

Data Engineering Home Assignment - Candivore

Question 1: High-Throughput Streaming Architecture for Large-Scale Events

In this architecture, I am tasked with designing a system for processing large-scale gaming events using a high-throughput, distributed, real-time streaming architecture. I will use Apache Kafka for real-time event ingestion and Apache Spark Streaming for processing. Data will be stored in Snowflake for future analysis. This architecture ensures scalability, fault tolerance, and low-latency processing.

a. Key Components:

1. Event Producers (Game Clients/Servers): These are the systems generating player events like logins, purchases, and gameplay.
2. Message Broker (Apache Kafka): A distributed, fault-tolerant message broker that stores and streams real-time events.
3. Stream Processing Framework (Apache Spark Streaming): Real-time event processing to filter, aggregate, and enrich the data.
4. Data Warehouse (Snowflake): Storing processed event data for analysis and machine learning models.
5. Monitoring and Alerting (Prometheus/Grafana): Monitoring system health, performance, and alerts in real-time.

b. Description of Each Component's Function:

1. Event Producers: These systems (game clients and servers) generate real-time player events and send them to Kafka for ingestion.
2. Apache Kafka: Kafka acts as the message broker, ingesting, storing, and streaming the event data into partitions to enable scalable, parallel processing. It ensures high availability and fault tolerance through partition replication.

3. Apache Spark Streaming: Spark processes the event streams in real time by reading the data from Kafka, performing operations such as filtering, aggregating, and stateful processing.
4. Snowflake: A cloud-based data warehouse that stores the processed data, allowing for advanced analytics, reporting, and machine learning workloads.
5. Prometheus/Grafana: Monitoring tools that track the health and performance of the entire pipeline, with Grafana providing visual dashboards and alerting based on predefined thresholds.

c. Advantages and Potential Drawbacks:

1. Event Producers: High flexibility, but requires optimized serialization (e.g., Avro, Protobuf) to handle large data efficiently.
2. Kafka: High-throughput and scalable, but requires careful configuration and operational overhead at scale (partition management and replication tuning).
3. Spark Streaming: Scalable, low-latency processing, but stateful operations increase complexity and memory usage. Requires tuning for large-scale data.
4. Snowflake: Easily scalable and optimized for analytical queries, but can be costly with large datasets.
5. Prometheus/Grafana: Provides real-time monitoring, but requires proper configuration to avoid alert fatigue.

d. Understanding of Distributed Systems and Real-Time Processing:

This architecture leverages distributed systems concepts, ensuring it can scale horizontally to handle large volumes of data. Apache Kafka handles partitioned, replicated event streaming, ensuring high availability. Apache Spark Streaming distributes the workload across a cluster of nodes, allowing real-time processing with fault tolerance through mechanisms like checkpointing. Snowflake provides distributed query execution for large datasets, while Prometheus and Grafana provide a distributed monitoring solution to ensure system health. Together, these components allow for real-time data ingestion and processing at scale, ensuring low latency, fault tolerance, and high throughput.

