# Streaming Ingest with Kafka and Spark Streaming
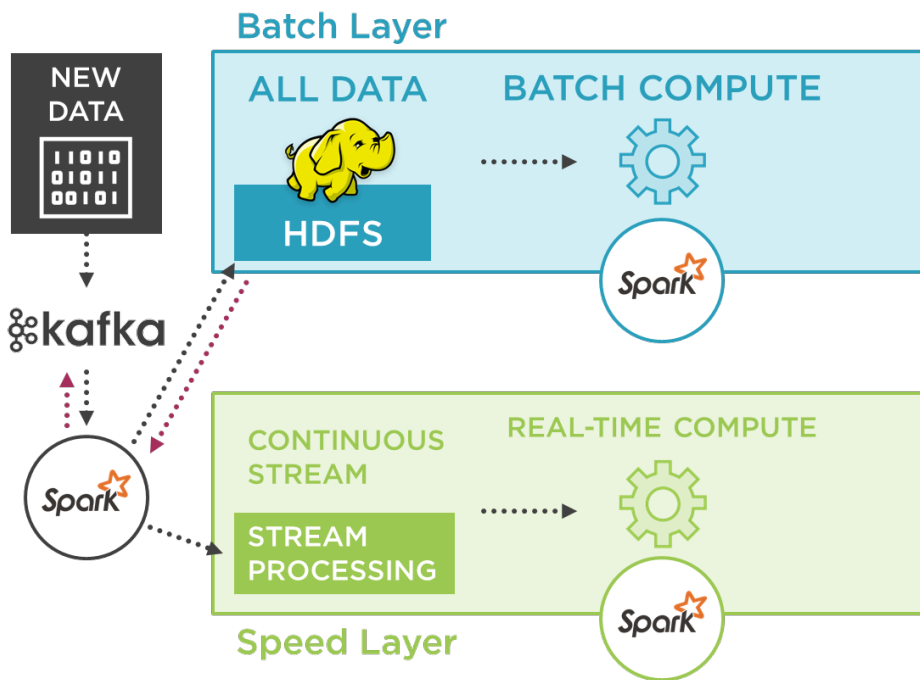
**Ahmad Alkilani**
DATA ARCHITECT

@akizl

# Streaming Ingest with Kafka and Spark Streaming



**Batch Layer**

NEW DATA

ALL DATA    BATCH COMPUTE

HDFS

Spark

CONTINUOUS STREAM    REAL-TIME COMPUTE

STREAM PROCESSING
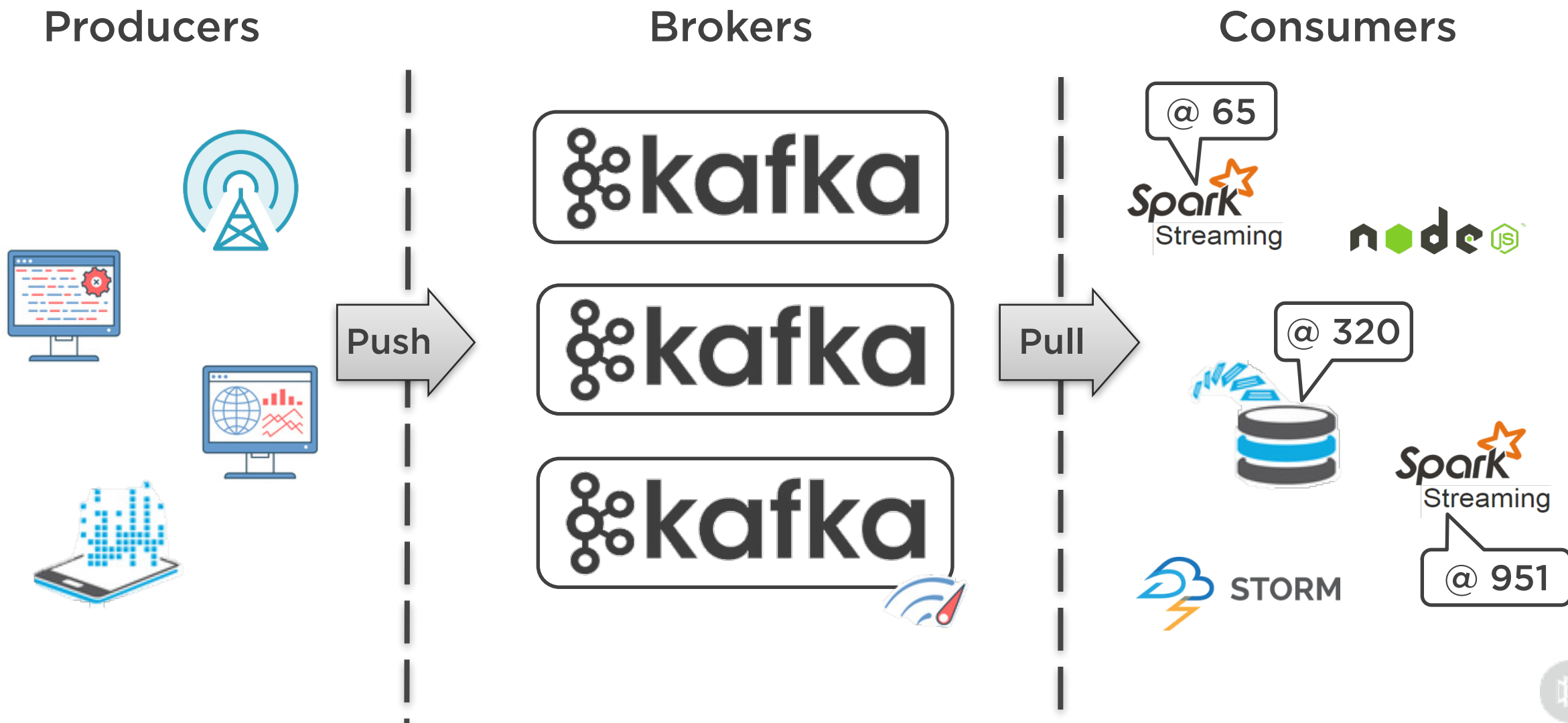
Spark

**Speed Layer**

- **Introduction to Kafka**
  - Architecture
  - Producers and Consumers

- **Create a Kafka Producer**

- **Spark Streaming Integration with Kafka**

- **Integrate Batch and Streaming**
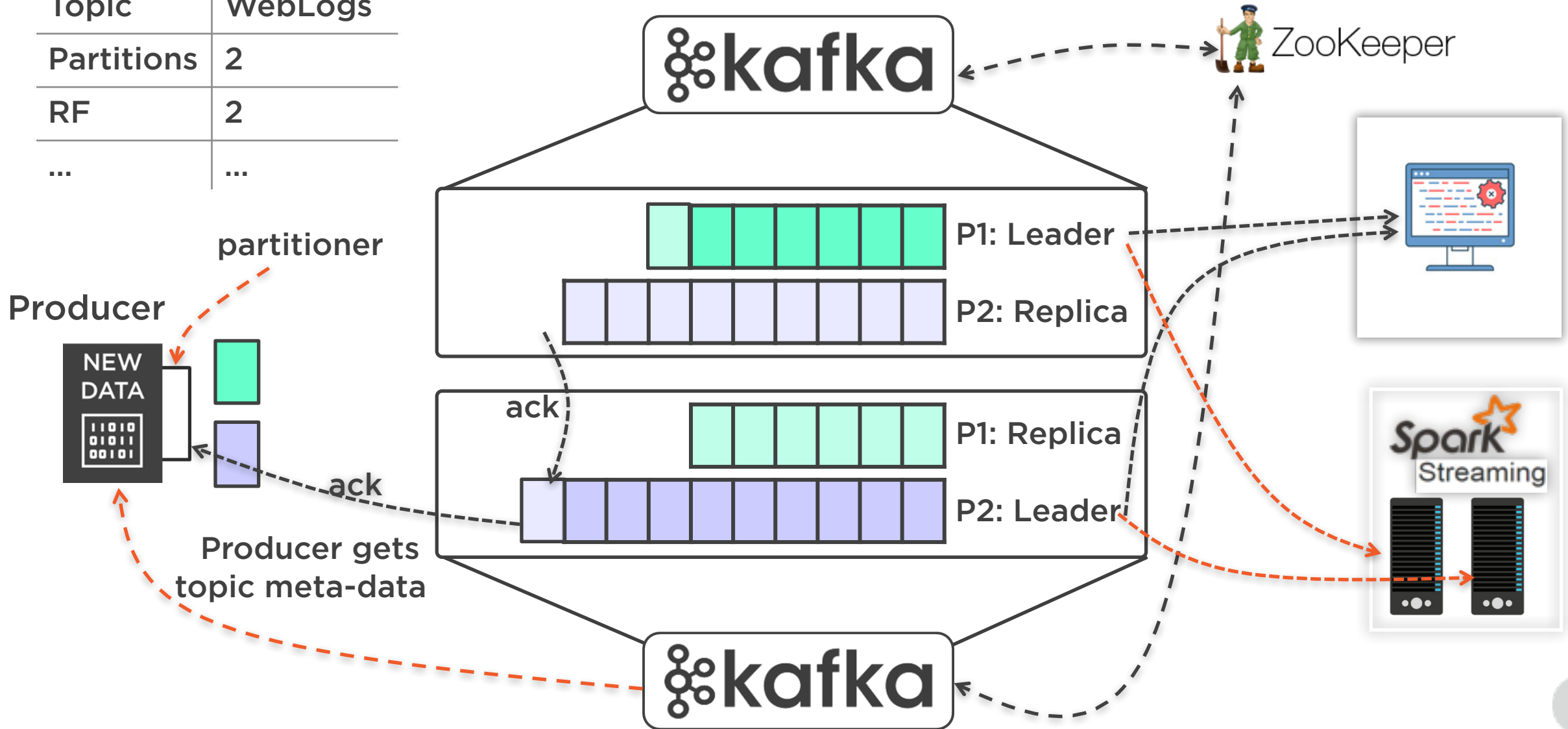
# kafka

Distributed publish-subscribe messaging system

# The Kafka Broker

| Topic | WebLogs |
|-------|---------|
| Partitions | 2 |
| RF | 2 |
| ... | ... |

partitioner

Producer

NEW DATA
1010
0111
00101

Producer gets
topic meta-data

ack

ack

ZooKeeper

P1: Leader

P2: Replica

P1: Replica

P2: Leader

Spark Streaming

# Partition Assignment & Consumers

# Kafka Consumers

## Topics

| weblogs | |
|---|---|
| Partitions | 3 |
| RF | 1 |

| telemetry | |
|---|---|
| Partitions | 3 |
| RF | 2 |

### Broker 1

- P1 (LR)
- P2 (LR)
- P3

### Broker 2

- P1 (LR)
- P2
- P3 (LR)

### Broker 3

- P1
- P2 (LR)
- P3 (LR)

C1

P1,P3  P2

1,P3

**Consumer Group A**

C1  C2  C3  C4

**Consumer Group B**

Consumer Group A

Zookeeper

C1 C2 C3

Partition 1

Partition 2

kafka

Partition 3

C1 C2 C3

Consumer Group B

pluralsight

# Messaging Models

# Messaging Models

## Queue Semantics

| Broker 1 | Broker 2 | Broker 3 |
|----------|----------|----------|
| P1 (LR)  | P3 (LR)  | P2 (LR)  |

**Topics**

weblogs

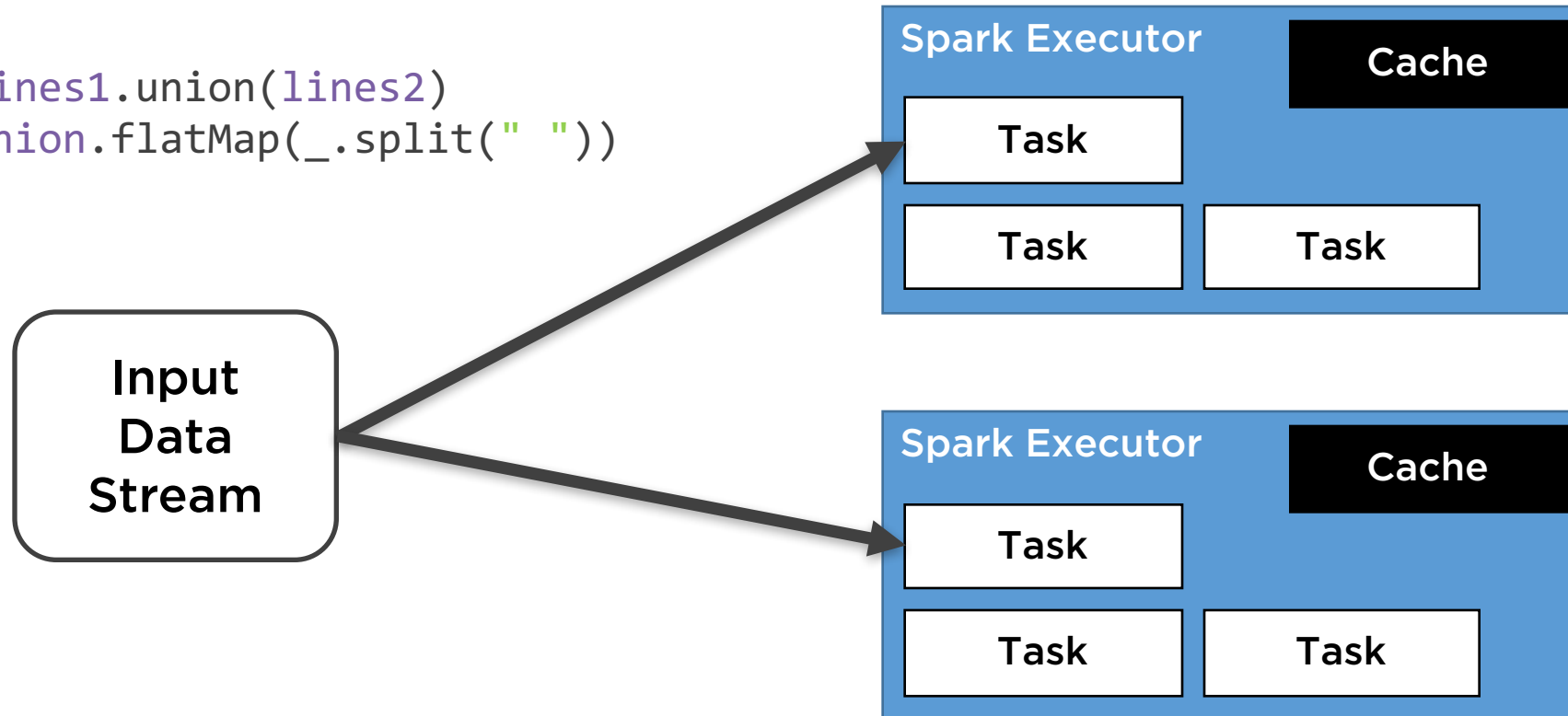| Partitions | 3 |
|------------|---|
| RF         | 1 |

C1  C2  C3

CG: **MyQueue**

# Receiver Model

```scala
val lines1 = ssc.socketTextStream("localhost", 9999)
val lines2 = ssc.socketTextStream("localhost", 9998)

val linesUnion = lines1.union(lines2)
val words = linesUnion.flatMap(_.split(" "))
```

# Spark Kafka Integration

**Kafka Consumer APIs**

## Spark Streaming Kafka Integration

**High-Level API**

### Receiver Approach
- Receivers to receive data
- Data stored in Spark executors
- Zero-data loss requires write-ahead log
- Allows for at-least-once semantics
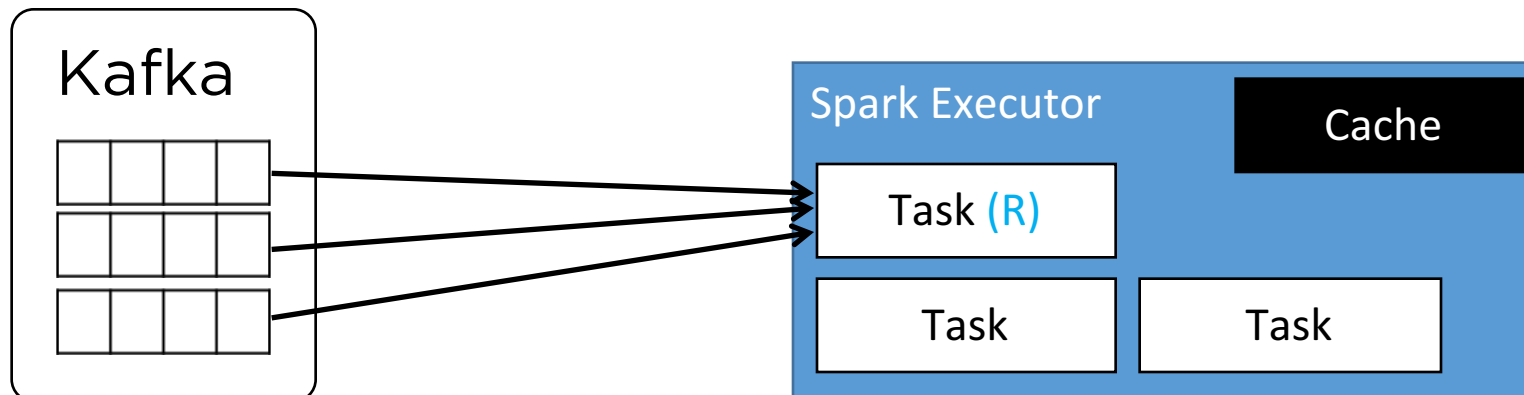
**Simple API**

### Direct Approach
- No receivers. Queries Kafka each batch for offset range
- Simplifies parallelism at the expense of latency
- Zero-data loss without write-ahead log; relies on Kafka's retention to replay messages. Better at processing larger datasets
- Allows for exactly-once semantics

# Receiver-based Approach

## Option 1: Create a single Kafka stream

```scala
val kafkaStream = KafkaUtils.createStream[String, String, StringDecoder, StringDecoder](
  ssc, kafkaParams, Map(topic -> 1), StorageLevel.MEMORY_AND_DISK)
  .map(_._2)
```
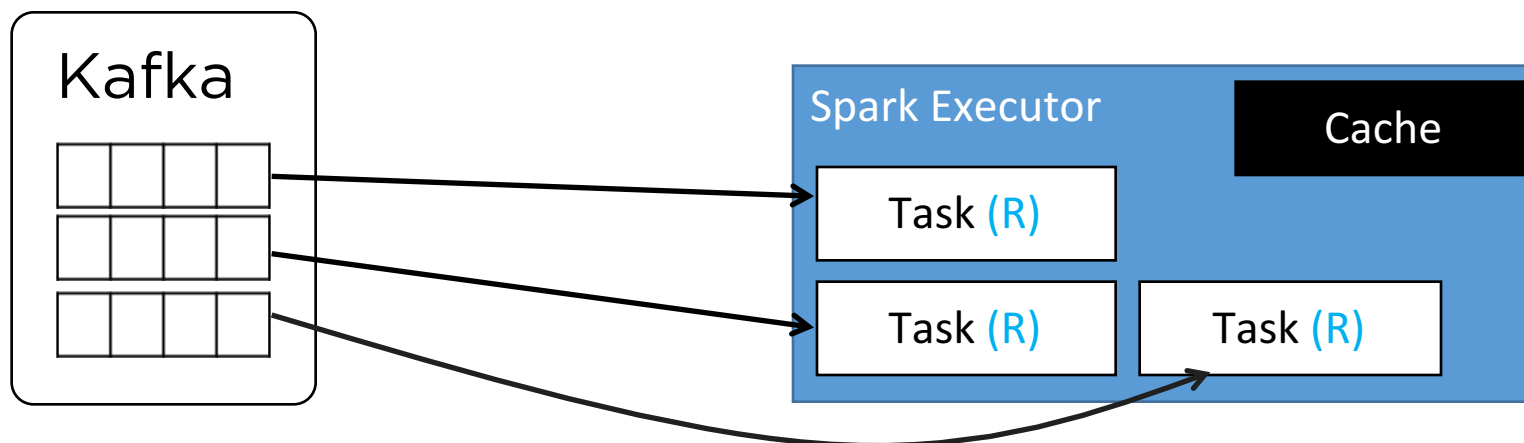
# Receiver-based Approach

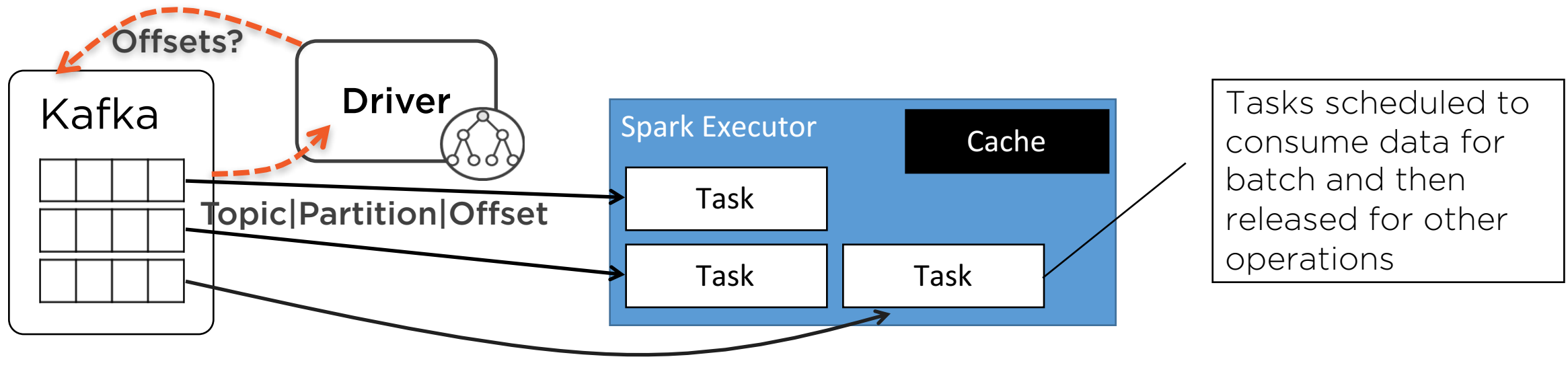## Option 2: Create a Kafka stream per topic-partition

```scala
val receiverCount = 3
val kafkaStreams = (1 to receiverCount).map { _ =>
  KafkaUtils.createStream[String, String, StringDecoder, StringDecoder](
    ssc, kafkaParams, Map(topic -> 1), StorageLevel.MEMORY_AND_DISK)
}
val kafkaStream = ssc.union(kafkaStreams)
  .map(_._2)
```

# Direct Approach

**Driver determines offsets since last batch**

```scala
val params = Map(
  "metadata.broker.list" -> "localhost:9092",
  "group.id" -> "lambda",
  "auto.offset.reset" -> "smallest"
)

KafkaUtils.createDirectStream[String, String, StringDecoder, StringDecoder](ssc, params, Set(topic))
  .map(_._2)
```
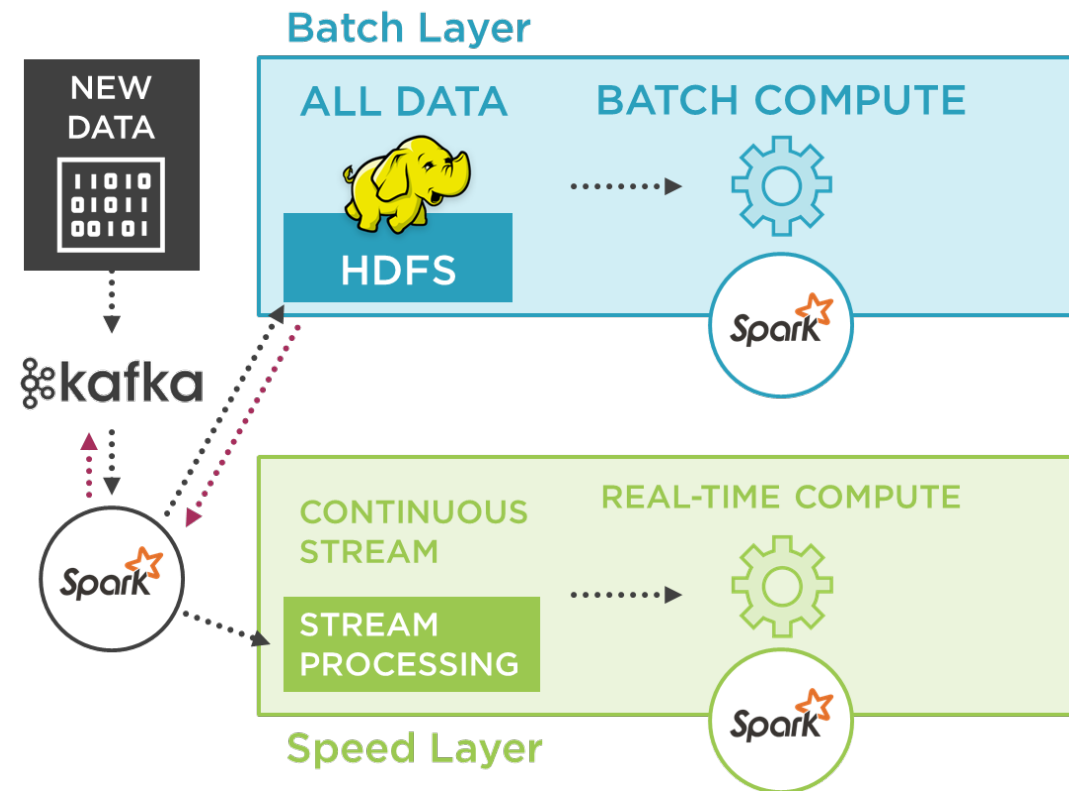
# Demo

## Save Data from Kafka to HDFS

## Build Resiliency into the Application
– Recover from complete failures
– Allow for application updates

# Kafka Direct Stream to HDFS

**HDFS**

../KafkaTopic/KafkaPartition

↳ **data, fromOffset, untilOffset**

Direct Kafka stream means there's a 1-1 mapping between
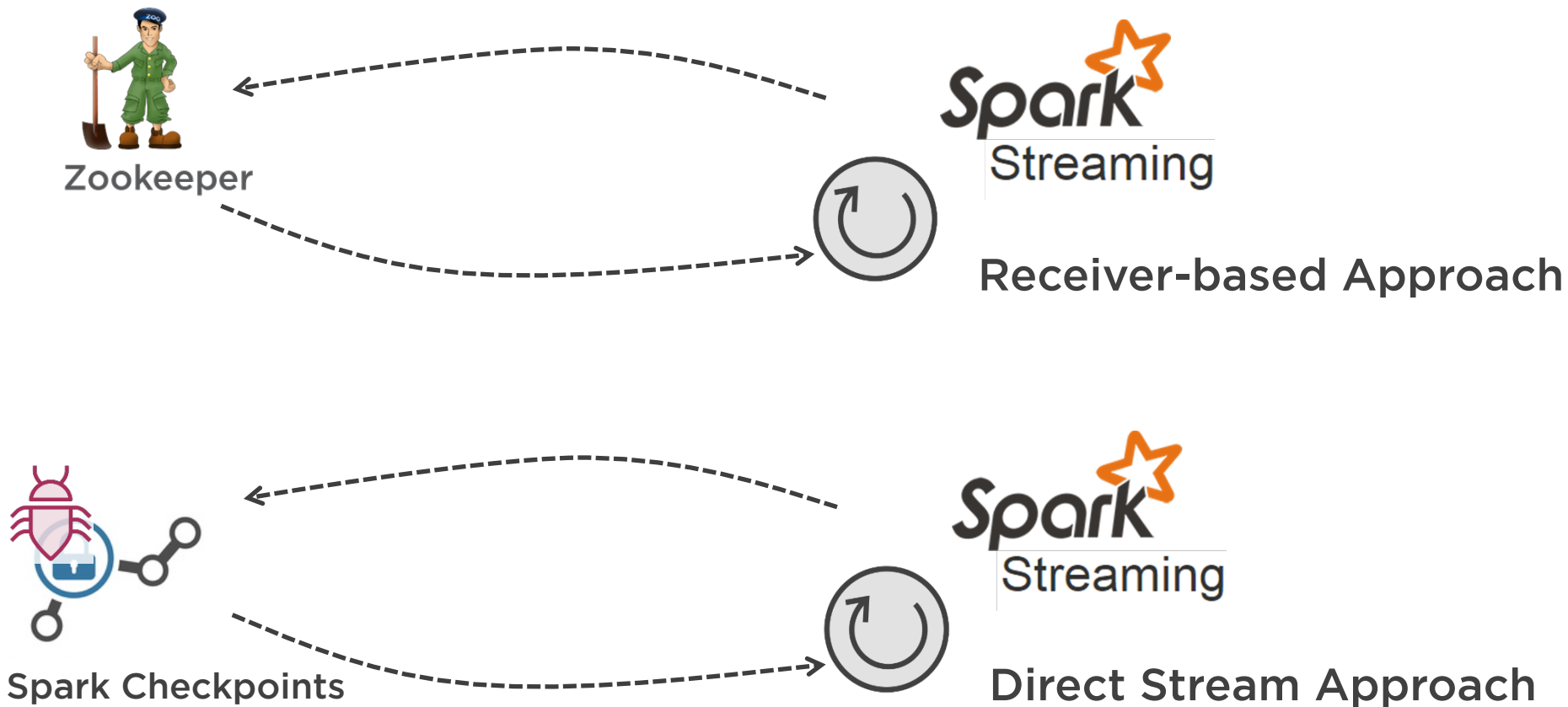Kafka partition and Spark partition

**HasOffsetRanges**

RDD

val offsetRanges = rdd.asInstanceOf[HasOffsetRanges].offsetRanges

offsetRanges(partitionNumber)
      .topic
      .partition
      .fromOffset
      .untilOffset

# Streaming Resiliency



**Zookeeper**

**Receiver-based Approach**

**Spark Checkpoints**

**Direct Stream Approach**

# Summary

- **Apache Kafka**
  - Broker
  - Producer
  - Consumers and Partitions

- **Spark Streaming**
  - Receiver-based
  - Direct Stream

- **Resiliency**
  - Direct Stream Offsets
  - Recover from Upgrades

- **HDFS and Batch Layer Integration**