical Architecture of a Real-Time Stream Processing System

## 3. Performance Metrics and Scalability

Evaluating the efficiency of real-time data stream processing systems requires analyzing various performance metrics that determine the system's responsiveness, reliability, and adaptability. Given the high-speed nature of real-time data streams, these metrics play a crucial role in assessing how well a system can handle dynamic workloads while ensuring minimal processing delays. The key performance indicators for real-time stream processing include throughput, latency, fault tolerance, and scalability[3].
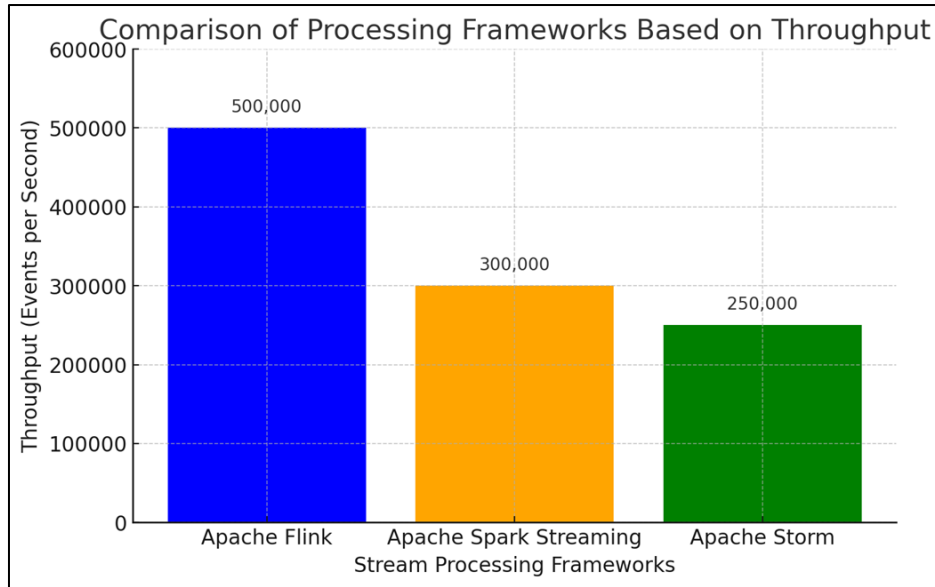
### 3.1. Throughput

Throughput refers to the number of events or messages processed per second by a real-time stream processing system. A higher throughput indicates that the system can handle large volumes of data efficiently without creating processing bottlenecks. Throughput is influenced by several factors, including the efficiency of the processing framework, the level of parallelism, and the underlying hardware or cloud infrastructure.

Modern distributed stream processing frameworks, such as Apache Flink and Apache Spark Streaming, leverage data partitioning and in-memory computation techniques to enhance throughput. The table below compares the throughput capabilities of different frameworks:

**Table 1** Performance Metrics of Stream Processing Frameworks

| Framework | Throughput (Events/sec) | Latency (ms) | Fault Tolerance |
|---|---|---|---|
| Apache Flink | 500K | 5 | High |
| Apache Spark Streaming | 300K | 20 | Medium |
| Apache Storm | 250K | 10 | High |

**Figure 2** Comparison of Processing Frameworks Based on Throughput

### 3.2. Latency

Latency is a critical metric in real-time data processing as it measures the time taken to process and respond to incoming data. Lower latency is essential for applications requiring real-time decision-making, such as fraud detection, network monitoring, and industrial automation.

The processing latency of a stream processing system is affected by factors such as data buffering, event windowing strategies, and network overhead. Apache Flink, for example, is optimized for low-latency processing due to its efficient state management and event-time processing capabilities. Apache Spark Streaming, on the other hand, follows a micro-batching approach, which introduces additional processing delays compared to true stream processing frameworks like Flink and Storm.

### 3.3. Fault Tolerance

Fault tolerance is the system's ability to recover from failures without losing data or disrupting processing tasks. In distributed real-time processing environments, node failures, network disruptions, and hardware crashes are common, making resilience a critical design consideration.

Most modern stream processing frameworks implement checkpointing and state replication mechanisms to enhance fault tolerance:

- Apache Flink ensures fault tolerance by using distributed snapshots and stateful checkpointing, enabling seamless recovery in the event of failures.
- Apache Spark Streaming provides fault tolerance by leveraging Resilient Distributed Datasets (RDDs) to recompute lost data.
- Apache Storm achieves fault tolerance by redistributing failed tasks and ensuring at-least-once or exactly-once message processing guarantees.

A highly fault-tolerant system ensures data consistency, prevents downtime, and minimizes losses in critical applications such as real-time financial trading and cybersecurity monitoring.

### 3.4. Scalability

Scalability defines the ability of a real-time data stream processing system to handle increasing data loads without compromising performance. As data volumes grow, a scalable system should seamlessly allocate additional resources, distribute workloads across multiple nodes, and maintain processing efficiency.

Scalability is achieved through:

- Horizontal Scaling (Scaling Out): Adding more nodes to distribute processing load.
- Vertical Scaling (Scaling Up): Enhancing existing hardware with more memory, CPUs, or storage capacity.
- Elastic Scaling: Dynamic resource allocation based on real-time demand, commonly implemented in cloud-based environments.

Apache Flink and Spark Streaming both support highly scalable architectures by leveraging distributed computing models, allowing them to handle fluctuating workloads efficiently. Apache Storm, while scalable, requires more manual configurations to optimize resource allocation.

## 3.5. Comparative Analysis of Stream Processing Frameworks

The choice of a real-time stream processing framework depends on specific use-case requirements, balancing trade-offs between throughput, latency, fault tolerance, and scalability.

- Apache Flink is ideal for applications requiring high throughput, low latency, and robust fault tolerance, such as real-time analytics and predictive maintenance.
- Apache Spark Streaming is well-suited for scenarios where batch and streaming workloads need to be integrated, such as ETL pipelines and business intelligence applications.
- Apache Storm excels in low-latency event-driven processing, making it effective for network monitoring and real-time alerting systems.

Performance metrics such as throughput, latency, fault tolerance, and scalability play a crucial role in evaluating the efficiency of real-time data stream processing systems. By understanding these factors, organizations can choose the most suitable framework based on their processing needs and workload characteristics. The growing demand for real-time analytics continues to push advancements in stream processing technologies, paving the way for highly efficient, scalable, and resilient data processing architectures.

## 4. Challenges in Large-Scale Systems

While real-time data stream processing offers significant advantages, it also presents several challenges, particularly in large-scale environments where vast amounts of data are generated, transmitted, and processed continuously. These challenges must be addressed to ensure optimal performance, reliability, and security. Some of the key issues include data skew and load balancing, fault recovery, scalability concerns, and security and privacy risks[4].

### 4.1. Data Skew and Load Balancing

One of the major obstacles in large-scale real-time data processing is data skew, where an uneven distribution of data leads to imbalanced workloads across processing nodes. In distributed architectures, certain nodes may receive disproportionately higher data loads than others, causing bottlenecks and reducing overall system efficiency.

Causes of Data Skew:

- Hotspot Data Sources: Certain data sources, such as high-traffic websites or popular IoT devices, generate significantly more data than others.
- Non-Uniform Key Distribution: When partitioning data based on keys (e.g., user ID, region), some keys may have significantly more records than others.
- Varying Processing Complexity: Some data elements require more complex computations, increasing processing times on specific nodes.

Mitigation Strategies:

- Implementing dynamic load balancing mechanisms that redistribute workloads in real-time based on processing capacity.
- Using data replication and partitioning techniques to spread data more evenly across multiple processing nodes.

- Employing adaptive scheduling algorithms that allocate resources based on real-time workload distribution.

## 4.2. Fault Recovery

Ensuring system reliability in the face of hardware failures, software crashes, or network disruptions is a critical challenge in real-time data processing. Since these systems operate continuously, any failure can lead to data loss, inconsistency, or service interruptions.

Key Fault Recovery Challenges:

- Stateful Processing Failures: Many real-time applications, such as financial transaction processing, require maintaining state across data streams. A failure in one node can disrupt dependent computations.
- Message Loss or Duplication: Network failures may lead to lost or duplicated messages, affecting data accuracy.
- Checkpointing Overhead: Frequent state checkpointing improves fault tolerance but can introduce performance overhead.

Mitigation Strategies:

- Checkpointing and Rollback Mechanisms: Frameworks like Apache Flink use distributed snapshots to restore states in case of failures.
- Exactly-Once Processing Guarantees: Techniques such as idempotent processing and transactional updates ensure data consistency.
- Redundant Processing Pipelines: Deploying backup nodes and replication strategies helps maintain continuity during failures.

## 4.3. Scalability Concerns

As the volume of real-time data continues to grow, scalability becomes a significant concern. Large-scale systems must efficiently handle increased data loads while maintaining performance.

Challenges in Scaling Real-Time Systems:

- Horizontal Scaling Complexity: Adding more nodes requires efficient data partitioning and resource orchestration.
- Resource Allocation Overheads: Dynamic allocation of CPU, memory, and storage based on fluctuating workloads can introduce overhead.
- Latency-Sensitive Workloads: Increasing the number of nodes may introduce additional network latencies and synchronization delays.

Mitigation Strategies:

- Elastic Scaling Techniques: Using cloud-based solutions that dynamically allocate resources based on demand.
- Efficient Data Partitioning Strategies: Implementing adaptive sharding techniques to balance workloads.
- Asynchronous Processing Models: Reducing inter-node dependencies by adopting event-driven architectures.

## 4.4. Security and Privacy Challenges

Real-time stream processing systems deal with continuous data ingestion from various sources, making them vulnerable to security threats, data breaches, and privacy concerns.

Key Security and Privacy Risks:

- Unauthorized Data Access: Data streams often include sensitive information such as financial transactions or personal user data.
- Man-in-the-Middle Attacks: Cybercriminals can intercept and manipulate real-time data streams.

- Data Anonymization Challenges: Ensuring privacy in streaming environments is difficult due to the continuous nature of data generation.
- Compliance with Regulations: Adhering to laws such as GDPR and CCPA requires strict data protection measures.

Mitigation Strategies:

- End-to-End Encryption: Encrypting data both in transit and at rest to prevent unauthorized access.
- Access Control Mechanisms: Implementing authentication and authorization protocols to restrict data access.
- Anomaly Detection Systems: Deploying AI-driven monitoring tools to detect suspicious activities in real time.
- Data Masking and Anonymization: Ensuring personally identifiable information (PII) is protected before processing.

Despite the transformative capabilities of real-time data stream processing, large-scale deployments face considerable challenges related to data distribution, fault recovery, scalability, and security. Addressing these issues requires advanced load-balancing techniques, fault-tolerant architectures, dynamic resource allocation, and robust security frameworks. By overcoming these challenges, organizations can unlock the full potential of real-time analytics and ensure reliable, efficient, and secure data processing.

## 5. Challenges in Large-Scale Systems

While real-time data stream processing offers significant advantages, it also presents several challenges, particularly in large-scale environments where vast amounts of data are generated, transmitted, and processed continuously. These challenges must be addressed to ensure optimal performance, reliability, and security. Some of the key issues include data skew and load balancing, fault recovery, scalability concerns, and security and privacy risks[5].

### 5.1. Data Skew and Load Balancing

One of the major obstacles in large-scale real-time data processing is data skew, where an uneven distribution of data leads to imbalanced workloads across processing nodes. In distributed architectures, certain nodes may receive disproportionately higher data loads than others, causing bottlenecks and reducing overall system efficiency.

Causes of Data Skew:

- Hotspot Data Sources: Certain data sources, such as high-traffic websites or popular IoT devices, generate significantly more data than others.
- Non-Uniform Key Distribution: When partitioning data based on keys (e.g., user ID, region), some keys may have significantly more records than others.
- Varying Processing Complexity: Some data elements require more complex computations, increasing processing times on specific nodes.

Mitigation Strategies:

- Implementing dynamic load balancing mechanisms that redistribute workloads in real-time based on processing capacity.
- Using data replication and partitioning techniques to spread data more evenly across multiple processing nodes.
- Employing adaptive scheduling algorithms that allocate resources based on real-time workload distribution.

### 5.2. Fault Recovery

Ensuring system reliability in the face of hardware failures, software crashes, or network disruptions is a critical challenge in real-time data processing. Since these systems operate continuously, any failure can lead to data loss, inconsistency, or service interruptions.

Key Fault Recovery Challenges:

- Stateful Processing Failures: Many real-time applications, such as financial transaction processing, require maintaining state across data streams. A failure in one node can disrupt dependent computations.
- Message Loss or Duplication: Network failures may lead to lost or duplicated messages, affecting data accuracy.
- Checkpointing Overhead: Frequent state checkpointing improves fault tolerance but can introduce performance overhead.

Mitigation Strategies:

- Checkpointing and Rollback Mechanisms: Frameworks like Apache Flink use distributed snapshots to restore states in case of failures.
- Exactly-Once Processing Guarantees: Techniques such as idempotent processing and transactional updates ensure data consistency.
- Redundant Processing Pipelines: Deploying backup nodes and replication strategies helps maintain continuity during failures.

## 5.3. Scalability Concerns

As the volume of real-time data continues to grow, scalability becomes a significant concern. Large-scale systems must efficiently handle increased data loads while maintaining performance.

Challenges in Scaling Real-Time Systems:

- Horizontal Scaling Complexity: Adding more nodes requires efficient data partitioning and resource orchestration.
- Resource Allocation Overheads: Dynamic allocation of CPU, memory, and storage based on fluctuating workloads can introduce overhead.
- Latency-Sensitive Workloads: Increasing the number of nodes may introduce additional network latencies and synchronization delays.

Mitigation Strategies:

- Elastic Scaling Techniques: Using cloud-based solutions that dynamically allocate resources based on demand.
- Efficient Data Partitioning Strategies: Implementing adaptive sharding techniques to balance workloads.
- Asynchronous Processing Models: Reducing inter-node dependencies by adopting event-driven architectures.

## 5.4. Security and Privacy Challenges

Real-time stream processing systems deal with continuous data ingestion from various sources, making them vulnerable to security threats, data breaches, and privacy concerns.

Key Security and Privacy Risks:

- Unauthorized Data Access: Data streams often include sensitive information such as financial transactions or personal user data.
- Man-in-the-Middle Attacks: Cybercriminals can intercept and manipulate real-time data streams.
- Data Anonymization Challenges: Ensuring privacy in streaming environments is difficult due to the continuous nature of data generation.
- Compliance with Regulations: Adhering to laws such as GDPR and CCPA requires strict data protection measures.

Mitigation Strategies:

- End-to-End Encryption: Encrypting data both in transit and at rest to prevent unauthorized access.
- Access Control Mechanisms: Implementing authentication and authorization protocols to restrict data access.
- Anomaly Detection Systems: Deploying AI-driven monitoring tools to detect suspicious activities in real time.

- Data Masking and Anonymization: Ensuring personally identifiable information (PII) is protected before processing.

Despite the transformative capabilities of real-time data stream processing, large-scale deployments face considerable challenges related to data distribution, fault recovery, scalability, and security. Addressing these issues requires advanced load-balancing techniques, fault-tolerant architectures, dynamic resource allocation, and robust security frameworks. By overcoming these challenges, organizations can unlock the full potential of real-time analytics and ensure reliable, efficient, and secure data processing.

## 6. Conclusion

Real-time data stream processing continues to evolve, driven by the growing demand for high-speed data analytics and low-latency decision-making. As organizations generate massive volumes of data from IoT devices, social media, and enterprise applications, new trends and technologies are emerging to enhance the scalability, efficiency, and intelligence of real-time processing systems. These advancements include AI-driven analytics, serverless computing, edge computing, blockchain security, and hybrid cloud architectures. Each of these developments contributes to making real-time data processing more responsive, cost-effective, and secure. One of the most significant advancements in real-time data processing is AI-driven analytics, where artificial intelligence and machine learning are integrated into streaming systems to enable automated insights. Traditional rule-based processing methods are being replaced by self-learning models capable of adapting to evolving data patterns. AI-powered predictive analytics is helping businesses detect anomalies, forecast trends, and make proactive decisions in various domains, such as financial fraud detection, healthcare diagnostics, and smart city traffic management. Additionally, natural language processing (NLP) techniques are being used in real-time sentiment analysis for social media monitoring and chatbot interactions. Another transformative trend is the adoption of serverless computing in stream processing. Serverless architectures eliminate the need for manual infrastructure management by automatically allocating computational resources based on workload demands. This model significantly improves cost-efficiency since resources are only consumed when necessary. Popular cloud platforms like AWS Lambda, Google Cloud Functions, and Azure Functions now support real-time stream processing by integrating with frameworks such as Apache Flink and Kafka Streams. The serverless approach simplifies deployment, reduces operational overhead, and enables event-driven execution, making it a preferred choice for organizations looking to scale their real-time applications effortlessly. With the proliferation of IoT devices and 5G networks, edge computing is becoming a critical component of real-time data processing. Instead of transmitting data to centralized cloud servers, edge computing allows processing to occur closer to the data source, reducing network latency and bandwidth usage. This is especially crucial for applications such as autonomous vehicles, industrial automation, and remote healthcare monitoring, where immediate responses are required. Edge processing not only accelerates decision-making but also enhances security by keeping sensitive data within local networks. As IoT ecosystems expand, edge computing is expected to play a pivotal role in distributed real-time analytics.

Security remains a major challenge in real-time data stream processing, and blockchain technology is being explored as a solution for ensuring data integrity and transparency. Blockchain can be used to create tamper-proof records of streaming data, making it ideal for financial transactions, supply chain management, and regulatory compliance. By leveraging decentralized validation, organizations can enhance trust in real-time data processing systems. Additionally, smart contracts—self-executing code based on predefined conditions—can automate processes based on incoming data, reducing the need for manual intervention in time-sensitive applications. As enterprises continue to generate vast amounts of streaming data, hybrid cloud architectures are emerging as a flexible solution for managing real-time processing workloads. Hybrid cloud models combine on-premises infrastructure, private clouds, and public cloud services to optimize performance, cost, and security. Organizations can store sensitive data on private clouds while utilizing public cloud resources for large-scale analytics. This approach enhances disaster recovery, fault tolerance, and regulatory compliance while ensuring that real-time data processing remains efficient and scalable.

## Compliance with ethical standards

*Disclosure of conflict of interest*

No conflict of interest to be disclosed.

## Reference

[1] Tönjes, Ralf, P. Barnaghi, M. Ali, A. Mileo, M. Hauswirth, F. Ganz, S. Ganea et al. "Real time iot stream processing and large-scale data analytics for smart city applications." In poster session, European Conference on Networks and Communications, p. 10. sn, 2014.

[2] Tantalaki, Nicoleta, Stavros Souravlas, and Manos Roumeliotis. "A review on big data real-time stream processing and its scheduling techniques." International Journal of Parallel, Emergent and Distributed Systems 35, no. 5 (2020): 571-601.

[3] Gulisano, Vincenzo, Ricardo Jimenez-Peris, Marta Patino-Martinez, and Patrick Valduriez. "Streamcloud: A large scale data streaming system." In 2010 IEEE 30th International Conference on Distributed Computing Systems, pp. 126-137. IEEE, 2010.

[4] Gürcan, Fatih, and Muhammet Berigel. "Real-time processing of big data streams: Lifecycle, tools, tasks, and challenges." In 2018 2nd International Symposium on Multidisciplinary Studies and Innovative Technologies (ISMSIT), pp. 1-6. IEEE, 2018.

[5] Shahrivari, Saeed. "Beyond batch processing: towards real-time and streaming big data." Computers 3, no. 4 (2014): 117-129.