# Streaming data analytics for web application

## Group 4

### October 01, 2025

## Contents

# Streaming data analytics for web application

## Group Information

| Name | Bits ID |
|------|---------|
| Balaji O M | 2024mt03025 |
| Balasubramaniyan | 2024mt03053 |

| Name | Bits ID |
| --- | --- |
| Deep Pokala | 2024mt03042 |
| Jagriti Sharma | 2024mt03116 |

## Abstract

This project implements a comprehensive real-time data streaming pipeline using Apache Flume, Apache Kafka, and Python-based producer and consumer applications. The system demonstrates end-to-end event processing capabilities where a Python producer generates synthetic JSON events, writes them to a log file, which is then tailed by Apache Flume and forwarded to Apache Kafka. A Python consumer subscribes to the Kafka topic and performs real-time analytics on the streaming data.

The architecture showcases modern streaming data processing patterns, containerization with Docker, and resilient service integration with proper error handling and retry mechanisms. The implementation serves as a foundation for understanding distributed streaming systems and can be extended for production-scale data processing scenarios.

## Architecture Diagram

*Figure 1: End-to-end streaming data pipeline architecture showing data flow from producer through Flume to Kafka and consumer processing.*

## Implementation Details

### 1. System Components

### 1.1 Python Producer Service

- **Technology**: Python 3.11 with Faker library
- **Functionality**: Generates synthetic e-commerce events with product information, prices, and timestamps
- **Output**: JSON-formatted events written to `/data/logs/input.log`
- **Configuration**: Configurable production rate via `PRODUCE_RATE_PER_SEC` environment variable

### 1.2 Apache Flume Agent

- **Version**: Apache Flume 1.9.0
- **Source**: Exec source with `tail -F` command to monitor log file
- **Channel**: Memory channel for high-throughput processing
- **Sink**: Kafka sink configured to publish to `events` topic
- **Configuration**: Custom Docker image built from OpenJDK 8 base
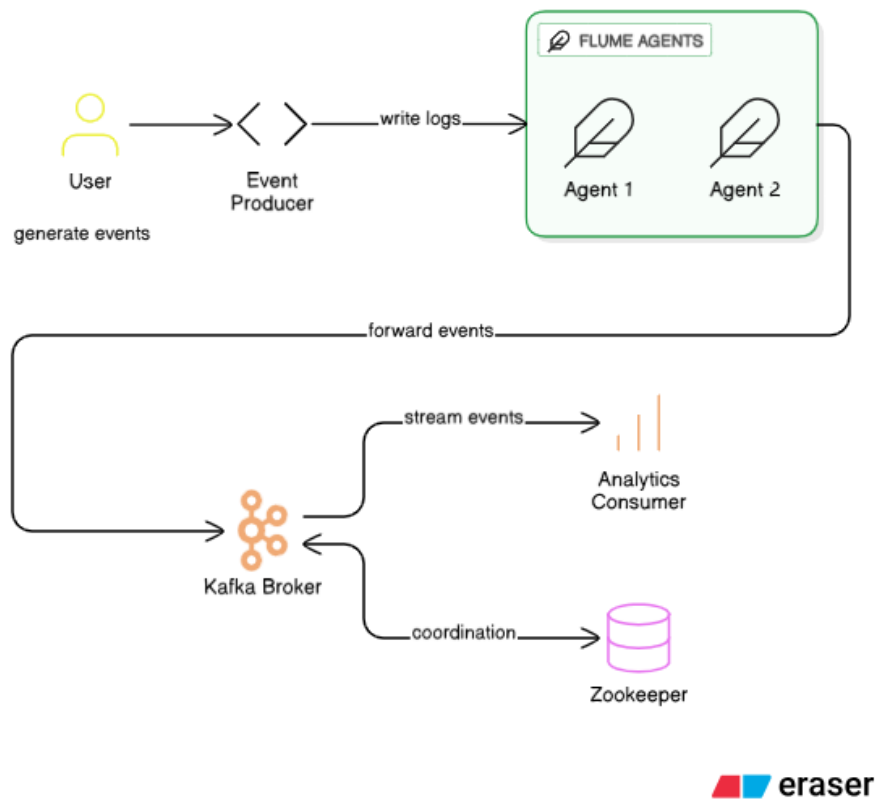
FLUME AGENTS

User

generate events

Event
Producer

write logs

Agent 1    Agent 2

forward events

stream events

Analytics
Consumer

Kafka Broker

coordination

Zookeeper

eraser

Figure 1: Streaming Architecture

**1.3 Apache Kafka Cluster**

- **Version**: Apache Kafka 3.6 with Zookeeper 3.9
- **Topic**: `events` (auto-created with single partition)
- **Configuration**: PLAINTEXT protocol for internal communication
- **Listeners**: Configured for both internal container and external access

**1.4 Python Consumer Service**

- **Technology**: Python 3.11 with kafka-python library
- **Functionality**: Subscribes to Kafka topic and performs real-time analytics
- **Processing**: Price categorization (low/medium/high buckets)
- **Resilience**: Connection retry logic with configurable backoff

**2. Containerization Strategy**

**2.1 Docker Compose Orchestration**

- **Services**: 5 containerized services (Zookeeper, Kafka, Flume, Producer, Consumer)
- **Networking**: Custom Docker network for service discovery
- **Volumes**: Shared volume for log file access between producer and Flume
- **Dependencies**: Proper service startup ordering with `depends_on`
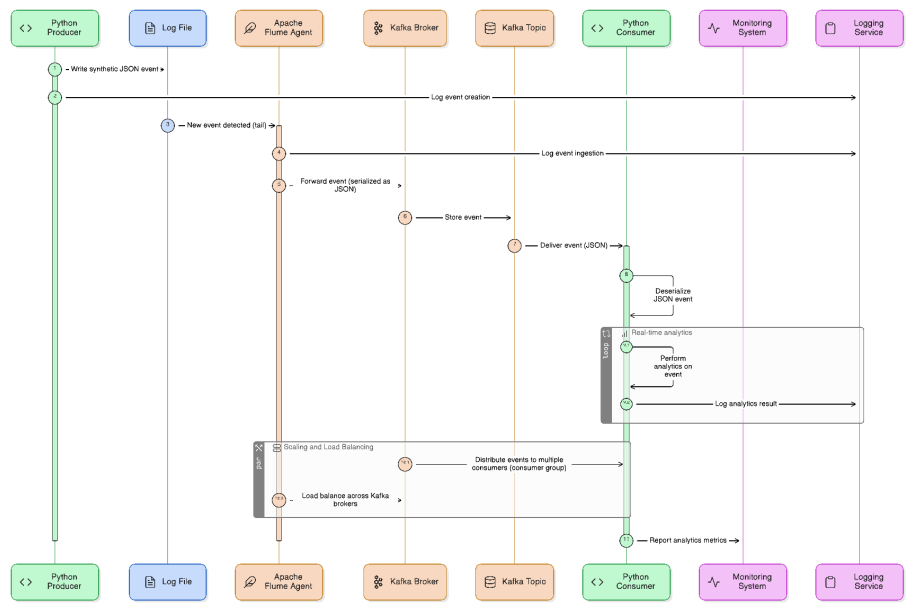
**2.2 Custom Flume Image**

- **Base**: OpenJDK 8 JRE slim
- **Installation**: Automated Flume 1.9.0 download and setup
- **Security**: Non-root user execution with proper permissions
- **Configuration**: Embedded flume.conf for Kafka integration

**3. Data Flow Architecture**

1. **Event Generation**: Producer creates JSON events with product data
2. **File Persistence**: Events written to shared log file via bind mount
3. **Event Ingestion**: Flume tails the file and buffers events in memory channel
4. **Message Publishing**: Flume publishes events to Kafka `events` topic
5. **Event Consumption**: Consumer subscribes and processes events in real-time
6. **Analytics**: Price-based categorization and event metadata extraction

**3.1 Sequence diagram** *Figure 2: Diagram showing data flow from producer through Flume to Kafka and consumer processing.*

Figure 2: Data Flow

## 4. Error Handling and Resilience

- **Consumer Retry Logic**: Configurable connection attempts with exponential backoff
- **Service Dependencies**: Proper startup ordering to prevent connection failures
- **Health Monitoring**: Container health checks and logging integration
- **Graceful Shutdown**: Proper resource cleanup on service termination

## Scripts/Commands for Streaming Integration

### 1. Environment Setup

```
# Clone repository and navigate to project directory
git clone <repository-url>
cd stream-processing-assignment-1

# Create necessary directories
mkdir -p data/logs
touch data/logs/input.log
```

### 2. Service Management

```
# Start all services
make up
# or
docker compose up -d --build

# Stop all services
make down
# or
docker compose down -v

# View service logs
make logs
# or
docker compose logs -f --tail=200

# Check service status
make ps
# or
docker compose ps
```

### 3. Data Management

```
# Reset log file (clean slate)
make clean
```

```
# Create Kafka topic manually (if needed)
make topic
# or
docker exec -it kafka kafka-topics.sh --bootstrap-server localhost:9092 \
  --create --if-not-exists --topic events --replication-factor 1 --partitions 1
```

## 4. Monitoring and Debugging

```
# List Kafka topics
docker exec -it kafka kafka-topics.sh --bootstrap-server localhost:9092 --list

# Consume messages directly from Kafka
docker exec -it kafka kafka-console-consumer.sh \
  --bootstrap-server localhost:9092 --topic events --from-beginning

# Check log file contents
tail -f data/logs/input.log

# Monitor specific service logs
docker compose logs -f producer
docker compose logs -f flume
docker compose logs -f consumer
```

## 5. Integration with External Platforms

### 5.1 Kafka Connect Integration

```
# Example: Connect to external database
# Create connector configuration
cat > kafka-connect-jdbc.json << EOF
{
  "name": "jdbc-sink-connector",
  "config": {
    "connector.class": "io.confluent.connect.jdbc.JdbcSinkConnector",
    "connection.url": "jdbc:postgresql://postgres:5432/streaming_db",
    "topics": "events",
    "auto.create": "true",
    "key.converter": "org.apache.kafka.connect.json.JsonConverter",
    "value.converter": "org.apache.kafka.connect.json.JsonConverter"
  }
}
EOF

# Deploy connector
curl -X POST -H "Content-Type: application/json" \
```

```
    --data @kafka-connect-jdbc.json \
  http://localhost:8083/connectors
```

**5.2 Elasticsearch Integration**

```
# Add Elasticsearch to docker-compose.yml
# Example configuration for logstash pipeline
input {
  kafka {
    bootstrap_servers => "kafka:9092"
    topics => ["events"]
    codec => "json"
  }
}

output {
  elasticsearch {
    hosts => ["elasticsearch:9200"]
    index => "streaming-events-%{+YYYY.MM.dd}"
  }
}
```

**5.3 Real-time Analytics with Apache Spark**

```python
# Spark Streaming job example
from pyspark.sql import SparkSession
from pyspark.sql.functions import *

spark = SparkSession.builder \
    .appName("StreamingAnalytics") \
    .getOrCreate()

df = spark \
    .readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", "localhost:9092") \
    .option("subscribe", "events") \
    .load()

# Process streaming data
processed_df = df.select(
    from_json(col("value").cast("string"), schema).alias("data")
).select("data.*")

# Write to output sink
query = processed_df.writeStream \
```

```
            .outputMode("append") \
            .format("console") \
            .start()
```

## 6. Production Deployment Commands

```
# Scale consumer instances
docker compose up -d --scale consumer=3

# Deploy with external Kafka cluster
export KAFKA_BOOTSTRAP_SERVERS=external-kafka:9092
docker compose up -d producer consumer

# Monitor resource usage
docker stats

# Backup Kafka data
docker exec kafka tar -czf /tmp/kafka-backup.tar.gz /bitnami/kafka/data
```

## 7. Application screenshots

**Start the application**:



Figure 3: Make up

**Check if the services are running using docker ps**:



Figure 4: Docker ps

**Check producer service**:

**Check flume service**:

**Check consumer service**:

Figure 5: Producer service



Figure 6: Flume service



Figure 7: Consumer service

## Conclusion

This streaming data processing assignment successfully demonstrates the implementation of a modern, containerized data pipeline using industry-standard technologies. The project showcases several key concepts:

### Key Achievements

1. **End-to-End Pipeline**: Successfully implemented a complete data flow from event generation to real-time processing
2. **Containerization**: Leveraged Docker and Docker Compose for consistent, reproducible deployments
3. **Service Integration**: Demonstrated proper service orchestration with dependency management
4. **Resilience**: Implemented retry logic and error handling for production-ready applications
5. **Scalability**: Architecture supports horizontal scaling of consumer instances

### Technical Insights

- **Flume Integration**: Successfully configured Flume as a reliable data ingestion layer with Kafka sink
- **Kafka Configuration**: Proper listener configuration for both internal and external access
- **Python Ecosystem**: Leveraged kafka-python for robust consumer implementation
- **Monitoring**: Comprehensive logging and debugging capabilities

### Future Enhancements

1. **Schema Registry**: Implement Avro schemas for data validation and evolution
2. **Stream Processing**: Integrate Apache Flink or Kafka Streams for complex event processing
3. **Metrics**: Add Prometheus/Grafana monitoring for operational visibility
4. **Security**: Implement SASL/SSL authentication and authorization
5. **Data Quality**: Add validation and transformation layers

### Learning Outcomes

This project provided hands-on experience with: - Distributed streaming architectures - Container orchestration and service discovery - Event-driven system design patterns - Real-time data processing concepts - Production deployment considerations

The implementation serves as a solid foundation for understanding modern data streaming platforms and can be extended for various use cases including IoT data processing, real-time analytics, and event-driven microservices architectures.