

# Event Driven Data Architecture: Design and Implementation with Kinesis and Spark Streaming

Arjun Mantri

Independent Researcher, Bellevue, USA

Email: mantri.arjun[at]gmail.com

ORCID Number- 0009-0005-7715-0108

**Abstract:** This paper reviews the design and implementation of an event-driven data architecture using Amazon Kinesis and Apache Spark Streaming. The evolution of real-time data processing has enabled organizations to handle and analyze data more dynamically and responsively. Amazon Kinesis is highlighted for its robust data ingestion capabilities, while Apache Spark Streaming is noted for its high-throughput, fault-tolerant stream processing. Integrating these technologies allows the creation of a scalable, low-latency, and fault-tolerant system. The paper explores various case studies to illustrate practical applications and benefits across industries such as OTT streaming services, travel booking platforms, and social media networks. For example, Netflix employs this architecture to personalize content recommendations and monitor service quality, while Expedia uses it for real-time availability and pricing updates. LinkedIn leverages the architecture for monitoring user activities and detecting trends in real-time. Implementation details include setting up Kinesis for real-time data ingestion and configuring Spark Streaming for processing and analytics. The system's scalability is ensured by dynamically adjusting Kinesis shards and Spark executors, while fault tolerance is achieved through data replication and checkpointing mechanisms. The findings demonstrate that integrating Amazon Kinesis and Apache Spark Streaming creates a powerful, event-driven data architecture that significantly enhances operational efficiency and supports advanced analytics. This architecture is crucial for modern data-driven applications, providing organizations with the ability to build scalable, real-time data pipelines that enhance performance and support sophisticated data analysis.

**Keywords:** Real-time data processing, Event-driven architecture, Amazon Kinesis, Apache Spark Streaming, Scalable data pipelines.

## 1. Introduction

The advent of real-time data processing has revolutionized the way organizations handle and analyze data, enabling more responsive and dynamic applications (Zaharia et al., 2012). This paper reviews the design and implementation of an event-driven data architecture using Amazon Kinesis and Apache Spark Streaming. By leveraging Kinesis for data ingestion and Spark Streaming for real-time processing, this architecture supports high scalability, low latency, and robust fault tolerance. The findings from various case studies illustrate the practical applications and benefits of this architecture in industries such as OTT (over-the-top) streaming services, travel booking platforms, and social media networks.

## 2. Event Driven Data Architecture

### a) Overview

An event-driven data architecture processes data as it arrives in real-time, rather than processing it in batches. This approach is particularly useful for applications requiring low latency and high throughput. The architecture typically consists of three main components: data ingestion, real-time processing, and data storage/analytics (AWS, 2020a).

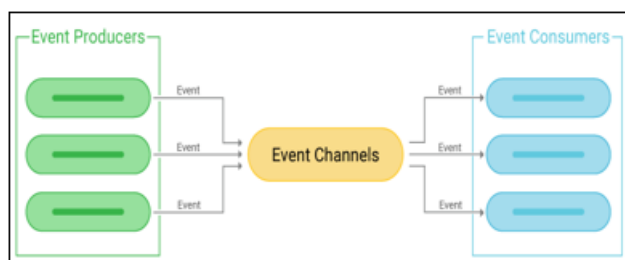


Figure 1: Event-Driven Data Architecture

### b) Data Ingestion with Amazon Kinesis

Amazon Kinesis is a fully managed service for real-time data streaming. It allows the collection, processing, and analysis of real-time data to provide timely insights and react quickly to new information. Kinesis comprises several components:

- **Kinesis Data Streams:** Enables real-time ingestion and processing of large streams of data records (AWS, 2020b).

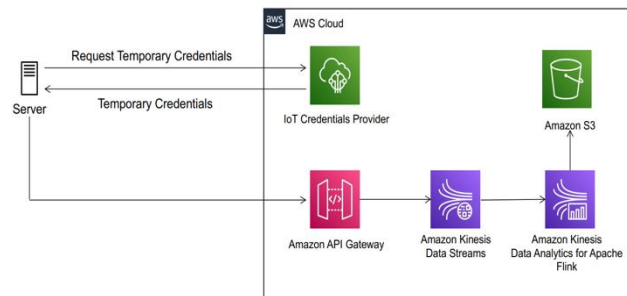


Figure 2: Data Ingestion with Amazon Kinesis

- **Kinesis Data Firehose:** Simplifies the process of loading streaming data into data lakes, data stores, and analytics services (AWS, 2020c).
- **Kinesis Data Analytics:** Allows the running of real-time SQL queries on streaming data (AWS, 2020d).

### c) Real-Time Processing with Apache Spark Streaming



Figure 3: Apache Spark Streaming

Volume 13 Issue 7, July 2024

Fully Refereed | Open Access | Double Blind Peer Reviewed Journal

[www.ijsr.net](http://www.ijsr.net)

Apache Spark Streaming extends the core Spark API to enable scalable, high-throughput, and fault-tolerant stream processing of live data streams (Karau et al., 2015). Spark Streaming can ingest data from various sources like Kafka, Flume, and Kinesis, process it using complex algorithms expressed with high-level functions like map, reduce, join, and window, and push out the processed data to file systems, databases, and live dashboards (Zaharia et al., 2012).

#### d) Scalability and Fault Tolerance

The integration of Kinesis and Spark Streaming ensures that the system can scale to handle high data volumes and maintain low latency. Kinesis can elastically scale the ingestion capacity, while Spark Streaming can scale the processing capacity by adding more nodes to the Spark cluster (Zaharia et al., 2013). Fault tolerance is achieved through Kinesis's data replication and Spark Streaming's support for checkpointing and replaying data (Karau et al., 2015).

### 3. Case Studies

#### a) OTT Streaming Services

In the OTT streaming industry, real-time data processing is crucial for delivering personalized content recommendations and monitoring service quality. For instance, Netflix uses a real-time data architecture to process user interactions and streaming logs. Kinesis ingests this data, and Spark Streaming processes it to generate insights and drive recommendations (Netflix Technology Blog, 2020).

#### b) Travel Booking Platforms

Travel booking platforms leverage real-time data to provide up-to-date availability and pricing information. Companies like Expedia use Kinesis to collect data from various sources, including user searches and booking transactions. Spark Streaming processes this data to update availability and pricing in real-time, enhancing user experience and operational efficiency (Expedia Engineering Blog, 2019).

#### c) Social Media Networks

Social media networks require real-time data processing to monitor user activities, detect trends, and moderate content. For example, LinkedIn uses a similar architecture to process user interactions and engagement metrics in real-time. Kinesis streams the data, and Spark Streaming processes it to generate analytics and support various features like feed ranking and spam detection (LinkedIn Engineering Blog, 2018).

### 4. Implementation Details

#### Setting Up Amazon Kinesis

- 1) **Create a Kinesis Stream:** Define the stream with the required shard capacity based on the expected data volume (AWS, 2020b).
- 2) **Ingest Data:** Use the Kinesis Producer Library (KPL) or the AWS SDK to send data records to the stream (AWS, 2020b).
- 3) **Configure Kinesis Data Firehose:** Set up the delivery stream to transfer data from Kinesis Data Streams to the target destinations such as S3, Redshift, or Elasticsearch (AWS, 2020c).

#### Configuring Apache Spark Streaming

- 1) **Set Up a Spark Cluster:** Deploy a Spark cluster on Amazon EMR or another supported platform (AWS, 2020e).
- 2) **Integrate with Kinesis:** Use the Spark-Kinesis integration library to read data from the Kinesis stream (Karau et al., 2015).
- 3) **Develop Spark Streaming Applications:** Write Spark applications to process the ingested data. This may include transformations, aggregations, and machine learning algorithms (Zaharia et al., 2012).

**Deploy and Monitor:** Deploy the Spark applications and use monitoring tools to ensure the system's performance and reliability (AWS, 2020e).

#### Benefits of the Architecture

##### a) High Scalability

The architecture can handle large volumes of data by scaling Kinesis shards and Spark executors. This elastic scalability ensures that the system can adapt to varying data loads without compromising performance (Zaharia et al., 2013).

##### b) Low Latency

The real-time processing capabilities of Spark Streaming, combined with the low-latency data ingestion of Kinesis, enable the system to deliver timely insights and support dynamic applications (Zaharia et al., 2012).

##### c) Robust Fault Tolerance

Both Kinesis and Spark Streaming offer robust fault-tolerant features. Kinesis replicates data across multiple availability zones, while Spark Streaming supports checkpointing and replay mechanisms to recover from failures (Karau et al., 2015).

##### d) Advanced Analytics

By integrating with Spark, the architecture supports advanced analytics, including machine learning and graph processing. This enables organizations to derive deeper insights from their data and build sophisticated data-driven applications (Zaharia et al., 2013).

### 5. Findings and Conclusion

The integration of Amazon Kinesis and Apache Spark Streaming creates a powerful event-driven data architecture that supports high scalability, low latency, and robust fault tolerance. Case studies from various industries demonstrate that this architecture can significantly enhance operational efficiency and support advanced analytics. Organizations can leverage this architecture to build scalable, real-time data pipelines, enabling more responsive and dynamic applications.

This review paper provides a comprehensive overview of the design and implementation of an event-driven data architecture using Amazon Kinesis and Apache Spark Streaming. The integration of these technologies enables organizations to handle and analyze data in real-time, supporting high scalability, low latency, and robust fault tolerance. The practical applications and benefits of this

architecture are illustrated through various case studies from different industries. The findings demonstrate that integrating Kinesis and Spark Streaming can significantly enhance operational efficiency and support advanced analytics, making it a valuable architecture for modern data-driven applications.

## References

- [1] AWS. (2020a). Real-Time Data Processing with Amazon Kinesis. Retrieved from <https://aws.amazon.com/whitepapers/>
- [2] AWS. (2020b). Kinesis Data Streams. Retrieved from <https://docs.aws.amazon.com/kinesis/>
- [3] AWS. (2020c). Kinesis Data Firehose. Retrieved from <https://docs.aws.amazon.com/kinesis/>
- [4] AWS. (2020d). Kinesis Data Analytics. Retrieved from <https://docs.aws.amazon.com/kinesis/>
- [5] AWS. (2020e). Using Apache Spark with Amazon EMR. Retrieved from <https://docs.aws.amazon.com/emr/latest/ReleaseGuide/emr-spark.html>
- [6] Expedia Engineering Blog. (2019). Real-Time Data Processing at Expedia. Retrieved from <https://techblog.expedia.com/>
- [7] Karau, H., Konwinski, A., Wendell, P., & Zaharia, M. (2015). Learning Spark: Lightning-Fast Big Data Analysis. O'Reilly Media.
- [8] LinkedIn Engineering Blog. (2018). Real-Time Data Processing at LinkedIn. Retrieved from <https://engineering.linkedin.com/>
- [9] Netflix Technology Blog. (2020). Real-Time Data Architecture at Netflix. Retrieved from <https://netflixtechblog.com/>
- [10] Zaharia, M., Das, T., Li, H., et al. (2012). Discretized Streams: An Efficient and Fault-Tolerant Model for Stream Processing on Large Clusters. HotCloud, 2012.
- [11] Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., & Stoica, I. (2013). Spark: Cluster Computing with Working Sets. HotCloud, 2013.