



# Apache Spark Streaming

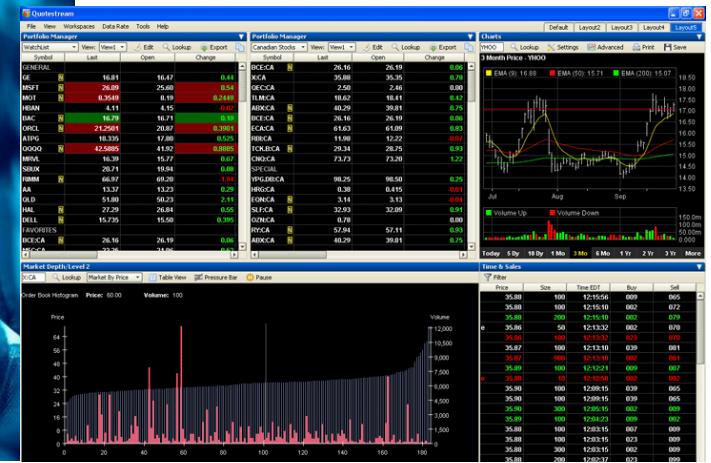
Diego J. Bodas Sagi

- ➊ Real Time Analytics
- ➋ Batch vs Stream
- ➌ Processing Flow
- ➍ Architectures
- ➎ Tools
- ➏ Spark Streaming
- ➐ Some concepts
- ➑ Resumen





# Real Time Analytics



PLANES:

**Minute-to-minute  
monitoring**



**Weather**

Detailed weather sensors are more accurate than the National Weather Service, telling the airline when to prepare de-icers and delay and cancel flights.

**Parts**

Even a five degree temperature variation may indicate a part needs to be replaced.

**Flight Plan**

Keeping an eye on a plane's path from the ground, and alerting pilots of any anomalies.





# Batch vs Stream processing

## Batch Processing

**High** Latency

Batch Processors

Static Files

## Stream Processing

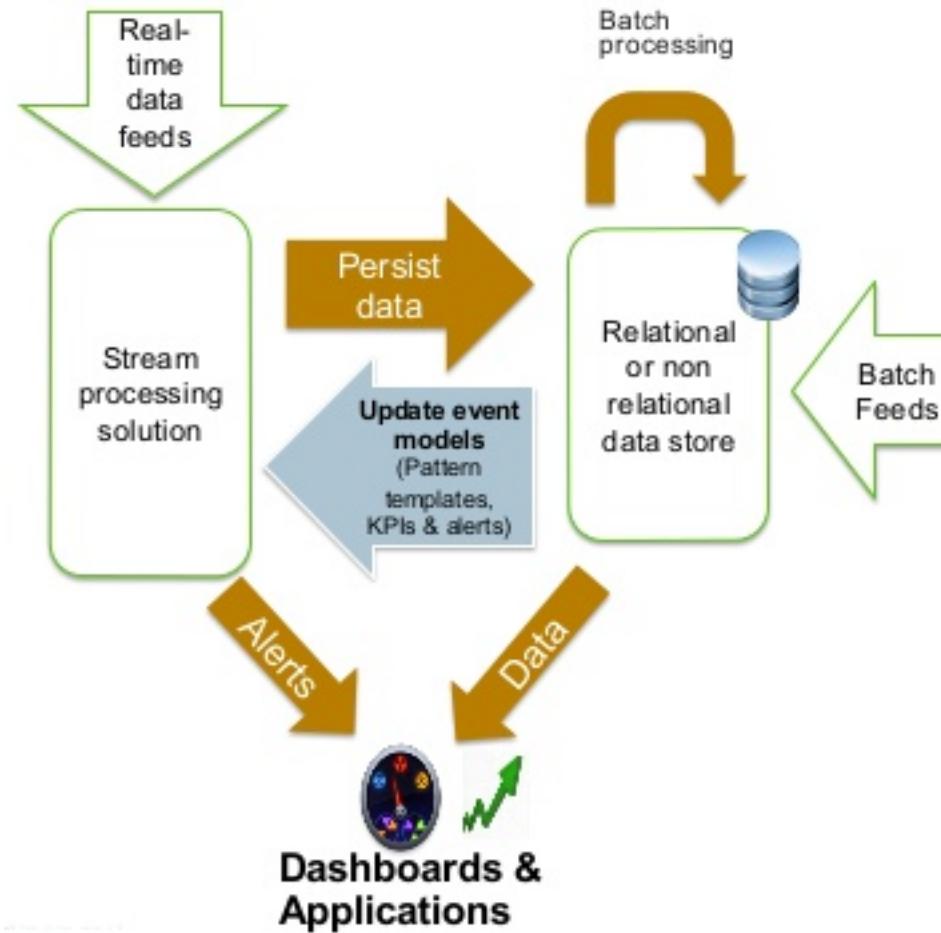
**Low** Latency

Stream Processors

Event Streams



# Processing Flow

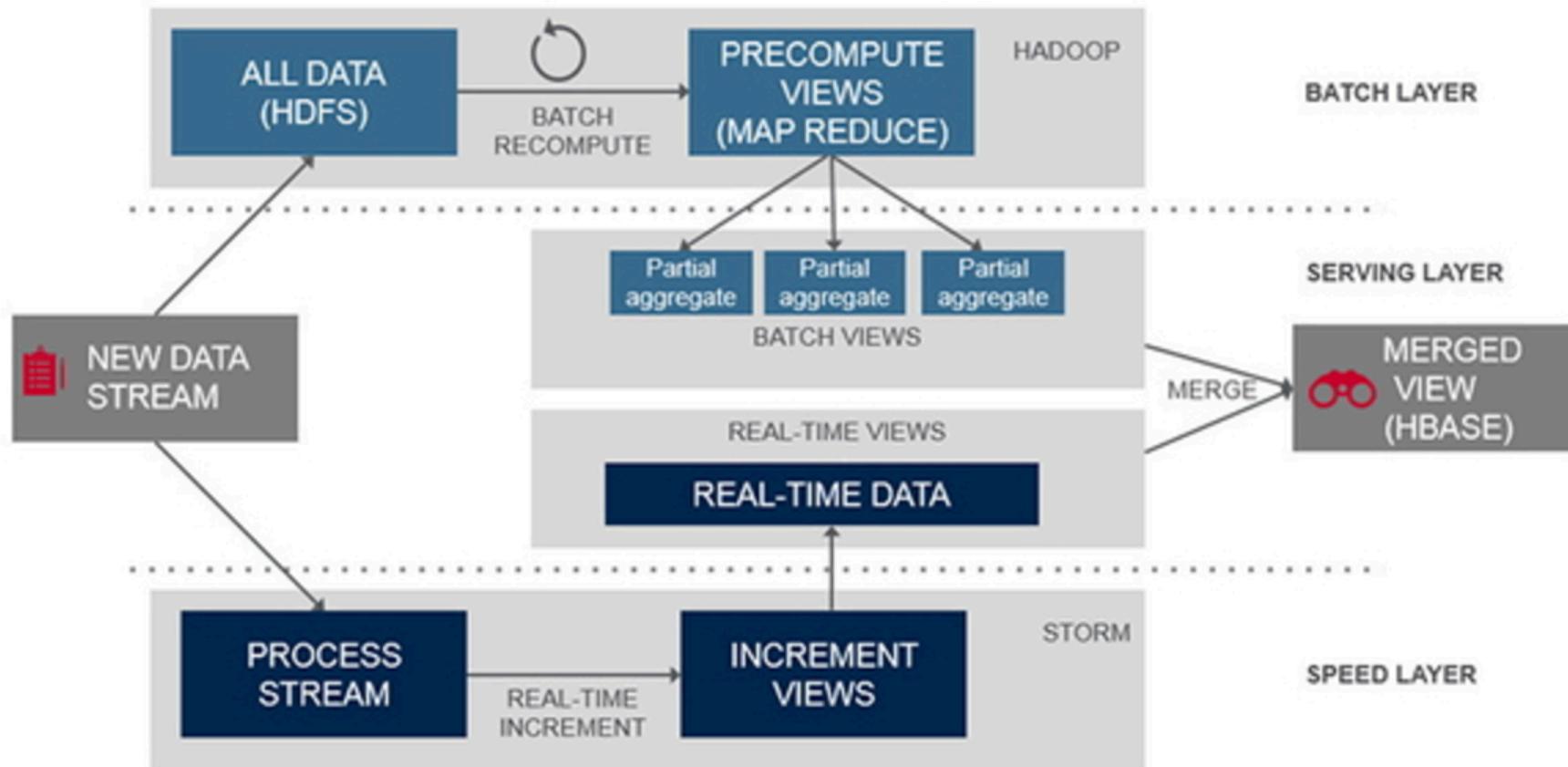


Source: Hortonworks





# Lambda architecture





## Streaming Patterns (by Cloudera)

- Stream ingestion: Involves low-latency persisting of events to HDFS, Apache HBase, and Apache Solr
- Near Real-Time (NRT): Event Processing with External Context: Takes actions like alerting, flagging, transforming, and filtering of events as they arrive. Actions might be taken based on sophisticated criteria, such as anomaly detection models. Common use cases, such as NRT fraud detection and recommendation, often demand low latencies under 100 milliseconds



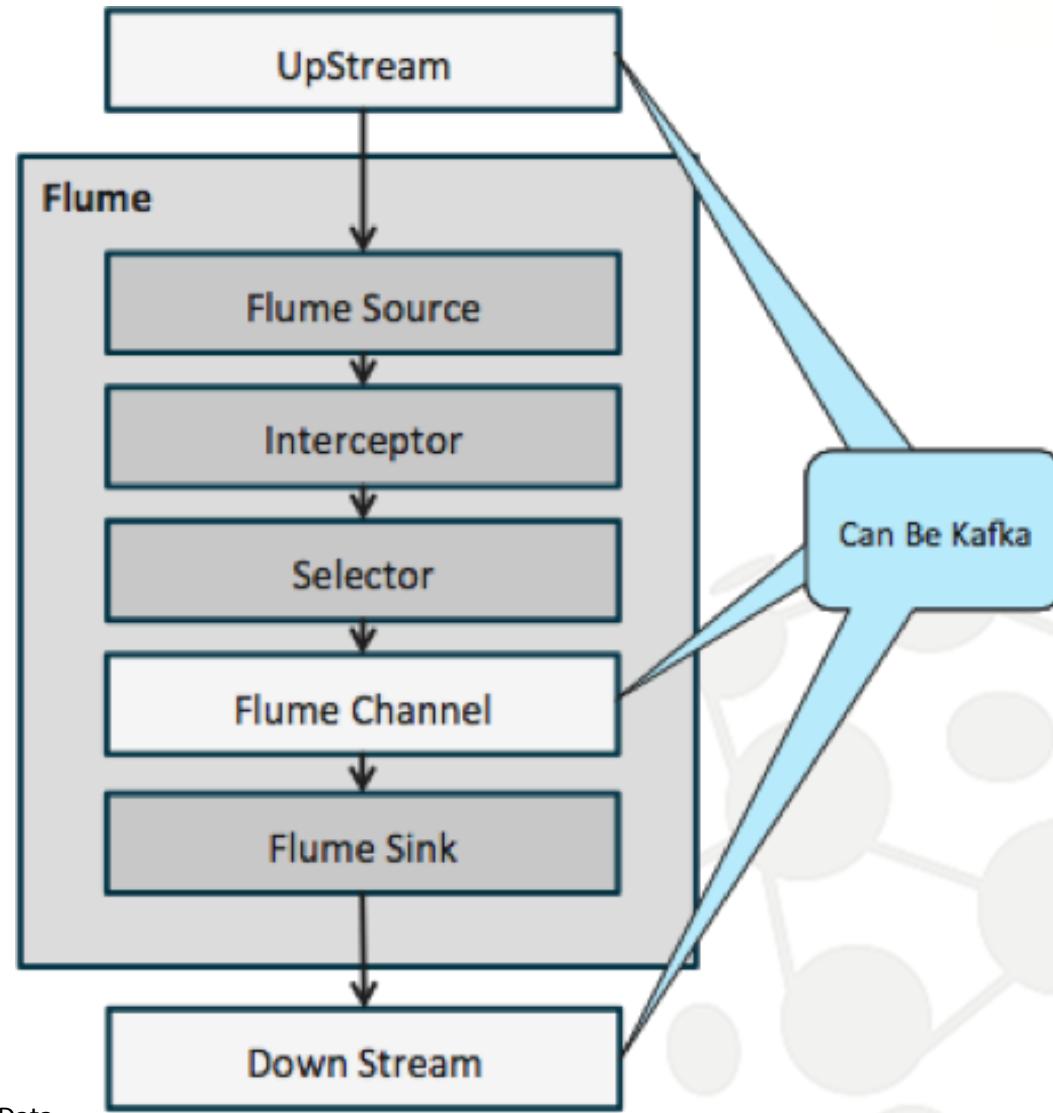


## Streaming Patterns (by Cloudera)

- NRT Event Partitioned Processing: Similar to NRT event processing, but deriving benefits from partitioning the data—like storing more relevant external information in memory. This pattern also requires processing latencies under 100 milliseconds
- Complex Topology for Aggregations or ML: The holy grail of stream processing: gets real-time answers from data with a complex and flexible set of operations. Here, because results often depend on windowed computations and require more active data, the focus shifts from ultra-low latency to functionality and accuracy.

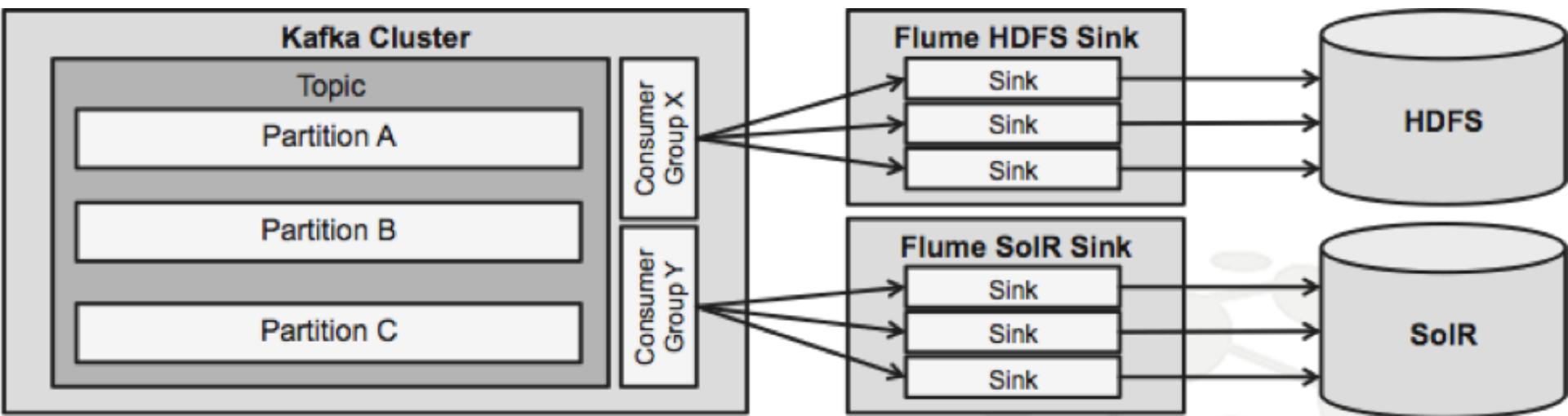


# Streaming Ingestion

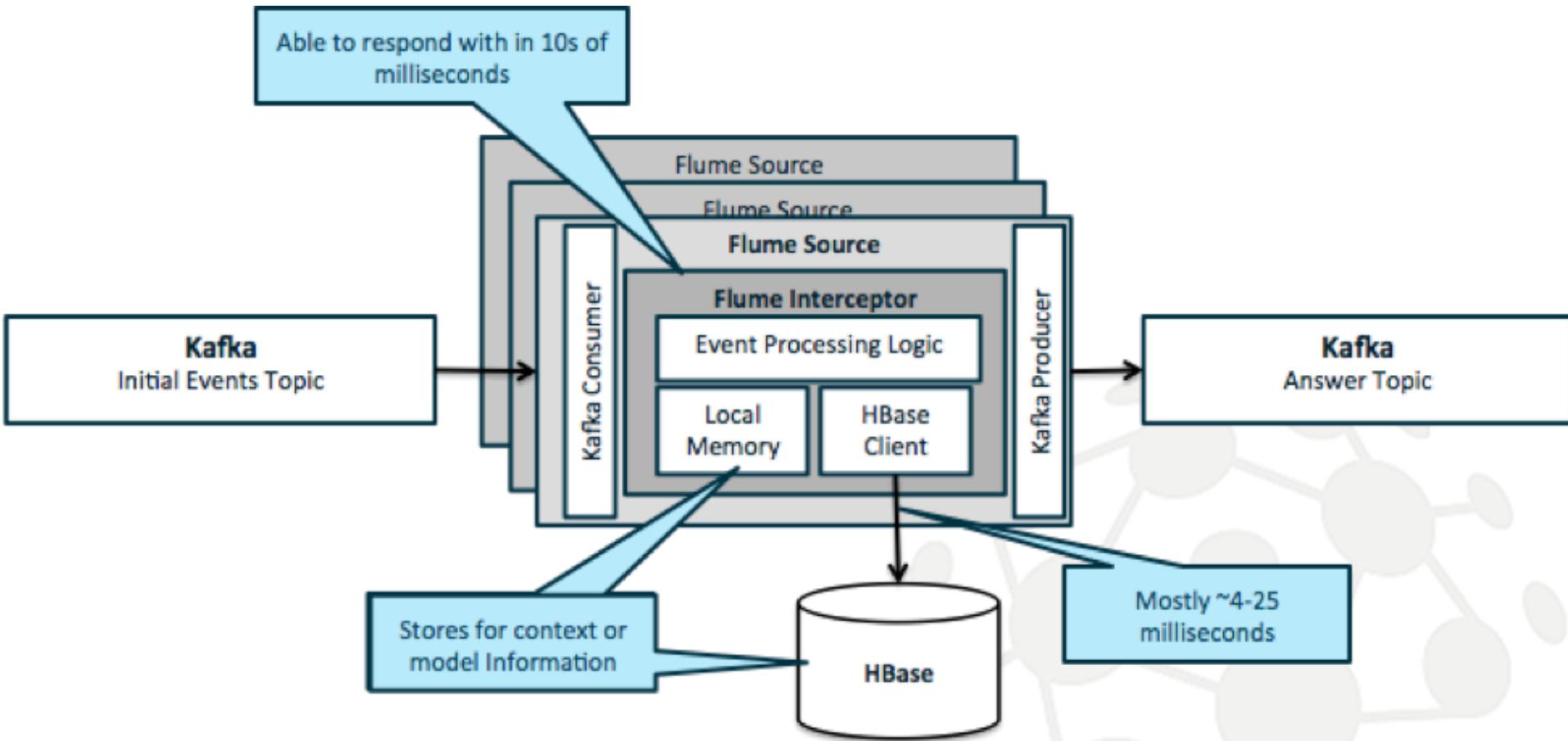




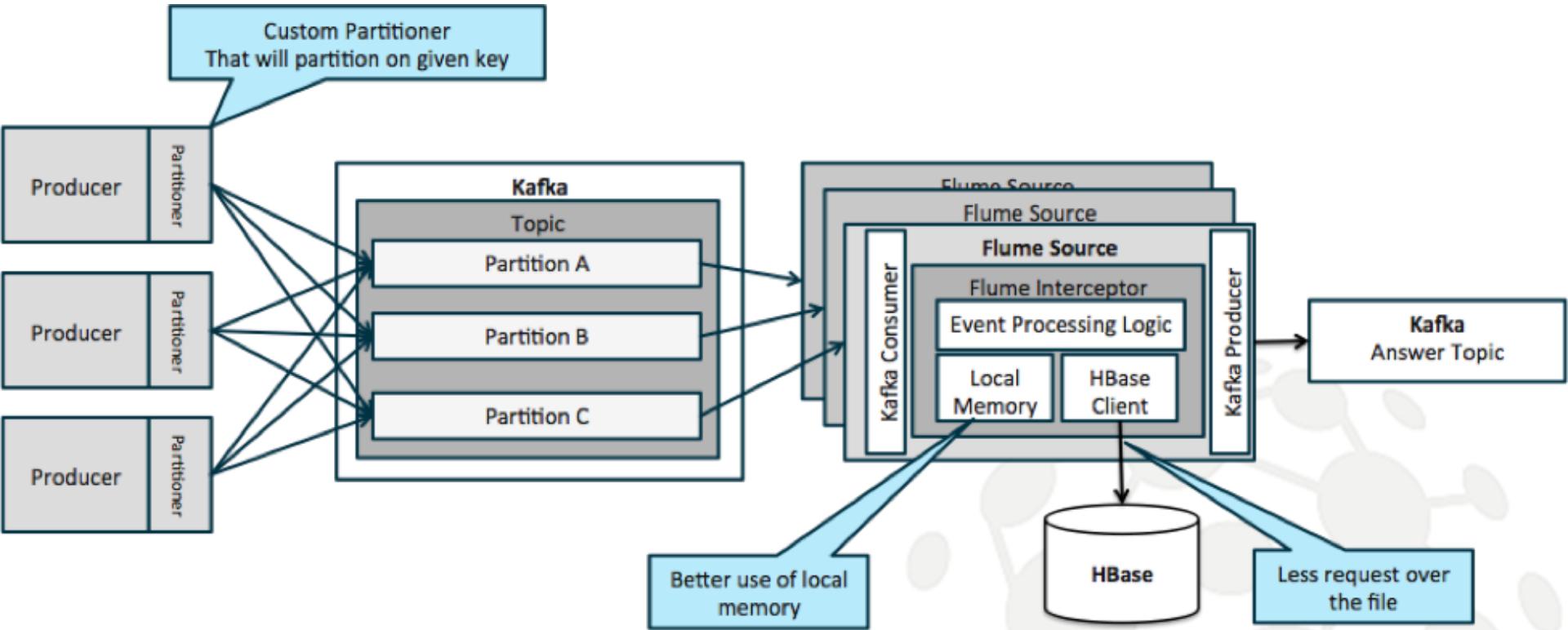
# Streaming Ingestion Architecture



# NRT Event Processing



# NRT Partitioned Event Processing



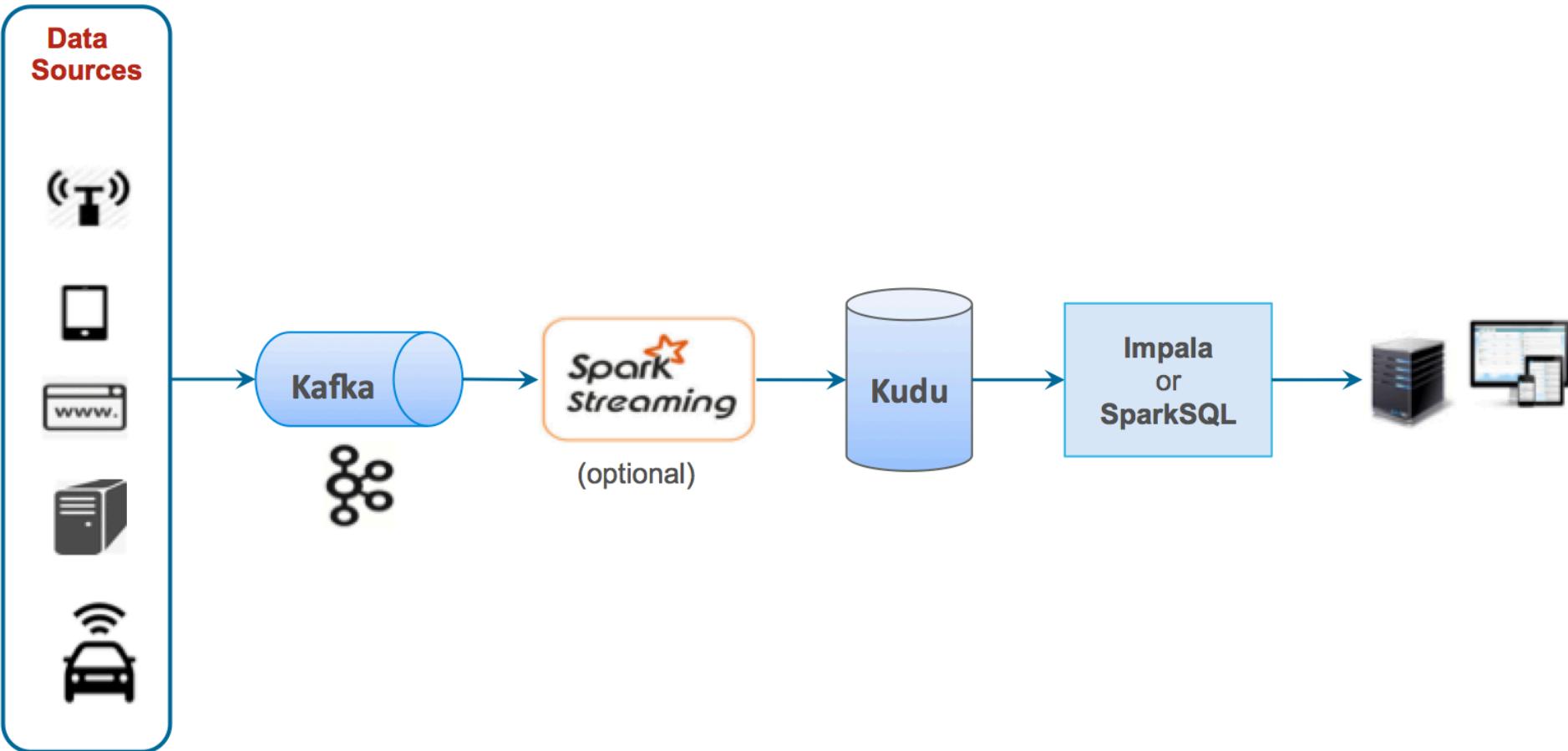


# Complex Topology for aggregations

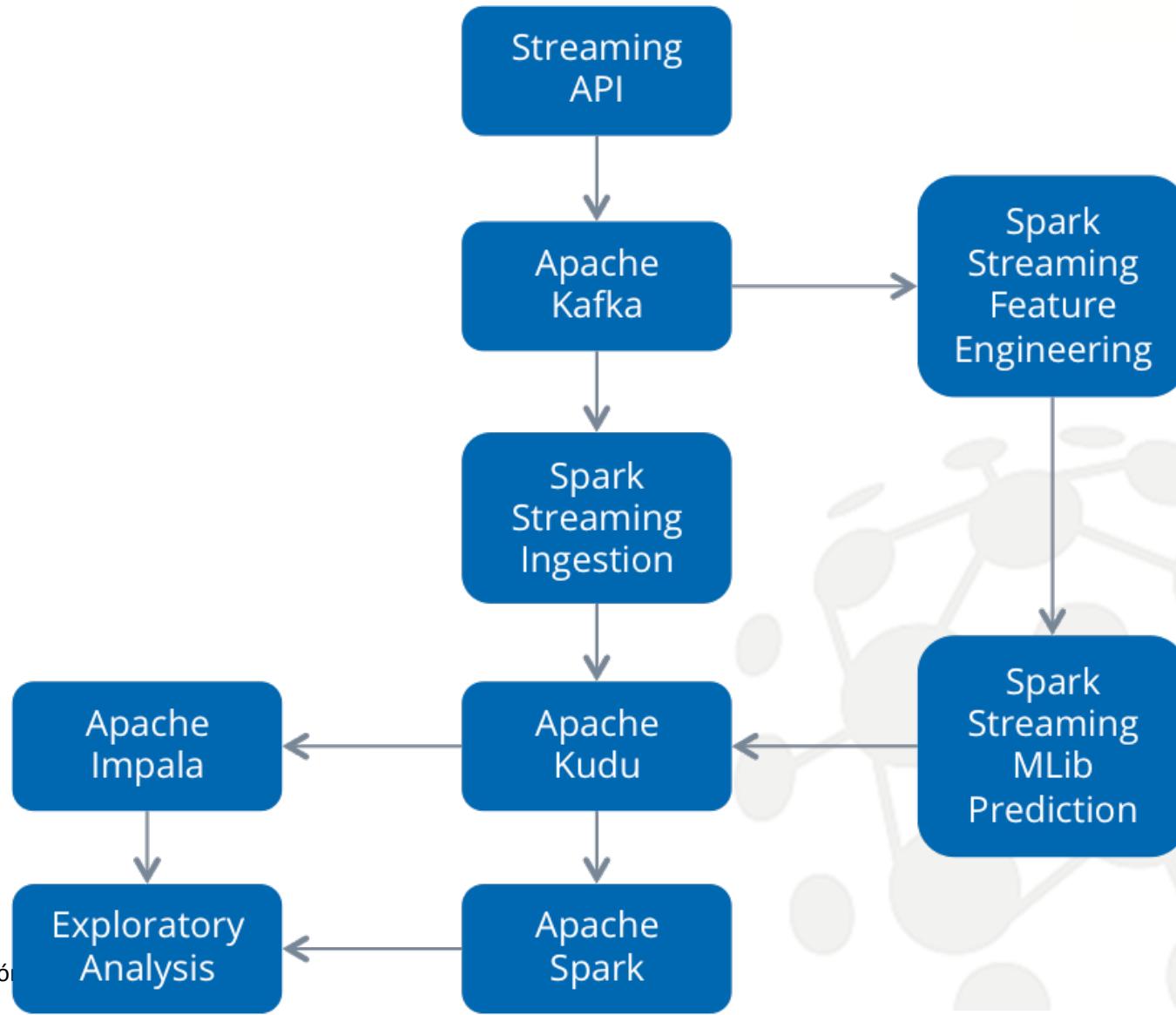




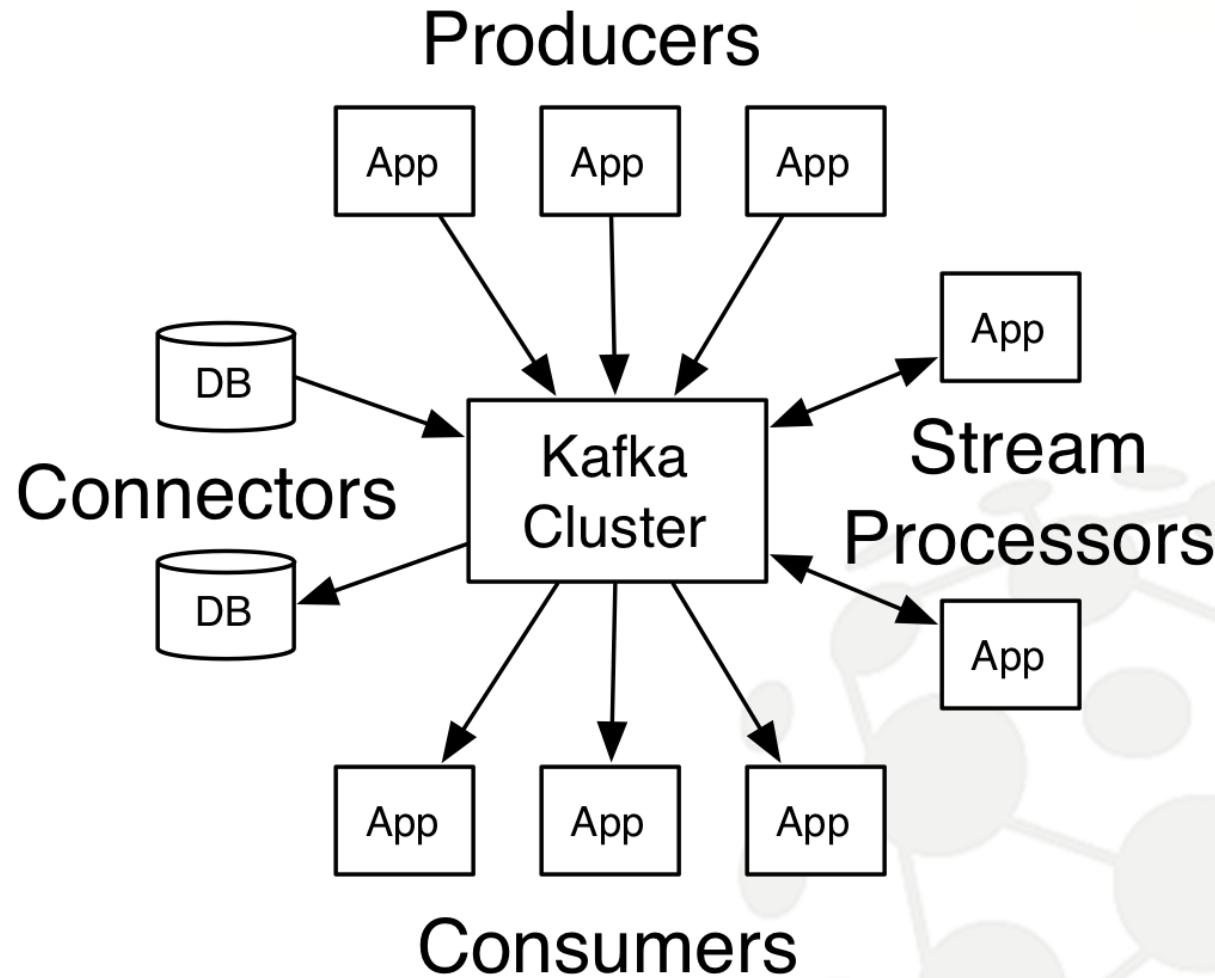
# Simple architecture



## A demo architecture

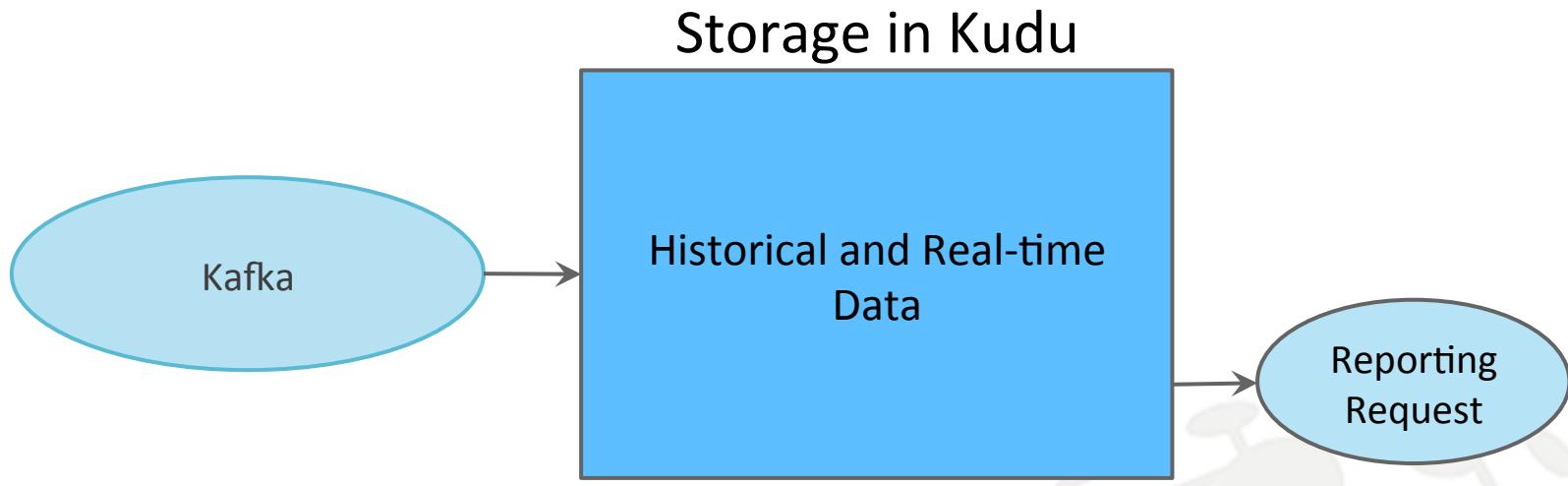


# Apache Kafka



<https://kafka.apache.org/documentation.html>

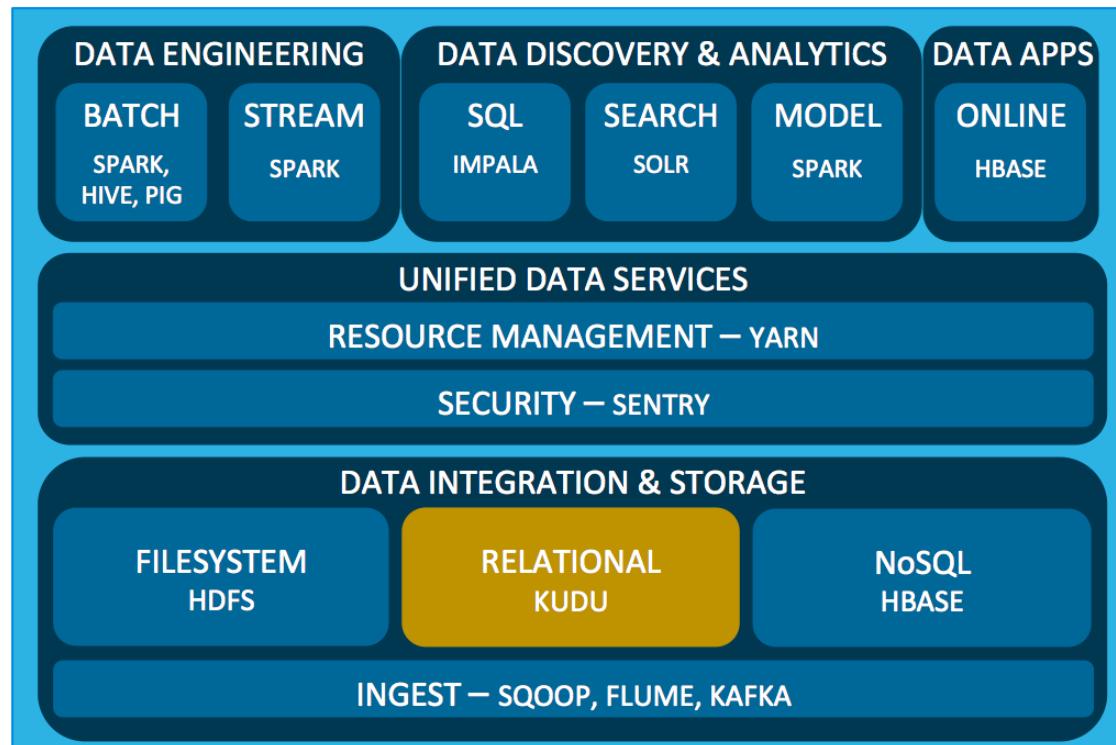






# ¿Qué es Apache Kudu?

- High-performance distributed storage engine



# Apache Storm

- Apache storm is a stream processing system build on top of HDFS
- Apache storm has it's own API's and do not use Map/Reduce
- It's a one message at time in core and micro batch is built on top of it(trident)
- Built by twitter



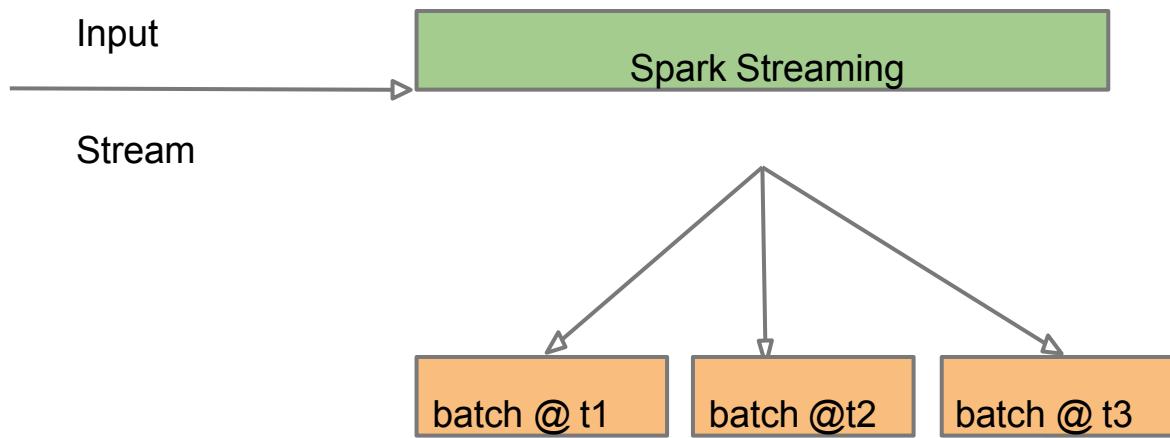
# Spark Streaming

- Spark streaming is a fast batch processing system
- Spark streaming collects stream data into small batch and runs batch processing on it
- Batch can be as small as 1s to as big as multiple hours
- Spark job creation and execution overhead is so low it can do all that under a sec
- These batches are called as DStreams





# Discretized streams (DStream)



Input stream is divided into multiple discrete batches. Batch is configurable.

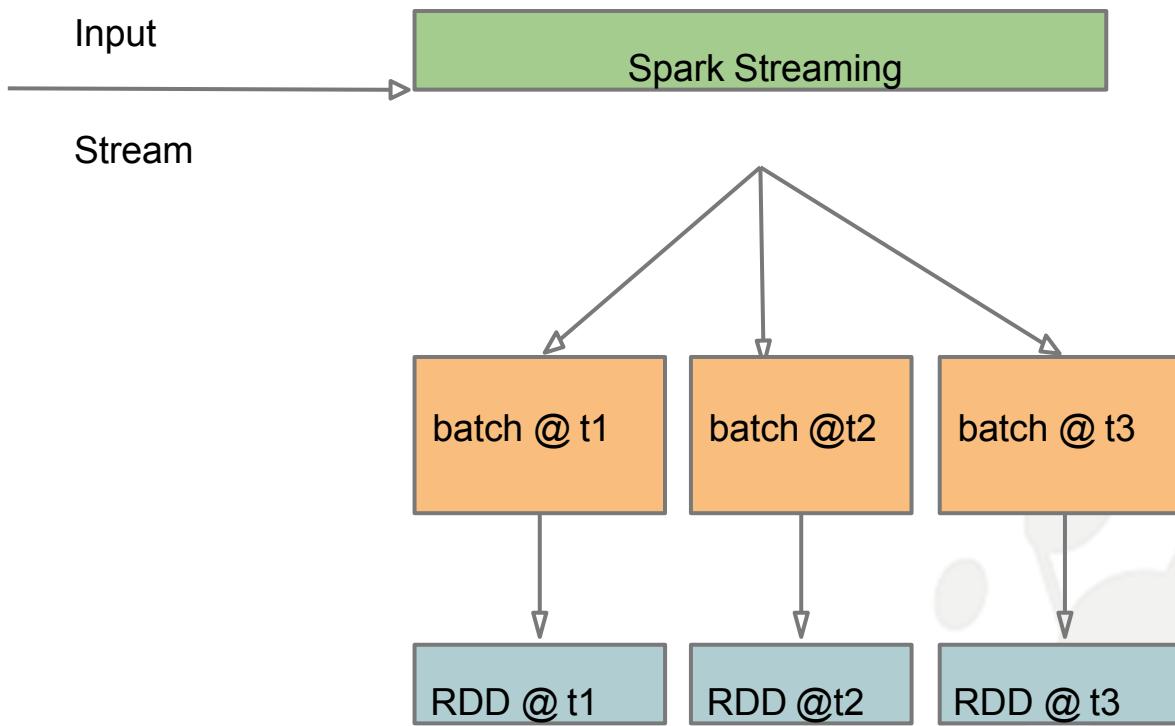


- Discretized streams
- Each batch of data is converted to small discrete batches
- Batch size can be from 1s - multiple mins
- DStream can be constructed from
  - Sockets
  - Kafka
  - HDFS
  - Custom receivers





# DStream to RDD





## Dstream to RDD

- Each batch of Dstream is represented as RDD underneath
- These RDD are replicated in cluster for fault tolerance
- Every DStream operation result in RDD transformation
- There are API's to access these RDD is directly
- Can combine stream and batch processing



# Socket stream

- Ability to listen to any socket on remote machines
- Need to configure host and port
- Both Raw and Text representation of socket available
- Built in retry mechanism



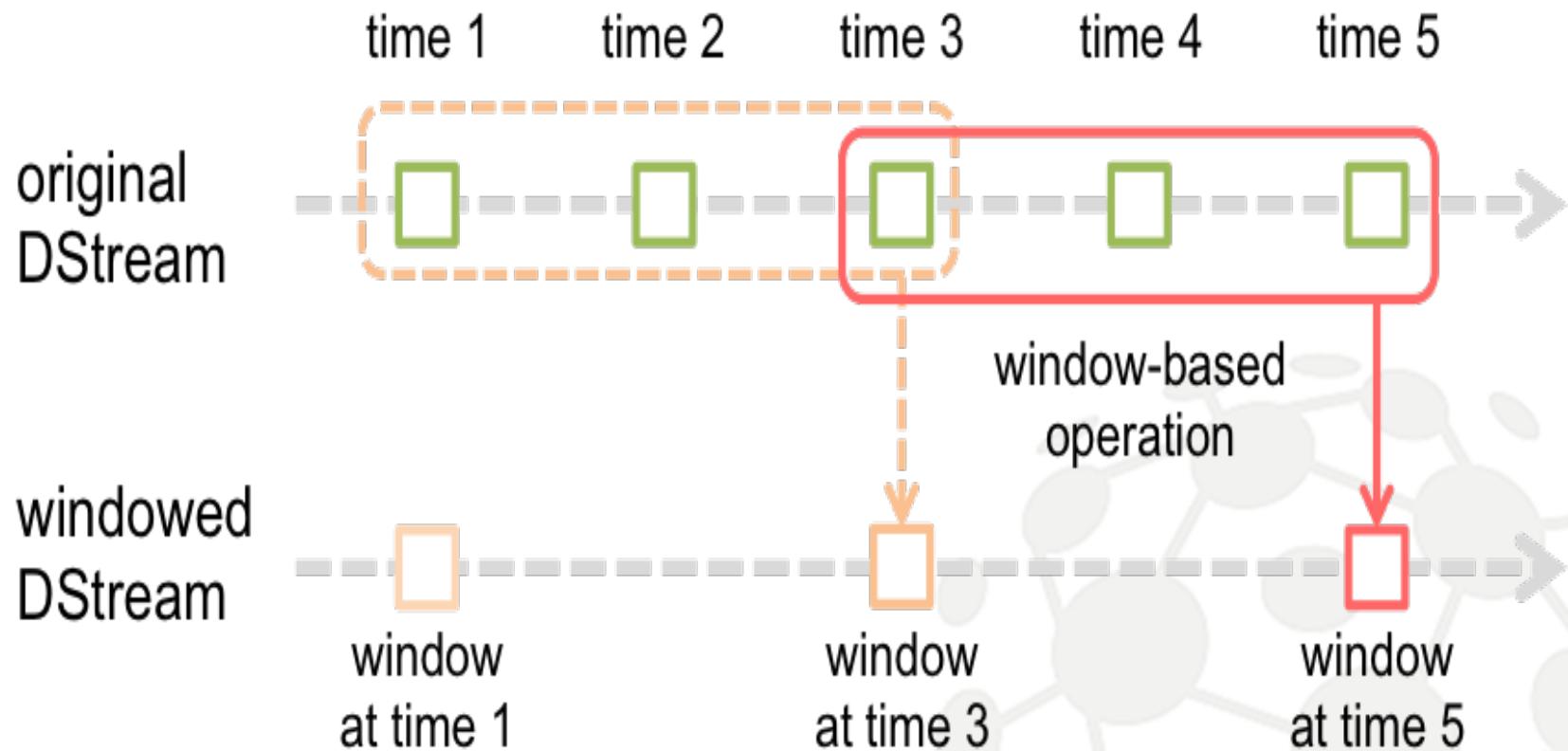
# File Stream

- File streams allows for track new files in a given directory on HDFS
- Whenever there is new file appears, spark streaming will pick it up
- Only works for new files, modification for existing files will not be considered
- Tracked using file creation time



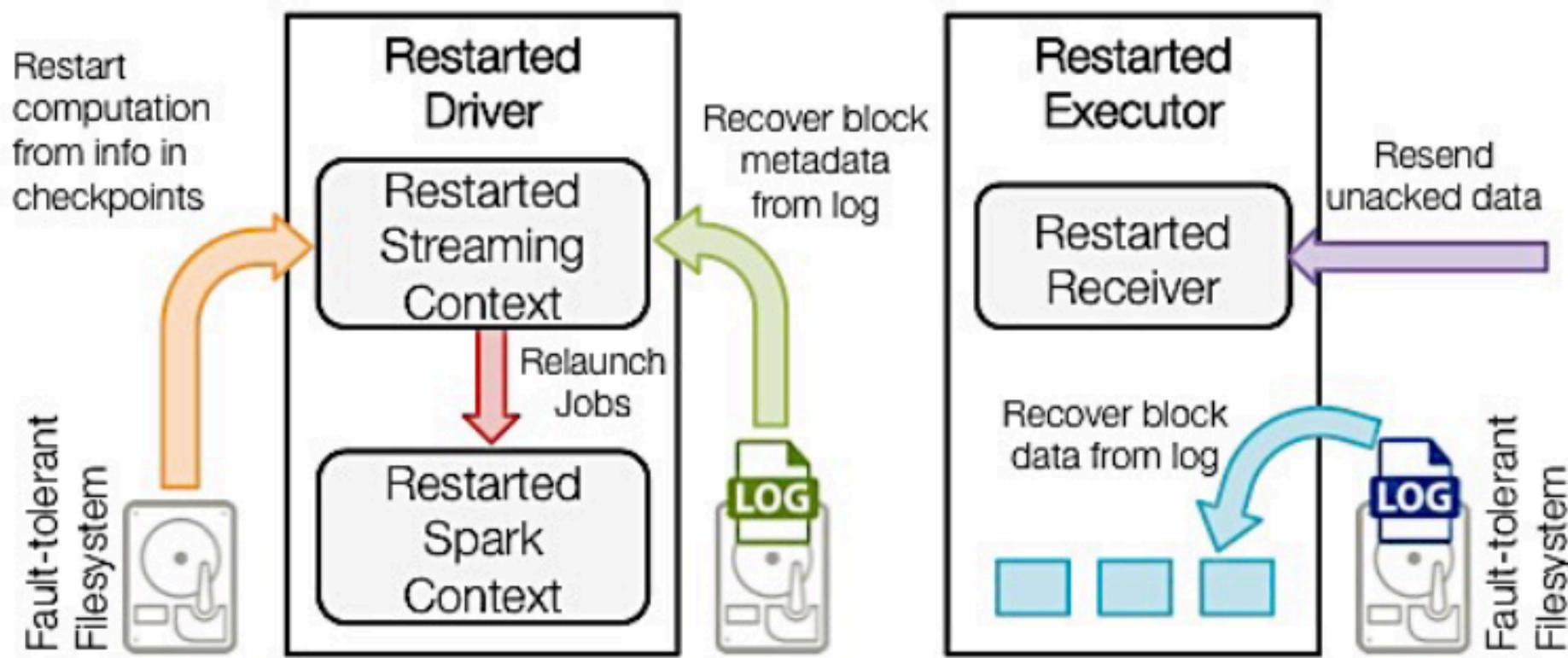


## Windows based operators



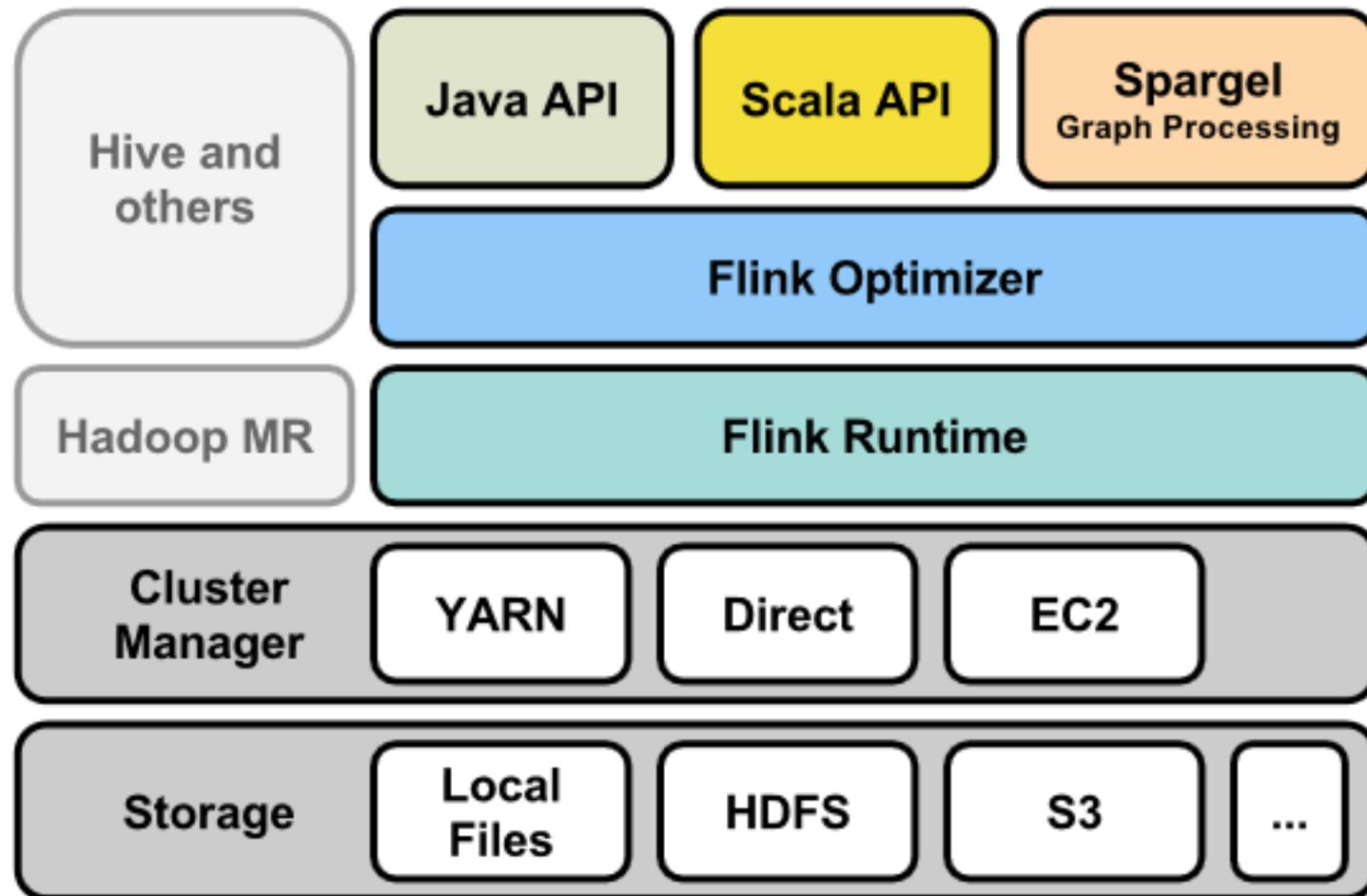


## Recovery





# La amezana para Spark





## Resumen

- En muchas situaciones se requiere un análisis en tiempo real
- Es común trabajar en arquitecturas lambda preparadas para sacar lo mejor de los mundos batch y real time
- Tenemos diversas herramientas disponibles, cada una de ellas con un propósito concreto

