# Scalable Stream Processing - Spark Streaming and Beam

Amir H. Payberah
payberah@kth.se
2022-09-27

https://id2221kth.github.io
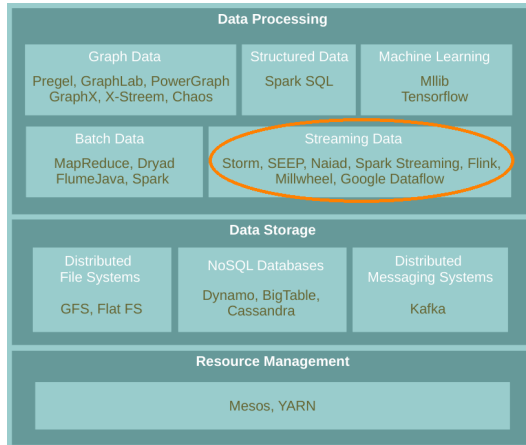
https://tinyurl.com/bdenpwc5

# 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

# Spark Streaming
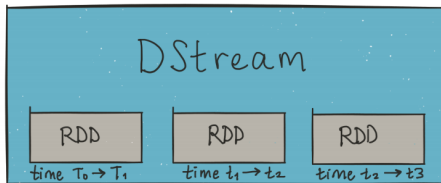
# Spark Streaming

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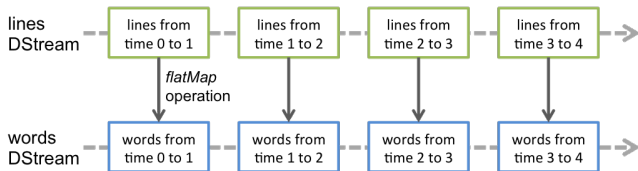
▶ **DStream**: sequence of RDDs representing a stream of data.

# DStream (2/2)

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

- 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)
```

- 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.

- 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.

- `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))
```
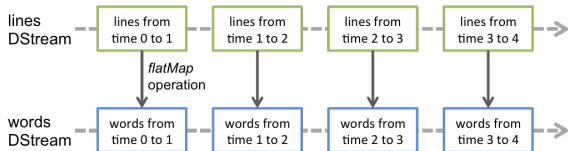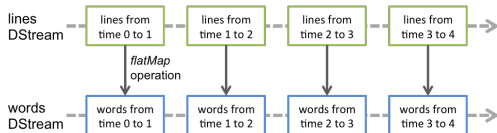
▶ Create a DStream that represents streaming data from a TCP source.

▶ 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.

```scala
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()
```

```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()
```
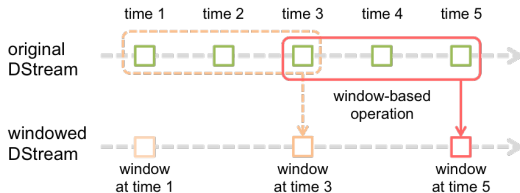
# Window Operations (1/2)

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

- 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.

- ▸ `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.
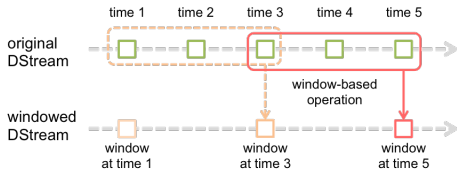
# Example - Word Count with Window

```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()
```

# What about States?

- Accumulate and aggregate the results from the start of the streaming job.

- Need to check the previous state of the RDD in order to do something with the current RDD.

- Spark supports stateful streams.

▶ It is mandatory that you provide a checkpointing directory for stateful streams.

```
val ssc = new StreamingContext(conf, Seconds(1))
ssc.checkpoint("path/to/persistent/storage")
```

- ▶ mapWithState
  - It is executed only on set of keys that are available in the last micro batch.

```
def mapWithState[StateType, MappedType](spec: StateSpec[K, V, StateType, MappedType]):
    DStream[MappedType]

StateSpec.function(updateFunc)
val updateFunc = (batch: Time, key: String, value: Option[Int], state: State[Int])
```

- ▶ Define the update function (partial updates) in StateSpec.

```scala
val ssc = new StreamingContext(conf, Seconds(1))
ssc.checkpoint(".")

val lines = ssc.socketTextStream(IP, Port)
val words = lines.flatMap(_.split(" "))
val pairs = words.map(word => (word, 1))

val stateWordCount = pairs.mapWithState(StateSpec.function(updateFunc))

val updateFunc = (key: String, value: Option[Int], state: State[Int]) => {
  val newCount = value.getOrElse(0)
  val oldCount = state.getOption.getOrElse(0)
  val sum = newCount + oldCount
  state.update(sum)
  (key, sum)
}
```
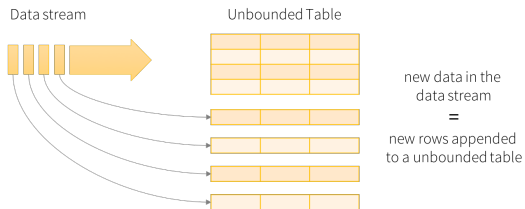
- The first micro batch contains a message a.

- `updateFunc = (key:  String, value:  Option[Int], state:  State[Int]) => (key, sum)`

- Input: `key = a`, `value = Some(1)`, `state = 0`

- Output: `key = a`, `sum = 1`

- The second micro batch contains messages `a` and `b`.

- `updateFunc = (key:  String, value:  Option[Int], state:  State[Int]) => (key, sum)`

- Input: `key = a`, `value = Some(1)`, `state = 1`

- Input: `key = b`, `value = Some(1)`, `state = 0`

- Output: `key = a`, `sum = 2`

- Output: `key = b`, `sum = 1`

▶ The third micro batch contains a message b.

▶ `updateFunc = (key:  String, value:  Option[Int], state:  State[Int]) => (key, sum)`

▶ Input: `key = b`, `value = Some(1)`, `state = 1`

▶ Output: `key = b`, `sum = 2`

# Structured Streaming

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



Data stream as an unbounded table

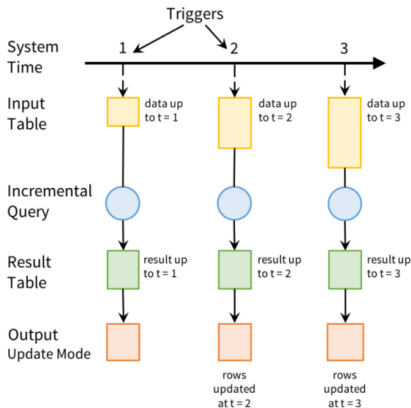# Programming Model (1/2)

- Defines a query on the input table, as a static table.
  - Spark automatically converts this batch-like query to a streaming execution plan.

- Specify triggers to control when to update the results.
  - Each time a trigger fires, Spark checks for new data (new row in the input table), and incrementally updates the result.

- 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.

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

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

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

- Define output sink and output mode.

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

```
val writer = counts.writeStream.format("console").outputMode("complete")
```

▶ Specify processing details.

```scala
\\ word count details
import org.apache.spark.sql.streaming._

val checkpointDir = "..."

val writer2 = writer
    .trigger(Trigger.ProcessingTime("1 second"))
    .option("checkpointLocation", checkpointDir)
```

- ▶ Start the query.

- ▶ streamingQuery represents an active query and can be used to manage the query.

```scala
val streamingQuery = writer2.start()
```

▶ Most of operations on DataFrame/Dataset are supported for streaming.

```scala
case class Call(action: String, time: Timestamp, id: Int)

val df: DataFrame = spark.readStream.json("s3://logs")
val ds: Dataset[Call] = df.as[Call]
```

▶ Selection and projection

```scala
df.select("action").where("id > 10") // using untyped APIs
ds.filter(_.id > 10).map(_.action) // using typed APIs
```

# Basic Operations (2/2)

- Aggregation

```
df.groupBy("action") // using untyped API
ds.groupByKey(_.action) // using typed API
```

- SQL commands

```
df.createOrReplaceTempView("dfView")
spark.sql("select count(*) from dfView") // returns another streaming DF
```

# Google Dataflow and Beam

- MillWheel is a framework for building low-latency data-processing applications.

- A dataflow graph of transformations (computations).

- Stream: unbounded data of (key, value, timestamp) records.
  - Timestamp: event-time

# Key Extraction Function and Computations

- Stream of (key, value, timestamp) records.
- Key extraction function: specified by the stream consumer to assign keys to records.
- Computation can only access state for the specific key.
- Multiple computations can extract different keys from the same stream.

# Persistent State

- Keep the states of the computations

- Managed on per-key basis

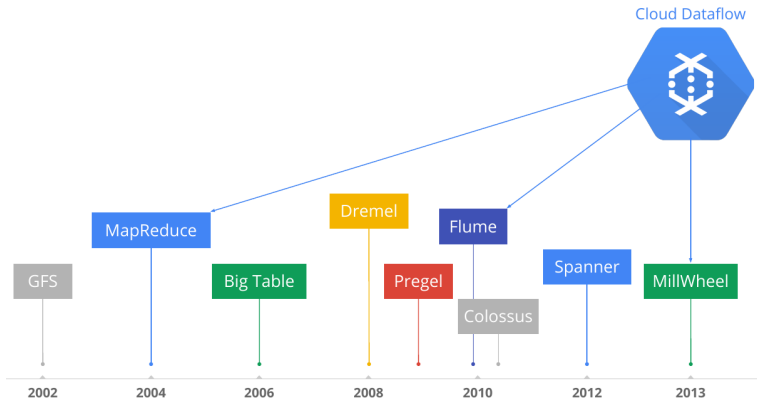- Stored in Bigtable or Spanner

- Common use: aggregation, joins, …

# What is Google Cloud Dataflow?

# Google Cloud Dataflow (1/2)

- Google managed service for unified batch and stream data processing.

# Google Cloud Dataflow (2/2)

- Open source Cloud Dataflow SDK

- Express your data processing pipeline using FlumeJava.

- If you run it in batch mode, it executed on the MapReduce framework.

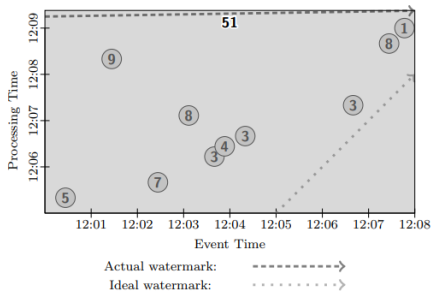- If you run it in streaming mode, it is executed on the MillWheel framework.

- Windowing determines where in event time data are grouped together for processing.

- Triggering determines when in processing time the results of groupings are emitted as panes.

▶ Batch processing
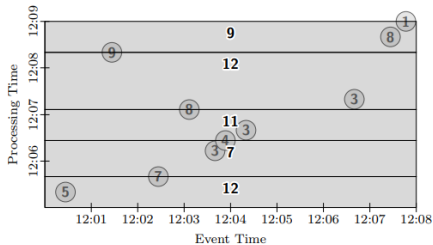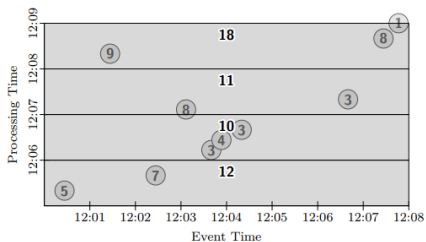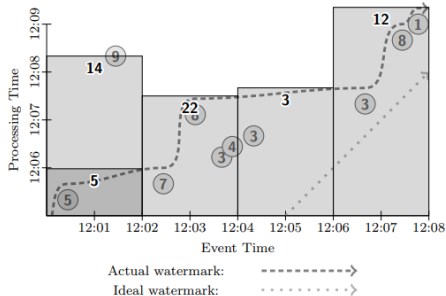
- Trigger at period (time-based triggers)
- Trigger at count (data-driven triggers)

- Fixed window, trigger at period (micro-batch)
- Fixed window, trigger at watermark (streaming)

# Where is Apache Beam?

- In 2016, Google Cloud Dataflow team announced its intention to donate the programming model and SDKs to the Apache Software Foundation.

- That resulted in the incubating project Apache Beam.

- Pipelines

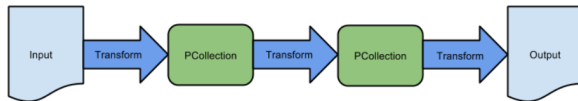- PCollections

- Transforms

- I/O sources and sinks
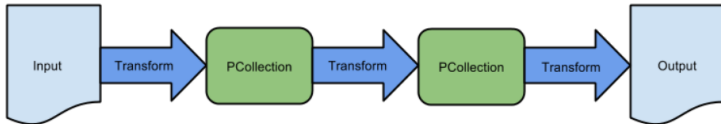
- A **pipeline** represents a **data processing job**.

- **Directed graph** of operating on data.

- A pipeline consists of **two** parts:
  - **Data** (`PCollection`)
  - **Transforms** applied to that data

# PCollections (1/2)

- A parallel collection of records

- Immutable

- Must specify bounded or unbounded

```java
// Create a Java Collection, in this case a List of Strings.
static final List<String> LINES = Arrays.asList("line 1", "line 2", "line 3");

PipelineOptions options = PipelineOptionsFactory.create();
Pipeline p = Pipeline.create(options);

// Create the PCollection
p.apply(Create.of(LINES)).setCoder(StringUtf8Coder.of())
```
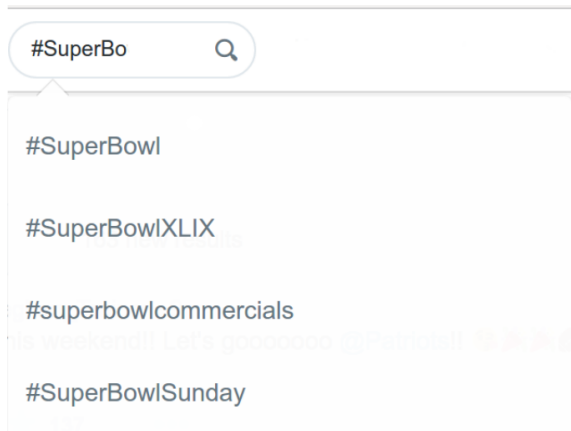
- A processing operation that transforms data

- Each transform accepts one (or multiple) `PCollections` as input, performs an operation, and produces one (or multiple) new `PCollections` as output.

- Core transforms: `ParDo, GroupByKey, Combine, Flatten`

# Example: HashTag Autocompletion (2/3)

# Example: HashTag Autocompletion (3/3)



```
Pipeline p = Pipeline.create();
p.begin()

    .apply(TextIO.Read.from("gs://…"))

    .apply(ParDo.of(new ExtractTags()))

    .apply(Count.perElement())

    .apply(ParDo.of(new ExpandPrefixes())

    .apply(Top.largestPerKey(3))

    .apply(TextIO.Write.to("gs://…"));

p.run();
```

# Summary

# Summary

- Spark
  - Mini-batch processing
  - DStream: sequence of RDDs
  - RDD and window operations
  - Structured streaming

- Google cloud dataflow
  - Pipeline
  - PCollection
  - Transforms

# References

- M. Zaharia et al., "Spark: The Definitive Guide", O'Reilly Media, 2018 - Chapters 20-23.

- M. Zaharia et al., "Discretized Streams: An Efficient and Fault-Tolerant Model for Stream Processing on Large Clusters", HotCloud'12.

- T. Akidau et al., "MillWheel: fault-tolerant stream processing at internet scale", VLDB 2013.

- 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.

- The world beyond batch: Streaming 102
  https://www.oreilly.com/ideas/the-world-beyond-batch-streaming-102

Questions?