# Data Stream Processing

Amir H. Payberah
`payberah@kth.se`
2025-09-16
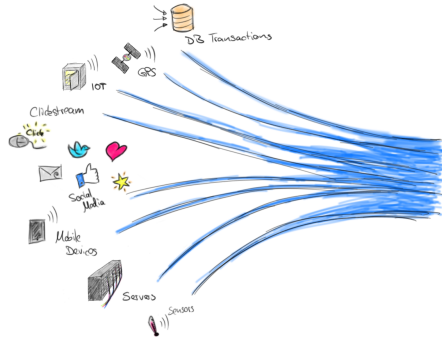
# Where Are We?



**Data Processing**

| Graph Data | Structured Data | Machine Learning |
|---|---|---|
| Pregel, GraphLab, PowerGraph GraphX, X-Streem, Chaos | Spark SQL | Mllib Tensorflow |

| Batch Data | Streaming Data |
|---|---|
| MapReduce, Dryad FlumeJava, Spark | Storm, SEEP, Naiad, Spark Streaming, Flink, Millwheel, Google Dataflow |

**Data Storage**

| Distributed File Systems | NoSQL Databases | Distributed Messaging Systems |
|---|---|---|
| GFS, Flat FS | Dynamo, BigTable, Cassandra | Kafka |

**Resource Management**

Mesos, YARN

# Stream Processing

- ▸ **Stream processing** is the **real-time** computation of **continuously** incoming data.

- ▸ The **input data** is **unbounded**: a **series of events**, no predetermined **beginning or end**.

- Data stream is unbound data, which is broken into a sequence of individual tuples.
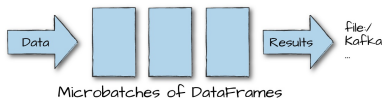
- A data tuple is the atomic data item in a data stream.

- Data stream is unbound data, which is broken into a sequence of individual tuples.

- A data tuple is the atomic data item in a data stream.

- Can be structured, semi-structured, and unstructured.
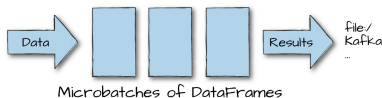
# Streaming Processing Patterns

▶ Micro-batch systems
  • Batch engines
  • Slicing up the unbounded data into a sets of bounded data, then process each batch.



Microbatches of DataFrames
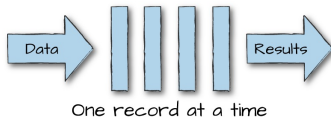
# Streaming Processing Patterns

▶ Micro-batch systems
  • Batch engines
  • Slicing up the unbounded data into a sets of bounded data, then process each batch.



Microbatches of DataFrames

▶ Continuous processing-based systems
  • Each node in the system continually listens to messages from other nodes and outputs new updates to its child nodes.



One record at a time

# Event and Processing Time

# Event Time vs. Processing Time

- Event time: the time at which events actually occurred.
  - Timestamps inserted into each record at the source.

- Prcosseing time: the time when the record is received at the streaming application.

- Windowing: dividing a continuous stream of data into discrete chunks (windows).

▶ Windowing: dividing a continuous stream of data into discrete chunks (windows).
  • Count-based, Time-based, etc.

- Windowing: dividing a continuous stream of data into discrete chunks (windows).
  - Count-based, Time-based, etc.

- Triggering: defines when a computation within a window should be performed.

# Windowing and Triggering

- Windowing: dividing a continuous stream of data into discrete chunks (windows).
  - Count-based, Time-based, etc.

- Triggering: defines when a computation within a window should be performed.
  - Count-based policy: the maximum number of tuples a window buffer can hold

# Windowing and Triggering

- Windowing: dividing a continuous stream of data into discrete chunks (windows).
  - Count-based, Time-based, etc.

- Triggering: defines when a computation within a window should be performed.
  - Count-based policy: the maximum number of tuples a window buffer can hold
  - Time-based policy: based on processing or event time period

- Two possibilites: tumbling and sliding

- Two possibilites: tumbling and sliding

- Tumbling window: when the buffer fills up, all the tuples are evicted.

- Two possibilites: tumbling and sliding

- Tumbling window: when the buffer fills up, all the tuples are evicted.

| | 1 | 2 1 | 3 2 1 | 4 3 2 1 | | 5 | 6 5 |

- Sliding window: when the buffer fills up, older tuples are evicted.

| | 1 | 2 1 | 3 2 1 | 4 3 2 1 | 5 4 3 2 | 6 5 4 3 |

- The system buffers up incoming data into windows until some amount of processing time has passed.

- E.g., five-minute fixed windows



[https://www.oreilly.com/ideas/the-world-beyond-batch-streaming-101]

▶ Reflect the times at which events actually happened.



[https://www.oreilly.com/ideas/the-world-beyond-batch-streaming-101]

▶ Reflect the times at which events actually happened.

▶ Handling out-of-order evnets.



[https://www.oreilly.com/ideas/the-world-beyond-batch-streaming-101]

▶ **Watermarking** helps a stream processing system to deal with lateness.

▶ **Watermarking** helps a stream processing system to deal with lateness.

▶ Watermarks flow as part of the data stream and carry a timestamp t.

▶ **Watermarking** helps a stream processing system to deal with lateness.

▶ Watermarks flow as part of the data stream and carry a timestamp t.

▶ A `W(t)` declares that event time has reached time t in that stream

# Triggering According to Event Time (2/2)

- **Watermarking** helps a stream processing system to deal with lateness.

- Watermarks flow as part of the data stream and carry a timestamp t.

- A `W(t)` declares that event time has reached time t in that stream

- A watermark is a threshold to specify how long the system waits for late events.

# Spark Streaming

▶ Run a streaming computation as a series of very small, deterministic batch jobs.

▶ Run a streaming computation as a series of very small, deterministic batch jobs.

- Chops up the live stream into batches of X seconds.

- Treats each batch as RDDs and processes them using RDD operations.

- Run a streaming computation as a series of very small, deterministic batch jobs.

  - Chops up the live stream into batches of X seconds.

  - Treats each batch as RDDs and processes them using RDD operations.

  - Discretized Stream Processing (DStream)

# DStream (1/2)

▶ **DStream**: sequence of RDDs representing a stream of data.

▶ **DStream**: sequence of RDDs representing a stream of data.

▶ Any operation applied on a DStream translates to operations on the underlying RDDs.

# StreamingContext

- StreamingContext is the main entry point of all Spark Streaming functionality.

```
val conf = new SparkConf().setAppName(appName).setMaster(master)
val ssc = new StreamingContext(conf, Seconds(1))
```

- The second parameter, Seconds(1), represents the time interval at which streaming data will be divided into batches.

- Socket connection
  - Creates a DStream from text data received over a TCP socket connection.

```
ssc.socketTextStream("localhost", 9999)
```

► Socket connection
  • Creates a DStream from text data received over a TCP socket connection.

```
ssc.socketTextStream("localhost", 9999)
```

► File stream
  • Reads data from files.

```
streamingContext.fileStream[KeyClass, ValueClass, InputFormatClass](dataDirectory)

streamingContext.textFileStream(dataDirectory)
```

# Input Operations

- **Socket** connection
  - Creates a DStream from text data received over a TCP socket connection.

```
ssc.socketTextStream("localhost", 9999)
```

- **File** stream
  - Reads data from files.

```
streamingContext.fileStream[KeyClass, ValueClass, InputFormatClass](dataDirectory)

streamingContext.textFileStream(dataDirectory)
```

- Connectors with external sources, e.g., Twitter, Kafka, Flume, Kinesis, ...

- Transformations on DStreams are still lazy!

- DStreams support many of the transformations available on normal Spark RDDs.

- Transformations on DStreams are still lazy!

- DStreams support many of the transformations available on normal Spark RDDs.

- Computation is kicked off explicitly by a call to the `start()` method.

- `map`: a new DStream by passing each element of the source DStream through a given function.

▶ `map`: a new DStream by passing each element of the source DStream through a given function.

▶ `reduce`: a new DStream of single-element RDDs by aggregating the elements in each RDD using a given function.

- `map`: a new DStream by passing each element of the source DStream through a given function.

- `reduce`: a new DStream of single-element RDDs by aggregating the elements in each RDD using a given function.

- `reduceByKey`: a new DStream of (K, V) pairs where the values for each key are aggregated using the given reduce function.

- First we create a `StreamingContex`

```scala
import org.apache.spark._
import org.apache.spark.streaming._

// Create a local StreamingContext with two working threads and batch interval of 1 second.
val conf = new SparkConf().setMaster("local[2]").setAppName("NetworkWordCount")
val ssc = new StreamingContext(conf, Seconds(1))
```

- Create a DStream that represents streaming data from a TCP source.
- Specified as hostname (e.g., localhost) and port (e.g., 9999).

```scala
val lines = ssc.socketTextStream("localhost", 9999)
```

► Use `flatMap` on the stream to split the records text to words.
► It creates a new DStream.

```
val words = lines.flatMap(_.split(" "))
```

- Map the `words` DStream to a DStream of (`word, 1`).
- Get the frequency of words in each batch of data.
- Finally, print the result.

```
val pairs = words.map(word => (word, 1))

val wordCounts = pairs.reduceByKey(_ + _)

wordCounts.print()
```

▶ Start the computation and wait for it to terminate.

```
// Start the computation
ssc.start()

// Wait for the computation to terminate
ssc.awaitTermination()
```

# Example - Word Count (6/6)

```scala
val conf = new SparkConf().setMaster("local[2]").setAppName("NetworkWordCount")
val ssc = new StreamingContext(conf, Seconds(1))

val lines = ssc.socketTextStream("localhost", 9999)
val words = lines.flatMap(_.split(" "))
val pairs = words.map(word => (word, 1))
val wordCounts = pairs.reduceByKey(_ + _)
wordCounts.print()

ssc.start()
ssc.awaitTermination()
```

# Window Operations (1/2)

- Spark provides a set of transformations that apply to a over a sliding window of data.

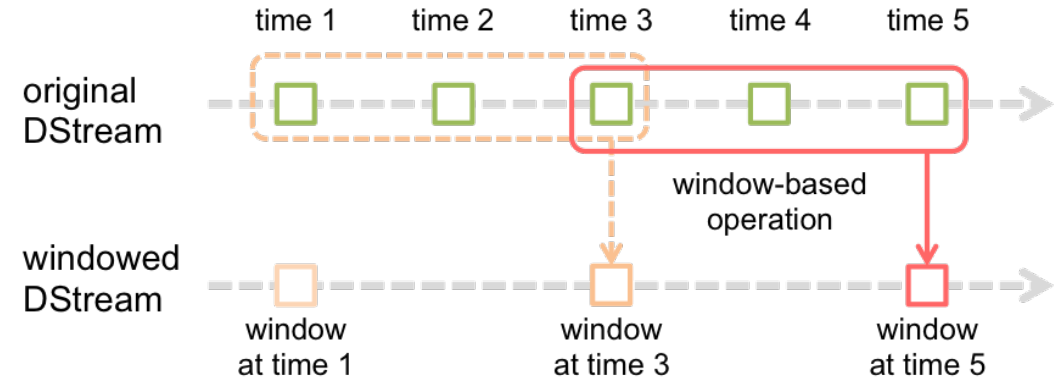

- Spark provides a set of transformations that apply to a over a sliding window of data.

- A window is defined by two parameters: window length and slide interval.

# Window Operations (1/2)

- Spark provides a set of transformations that apply to a over a sliding window of data.

- A window is defined by two parameters: window length and slide interval.

- A tumbling window effect can be achieved by making slide interval = window length

▶ `reduceByWindow(func, windowLength, slideInterval)`
  - Returns a new single-element DStream, created by aggregating elements in the stream over a sliding interval using func.

▶ `reduceByWindow(func, windowLength, slideInterval)`
- Returns a new single-element DStream, created by aggregating elements in the stream over a sliding interval using func.

▶ `reduceByKeyAndWindow(func, windowLength, slideInterval)`
- Called on a DStream of (K, V) pairs.
- Returns a new DStream of (K, V) pairs where the values for each key are aggregated using function func over batches in a sliding window.

# Example - Word Count with Window

```scala
val conf = new SparkConf().setMaster("local[2]").setAppName("NetworkWordCount")
val ssc = new StreamingContext(conf, Seconds(1))

val lines = ssc.socketTextStream("localhost", 9999)
val words = lines.flatMap(_.split(" "))
val pairs = words.map(word => (word, 1))
val windowedWordCounts = pairs.reduceByKeyAndWindow(_ + _, Seconds(30), Seconds(10))
windowedWordCounts.print()

ssc.start()
ssc.awaitTermination()
```

# Structured Streaming

▶ Treating a live data stream as a table that is being continuously appended.



Data stream as an unbounded table

User's batch-like query on input table

Incremental execution on streaming data

▶ Three output modes:

1. Append: only the new rows appended to the result table since the last trigger will be written to the external storage.

- Three output modes:

1. Append: only the new rows appended to the result table since the last trigger will be written to the external storage.

2. Complete: the entire updated result table will be written to external storage.

# Output Modes

- Three output modes:

1. Append: only the new rows appended to the result table since the last trigger will be written to the external storage.

2. Complete: the entire updated result table will be written to external storage.

3. Update: only the rows that were updated in the result table since the last trigger will be changed in the external storage.

- Define input sources.

- Use spark.readStream to create a DataStreamReader.

```scala
val spark = SparkSession.builder.master("local[2]").appName("appname").getOrCreate()

val lines = spark.readStream.format("socket")
    .option("host", "localhost")
    .option("port", 9999)
    .load()
```

- Transform data.

- E.g., below counts is a streaming DataFrame that represents the running word counts.

```
import org.apache.spark.sql.functions._

val words = lines.select(split(col("value"), " ").as("word"))

val wordCounts = words.groupBy("word").count()
```

- Define output sink and output mode.

- Use DataFrame.writeStream to define how to write the processed output data.

- Start the query.

```
val query = wordCounts.writeStream.format("console").outputMode("complete").start()

query.awaitTermination()
```

[https://spark.apache.org/docs/latest/structured-streaming-programming-guide.html]

# Data Stream Storage

► We need disseminate streams of events from various producers to various consumers.

- Messaging systems



Message

www.defit.org

# What is Messaging System?

- Messaging system is an approach to notify consumers about new events.

- Messaging system is an approach to notify consumers about new events.

- Messaging systems
  - Direct messaging
  - Message brokers

- Necessary in latency critical applications (e.g., remote surgery).
- A producer sends a message containing the event, which is pushed to consumers.

# Direct Messaging (1/2)

- Necessary in latency critical applications (e.g., remote surgery).
- A producer sends a message containing the event, which is pushed to consumers.
- Both consumers and producers have to be online at the same time.

▶ What happens if a consumer crashes or temporarily goes offline? (not durable)

- What happens if a consumer crashes or temporarily goes offline? (not durable)
- What happens if producers send messages faster than the consumers can process?

- What happens if a consumer crashes or temporarily goes offline? (not durable)
- What happens if producers send messages faster than the consumers can process?
  - Dropping messages
  - Backpressure

- What happens if a consumer crashes or temporarily goes offline? (not durable)
- What happens if producers send messages faster than the consumers can process?
  - Dropping messages
  - Backpressure

- We need message brokers that can log events to process at a later time.

[https://bluesyemre.com/2018/10/16/thousands-of-scientists-publish-a-paper-every-five-days]

- A message broker decouples the producer-consumer interaction.
- It runs as a server, with producers and consumers connecting to it as clients.

# Message Broker

- A message broker decouples the producer-consumer interaction.
- It runs as a server, with producers and consumers connecting to it as clients.
- Producers write messages to the broker, and consumers receive them by reading them from the broker.

# Message Broker

- A message broker decouples the producer-consumer interaction.
- It runs as a server, with producers and consumers connecting to it as clients.
- Producers write messages to the broker, and consumers receive them by reading them from the broker.
- Consumers are generally asynchronous.

▶ In typical message brokers, once a message is consumed, it is deleted.

- In typical message brokers, once a message is consumed, it is deleted.

- Log-based message brokers durably store all events in a sequential log.

# Logs-Based Message Broker

- In typical message brokers, once a message is consumed, it is deleted.

- Log-based message brokers durably store all events in a sequential log.

- A log is an append-only sequence of records on disk.

# Logs-Based Message Broker

- In typical message brokers, once a message is consumed, it is deleted.

- Log-based message brokers durably store all events in a sequential log.

- A log is an append-only sequence of records on disk.

- A producer sends a message by appending it to the end of the log.

- A consumer receives messages by reading the log sequentially.

# Kafka - A Log-Based Message Broker

▶ Kafka is a distributed, topic oriented, partitioned, replicated commit log service.

▶ Kafka is a distributed, topic oriented, partitioned, replicated commit log service.

- Kafka is a distributed, topic oriented, partitioned, replicated commit log service.

- Kafka is a distributed, topic oriented, partitioned, replicated commit log service.

- Kafka is a distributed, topic oriented, partitioned, replicated commit log service.

▶ Topics are queues: a stream of messages of a particular type

▶ Each message is assigned a sequential id called an offset.

▶ Topics are logical collections of partitions (the physical files).
  - Ordered
  - Append only
  - Immutable

▶ Ordering is only guaranteed within a partition for a topic.

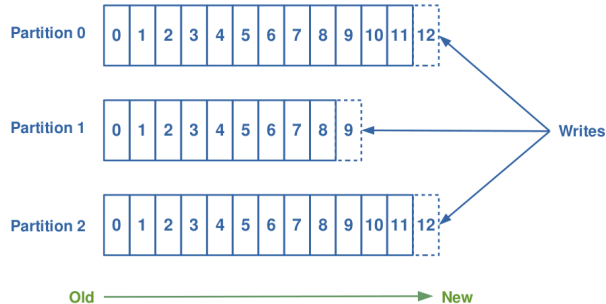- Ordering is only guaranteed within a partition for a topic.
- Messages sent by a producer to a particular topic partition will be appended in the order they are sent.

- Ordering is only guaranteed within a partition for a topic.

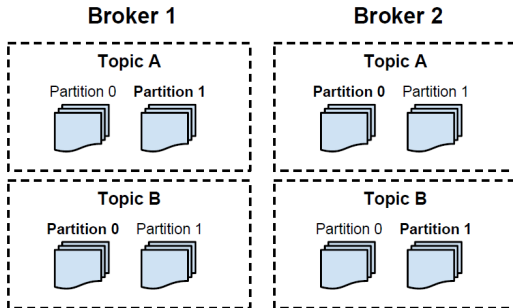- Messages sent by a producer to a particular topic partition will be appended in the order they are sent.

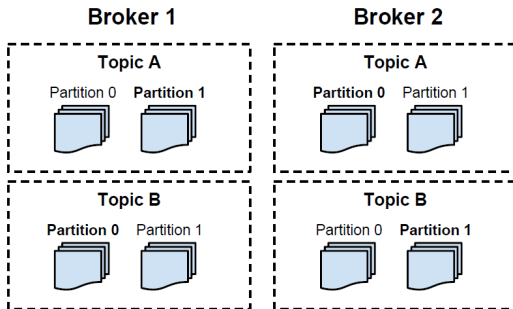- A consumer instance sees messages in the order they are stored in the log.

- Partitions of a topic are replicated: fault-tolerance

# Topics and Partition (5/6)

- Partitions of a topic are replicated: fault-tolerance

- A broker contains some of the partitions for a topic.

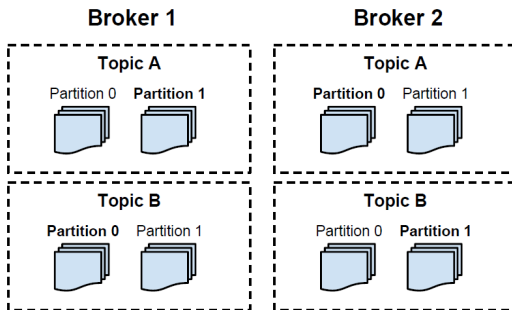# Topics and Partition (5/6)
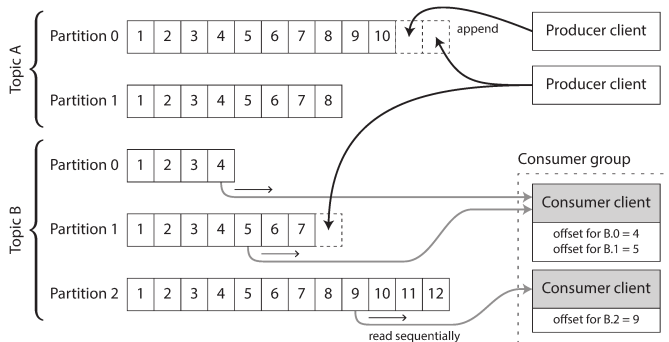
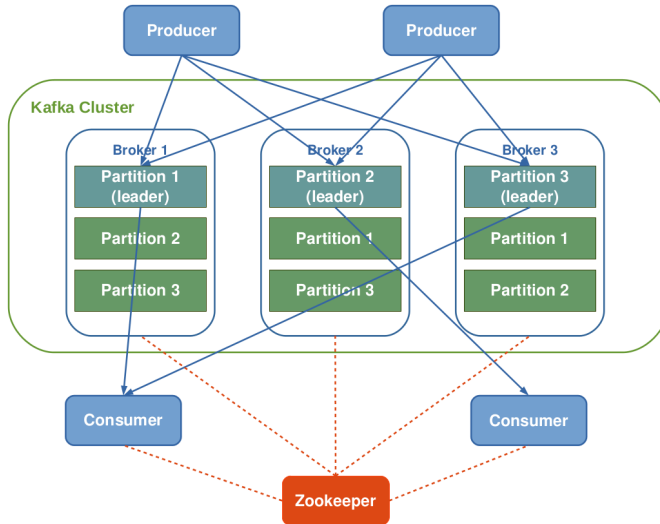- Partitions of a topic are replicated: fault-tolerance

- A broker contains some of the partitions for a topic.

- One broker is the leader of a partition: all writes and reads must go to the leader.

# Summary

# Summary

- Data stream, unbounded data, tuples

- Event-time vs. processing time

- Windowing and triggering

- Messaging system and partitioned logs

- Kafka: distributed, topic oriented, partitioned, replicated log service

- Spark streaming and structured streaming

# References

▶ J. Kreps et al., "Kafka: A distributed messaging system for log processing", NetDB 2011

▶ M. Zaharia et al., "Spark: The Definitive Guide", O'Reilly Media, 2018 - Chapter 20

▶ T. Akidau et al., "The dataflow model: a practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing", VLDB 2015.

▶ M. Fragkoulis et al., "A Survey on the Evolution of Stream Processing Systems", 2020

▶ T. Akidau, "The world beyond batch: Streaming 101", https://www.oreilly.com/ideas/the-world-beyond-batch-streaming-101

Questions?