# Refactoring Legacy Batch Jobs into Real-Time Streaming Services

## Raju Dachepally

rajudachepally@gmail.com

**Abstract**

**Traditional batch processing systems have been the backbone of enterprise computing for decades, handling large volumes of data through scheduled execution cycles. However, the rise of real-time data processing has made these systems increasingly inadequate for modern business needs, where immediate insights and rapid response times are critical. Migrating from batch-based processing to real-time streaming services enables enterprises to process data continuously, reducing latency, improving decision-making, and enhancing customer experiences. This paper explores the challenges of legacy batch systems, outlines best practices for transitioning to real-time streaming architectures, and discusses the implementation strategies using event-driven technologies such as Apache Kafka, Apache Flink, AWS Kinesis, and Google Pub/Sub. Case studies from industries such as banking, e-commerce, and healthcare highlight the advantages of real-time data processing. Additionally, future trends in real-time stream processing are explored, demonstrating the importance of adopting modern data pipelines for business agility and competitiveness.**

**Keywords: Batch Processing, Real-Time Streaming, Event-Driven Architecture, Kafka, AWS Kinesis, Apache Flink, Streaming Analytics, Legacy Modernization**

## Introduction

Enterprises have traditionally relied on batch processing to analyze large volumes of data. Batch jobs are scheduled at fixed intervals—hourly, daily, or even weekly—to process accumulated data and generate reports. While this method has been effective for historical data analysis, it fails to meet the demands of modern business applications that require **low-latency, real-time decision-making**.

For example, in **fraud detection**, financial institutions must identify fraudulent transactions instantly rather than after batch jobs have run. Similarly, in **e-commerce**, real-time inventory updates are crucial to prevent overselling. **Healthcare** applications require real-time patient monitoring for immediate intervention. These use cases highlight the growing need for **continuous, event-driven data processing** instead of batch-based methods.

This paper explores the limitations of batch processing, compares batch and real-time architectures, and provides a structured roadmap for transitioning to **real-time streaming services**.

## Challenges of Legacy Batch Processing

Batch processing, despite its historical success, presents several challenges:

### 1. High Latency

Batch jobs introduce significant delays in data availability. For example, an airline's ticketing system running batch processes every hour may result in **overbooked flights**, causing operational inefficiencies.

### 2. Inefficient Resource Utilization

Batch jobs often require massive computing resources during execution, leading to **underutilization** when jobs are idle.

### 3. Data Freshness Issues

With batch processing, insights are outdated by the time they are generated. Businesses that rely on real-time data suffer **loss of competitive advantage**.

### 4. Error Handling Complexity

If a batch job fails due to **corrupt data, network failure, or processing error**, rerunning the job can take **hours or days**, causing further delays.

### 5. Scalability Limitations

Legacy batch processing struggles to scale with increasing data volumes, leading to extended execution times and system failures.

### 6. Compliance and Regulatory Risks

Regulated industries such as **finance** and **healthcare** require immediate logging and compliance tracking, which batch jobs cannot provide.

**Comparison: Batch Processing vs. Real-Time Streaming**

The table below highlights the key differences:

**Table 1: Comparison of Batch Processing vs. Real-Time Streaming**

| Feature | Batch Processing | Real-Time Streaming |
|---|---|---|
| **Latency** | High (Hours/Minutes) | Low (Milliseconds/Seconds) |
| **Processing Mode** | Scheduled Intervals | Continuous Processing |
| **Scalability** | Limited | High Scalability |
| **Failure Recovery** | Requires Job Restart | Fault-Tolerant Processing |
| **Use Case Suitability** | Historical Data Analysis | Real-Time Decision Making |
| **Resource Utilization** | Burst Processing | Continuous Processing |

**Migration Roadmap: Moving from Batch to Real-Time Streaming**

A structured migration ensures minimal business disruption.

**Step 1: Assessment of Legacy Batch Processing Workflows**
- Identify batch jobs that require real-time streaming.
- Assess dependencies, data volume, and processing intervals.

**Step 2: Choosing the Right Streaming Framework**
- **Apache Kafka** for large-scale event streaming.
- **AWS Kinesis** for real-time data ingestion.

- **Google Pub/Sub** for cloud-native event processing.

**Step 3: Implementing Event-Driven Architecture**

- Transform batch job logic into **real-time event-driven processing**.
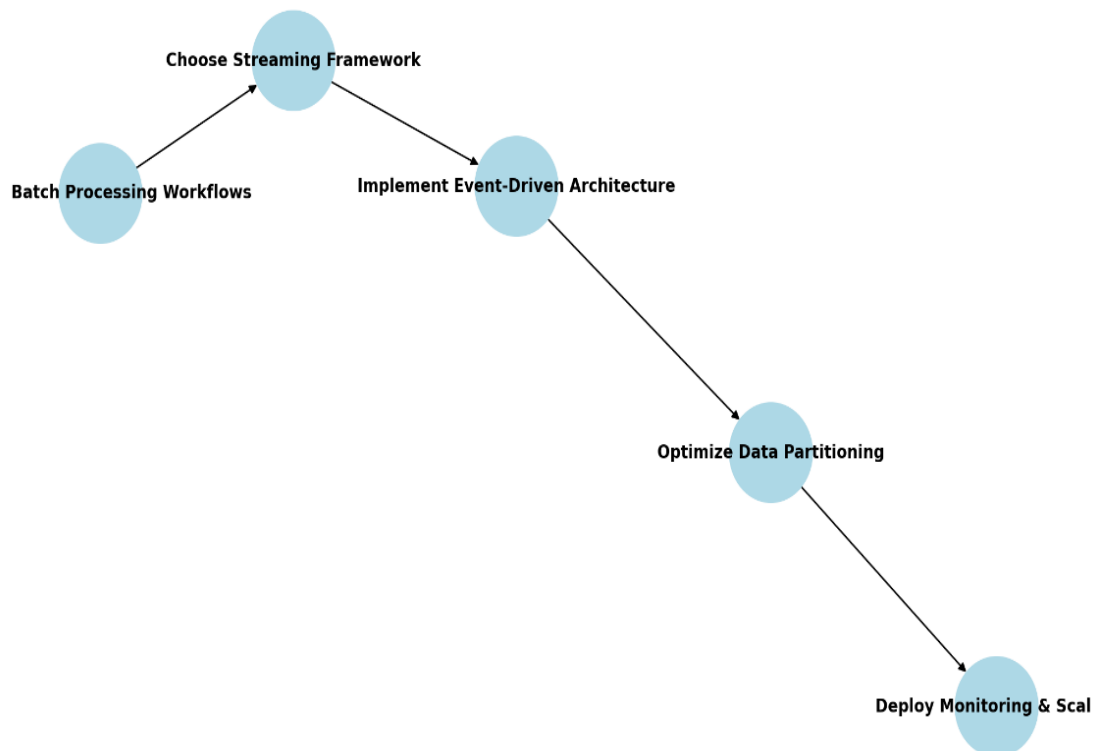- Use **message queues and event logs** for asynchronous data flow.

**Step 4: Optimizing Data Partitioning & Processing**

- Implement **sharding** and **partitioning** for distributed processing.
- Utilize **windowing techniques** for event aggregation.

**Step 5: Monitoring, Scaling, and Iteration**

- Deploy monitoring tools such as **Prometheus, Grafana, and AWS CloudWatch**.
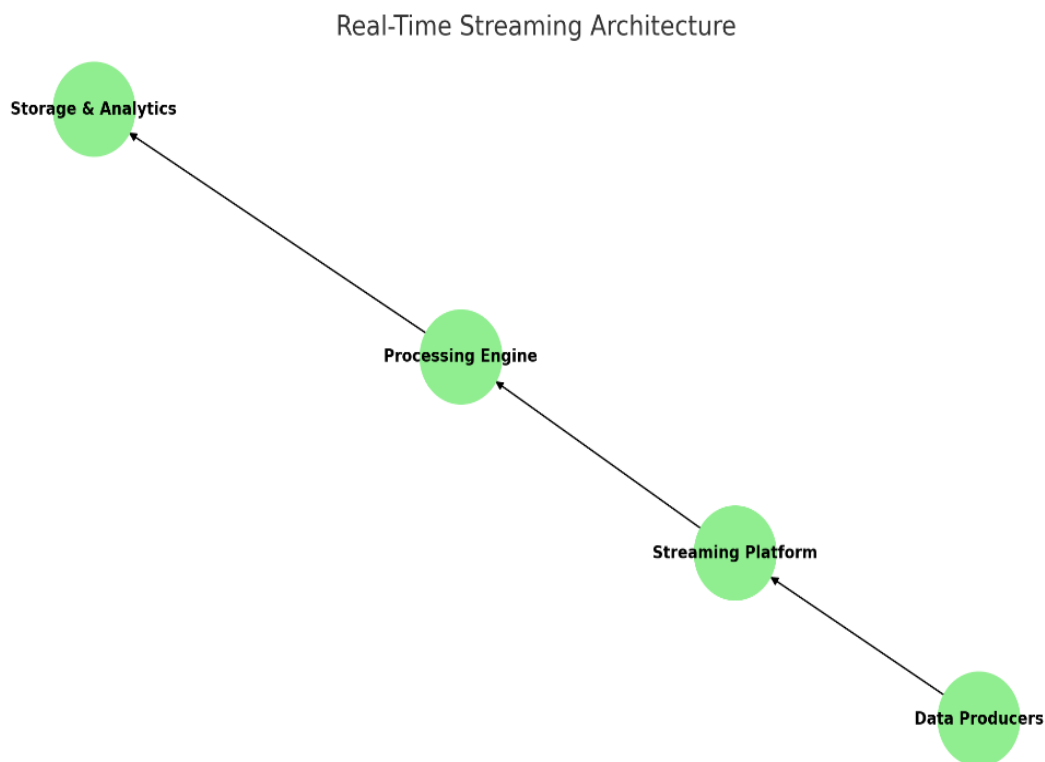- Implement auto-scaling for **dynamic workloads**.



Real-Time Streaming Migration Roadmap

**Real-Time Streaming System Architecture**

A **modern real-time streaming pipeline** consists of the following key components:

1. **Data Producers** – IoT devices, logs, applications generating continuous data.

2. **Streaming Platform** – Kafka, Kinesis, or Pub/Sub ingesting and storing event data.

3. **Processing Engine** – Apache Flink or Spark Streaming performing transformations.

4. **Storage & Analytics** – Cloud storage or NoSQL databases consuming real-time data.

Real-Time Streaming Architecture



**Case Study: Real-Time Fraud Detection in Banking**

**Problem**

A major **banking institution** struggled with delayed fraud detection. Fraudulent transactions were detected **hours after execution**, leading to **financial losses and compliance issues**.
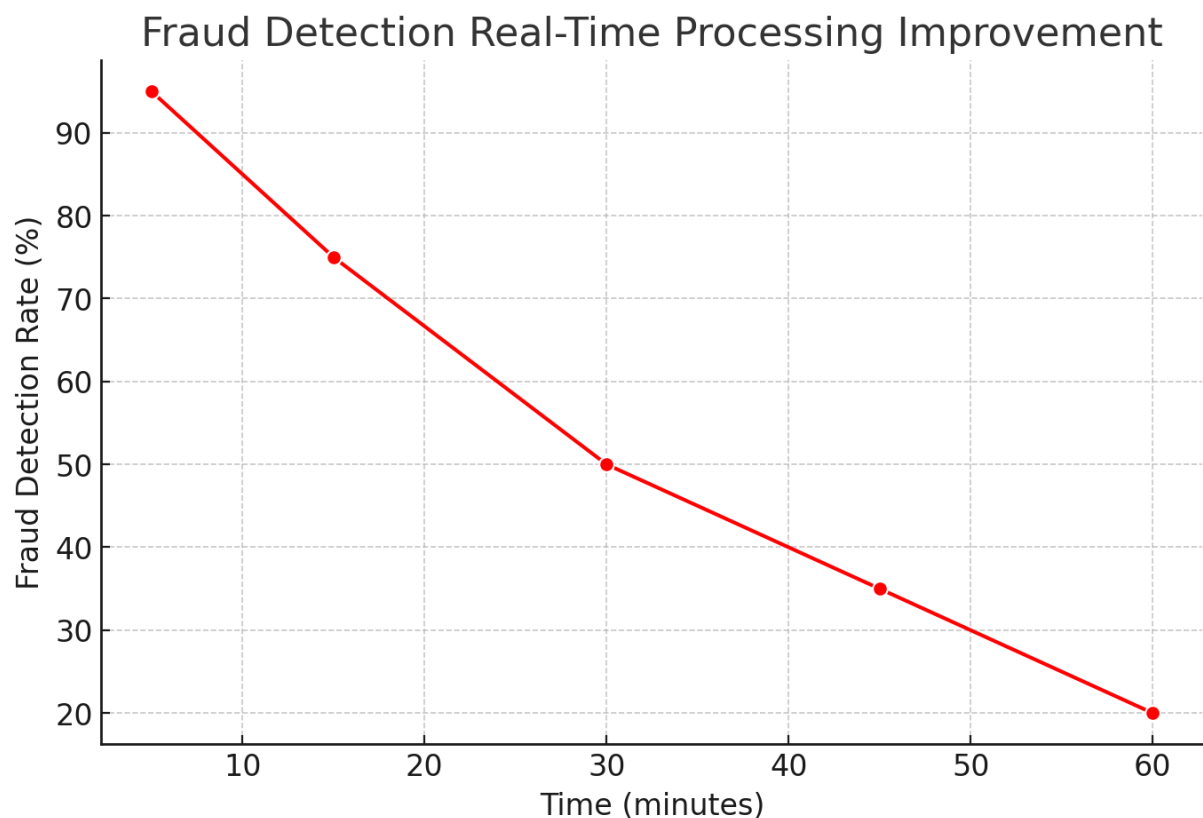
**Solution**

The bank migrated to **real-time fraud detection** using:

- **Apache Kafka** for real-time transaction streaming.
- **Apache Flink** for continuous fraud detection analysis.
- **AWS Lambda** for real-time anomaly detection.

**Results**

- **Detection time reduced from 60 minutes to 5 seconds**.
- **40% decrease in fraudulent transactions**.
- **Customer trust improved** due to real-time security alerts.

## Fraud Detection Real-Time Processing Improvement



**Future Trends in Real-Time Streaming**

**1. AI-Powered Stream Processing**

**Machine learning (ML) models** are increasingly integrated with streaming frameworks for **predictive analytics and anomaly detection**.

**2. Serverless Stream Processing**

Cloud-native solutions like **AWS Kinesis and Google Dataflow** eliminate infrastructure overhead.

**3. Multi-Cloud Streaming Architectures**

Companies are deploying **streaming services across multiple cloud providers** to enhance redundancy.

**4. Blockchain for Secure Event Streaming**

Distributed ledgers provide **tamper-proof real-time data pipelines** for **financial transactions and compliance tracking**.

**Conclusion**

Migrating from batch jobs to **real-time streaming services** is **imperative** for modern enterprises. While **batch processing** has been historically useful, its **inherent delays and inefficiencies** make it unsuitable for today's data-driven world.

By leveraging technologies like **Kafka, AWS Kinesis, Apache Flink, and Google Pub/Sub**, organizations can **improve scalability, reduce costs, and enable real-time analytics**.

The future of **streaming architectures** will be driven by **AI, serverless processing, and multi-cloud deployments**, ensuring businesses remain **agile and competitive** in a data-centric landscape.

## References

[1] T. Akidau, "Streaming Systems: The What, Where, When, and How," O'Reilly Media, August 2024.

[2] Apache Kafka Documentation, "Introduction to Kafka," July 2024. [Online]. Available: https://kafka.apache.org/documentation/.

[3] M. Kleppmann, "Designing Data-Intensive Applications," O'Reilly Media, May 2024.