

Lecture 5: Spark Streaming 1

Big Data Processing

1/2023
Thanh-Chung Dao Ph.D.

1

Agenda

- What is Spark Streaming
- Operation on DStreams

2

2



What is Spark Streaming

3

Email: info@ptu.edu.vn | Website: ptu.edu.vn

3

Spark Streaming

- Scalable, fault-tolerant stream processing system
- Receive data streams from input sources, process them in a cluster, push out to databases/dashboards

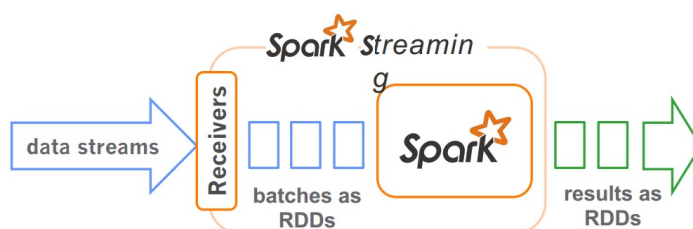


4

4

How does it work?

- The stream is treated as a **series** of very **small, deterministic batches** of data
- Spark treats each batch of data as RDDs and processes them using RDD operations
- Processed results are pushed out in batches

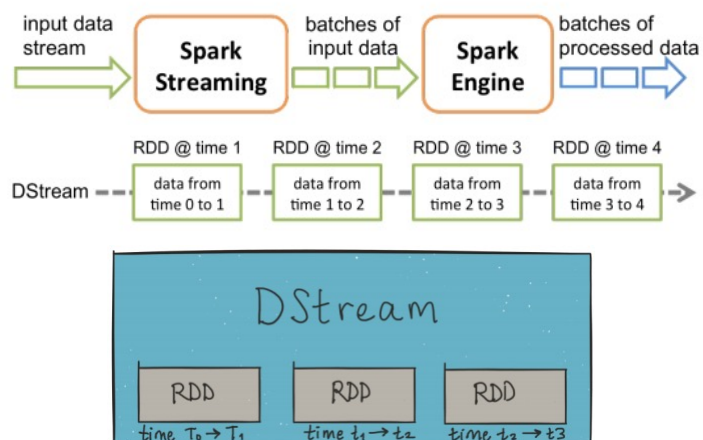


5

5

Discretized Stream (DStream)

- Sequence of RDDs representing a stream of data

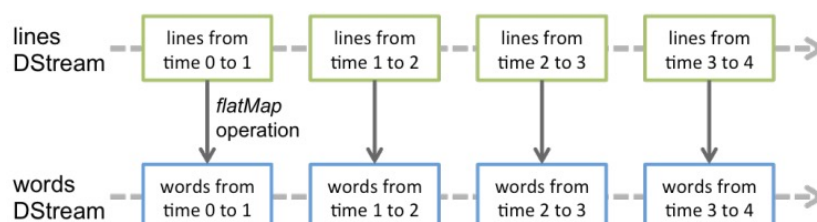


6

6

Discretized Stream (DStream)

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



7

7

StreamingContext

- The **main entry** point of all Spark Streaming functionality

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

- appName**: name of the application
- master**: a Spark, Mesos, or YARN cluster URL
- batchinterval**: time interval (in second) of each batch

8

8



Operation on DStreams

9

Email: info@ptu.edu.vn | Website: ptu.edu.vn

9

Operation on DStreams

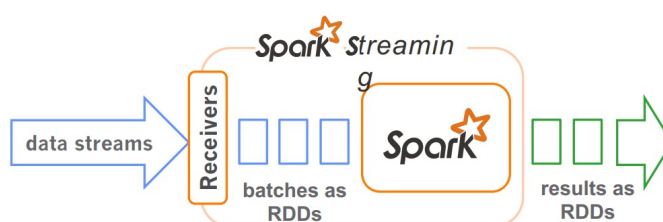
- Three categories
 - Input operation
 - Transformation operation
 - Output operation

10

10

Input Operations

- Every **input DStream** is associated with a **Receiver** object
- Two built-in categories of streaming sources:
 - Basic sources, e.g., file systems, socket connection
 - Advanced sources, e.g., Twitter, Kafka



11

11

Input Operations

- Basic sources
 - Socket connection

```
// Create a DStream that will connect to hostname:port
ssc.socketTextStream("localhost", 9999)
```

- File stream

```
streamingContext.fileStream[...] (dataDirectory)
```

- Advanced sources

```
val ssc = new StreamingContext(sparkContext, Seconds(1))
val tweets = TwitterUtils.createStream(ssc, auth)
```

12

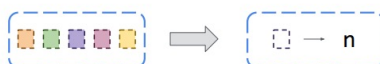
12

Transformation

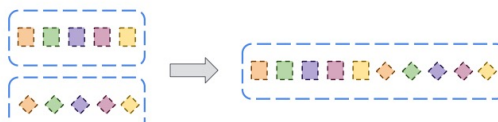
map,
flatMap,
filter



count,
reduce,
countByValue,
reduceByKey



union,
join
cogroup



13

13

Transformation

Transformation	Meaning
map (func)	Return a new DStream by passing each element of the source DStream through a function func
flatMap(func)	Similar to map, but each input item can be mapped to 0 or more output items
filter(func)	Return a new DStream by selecting only the records of the source DStream on which func returns true

14

14

Transformation

Transformation	Meaning
count	Return a new DStream of single-element RDDs by counting the number of elements in each RDD of the source DStream
countbyValue	Returns a new DStream of (K, Long) pairs where the value of each key is its frequency in each RDD of the source DStream.
reduce(func)	Return a new DStream of single-element RDDs by aggregating the elements in each RDD of the source DStream using a function func (which takes two arguments and returns one).
reduceByKey(func)	When called on a DStream of (K, V) pairs, return a new DStream of (K, V) pairs where the values for each key are aggregated using the given reduce function

15

15

Transformation

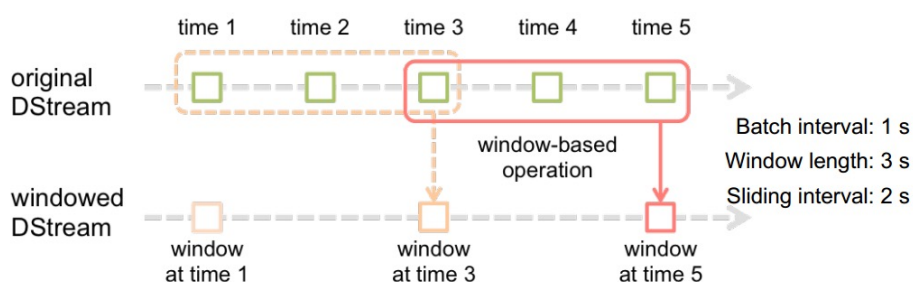
Transformation	Meaning
union(otherStream)	Return a new DStream that contains the union of the elements in the source DStream and otherDStream.
join(otherStream)	When called on two DStreams of (K, V) and (K, W) pairs, return a new DStream of (K, (V, W)) pairs with all pairs of elements for each key.

16

16

Window Operations

- Spark provides a set of transformations that apply to a sliding window of data
- A window is defined by: **window length** and **sliding interval**



17

17

Window Operations

- `window(windowLength, slideInterval)`
 - Returns a new DStream which is computed based on windowed batches
- `countByWindow(windowLength, slideInterval)`
 - Returns a sliding window count of elements in the stream.
- `reduceByWindow(func, windowLength, slideInterval)`
 - Returns a new single-element DStream, created by aggregating elements in the stream over a sliding interval using func.

18

18

Output Operation

- Push out DStream's data to external systems, e.g., a database or a file system

Operation	Meaning
print	Prints the first ten elements of every batch of data in a DStream on the driver node running the application
saveAsTextFiles	Save this DStream's contents as text files
saveAsHadoopFiles	Save this DStream's contents as Hadoop files.
foreachRDD(func)	Applies a function, func, to each RDD generated from the stream

19

19

Example

Word Count

```
val context = new StreamingContext(conf, Seconds(1))
val lines = context.socketTextStream(...)
val words = lines.flatMap(_.split(" "))
val wordCounts = words.map(x => (x, 1)).reduceByKey(_+_ )
wordCounts.print()
context.start()
```

Print the DStream contents on screen

Start the streaming job

 databricks

20

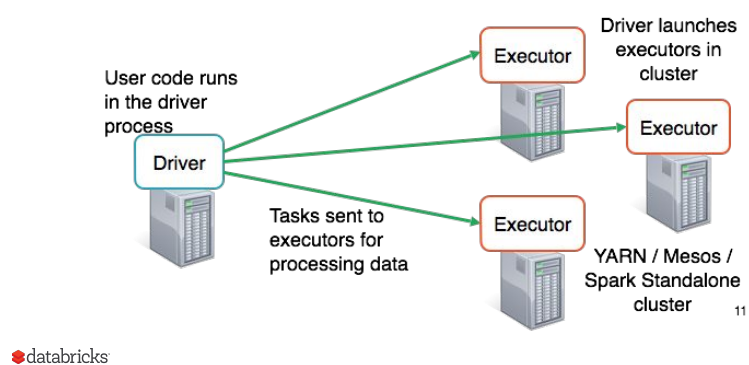
20

Lifecycle of a streaming app

21

21

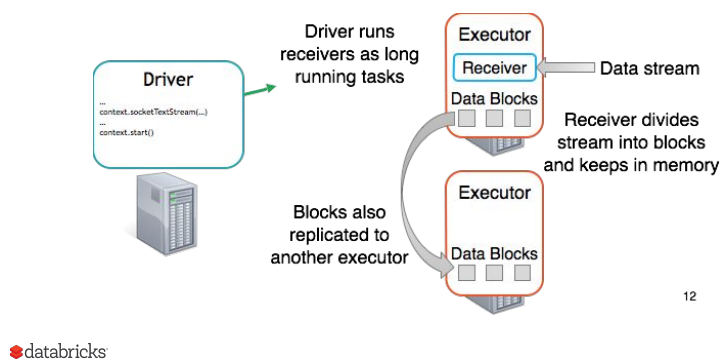
Execution in any Spark Application



22

22

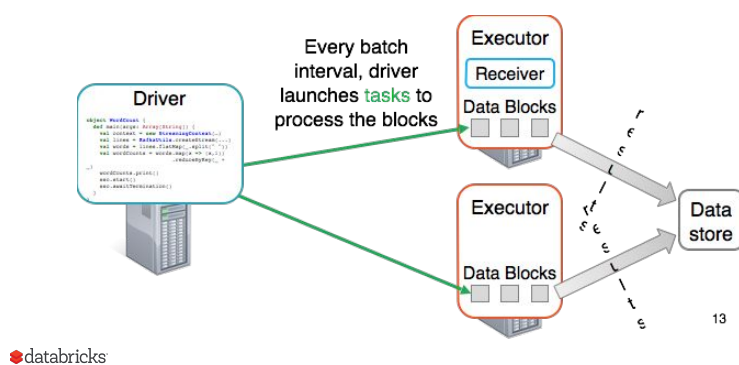
Execution in Spark Streaming: Receiving data



23

23

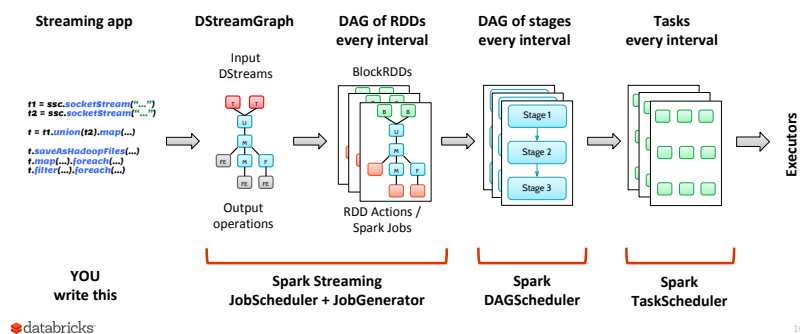
Execution in Spark Streaming: Processing data



24

24

End-to-end view

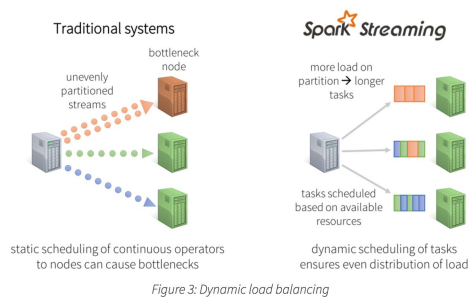


16

25

25

Dynamic Load Balancing



databricks

26

26

Fast failure and Straggler recovery

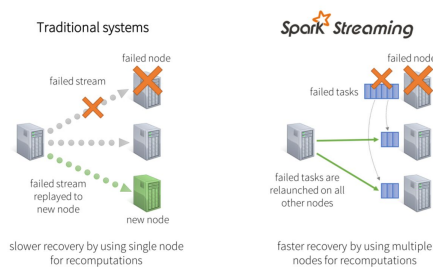


Figure 4: Faster failure recovery with redistribution of computation



27

27

Acknowledgement and References

Books:

- Holden Karau, Andy Konwinski, Patrick Wendell & Matei Zaharia. Learning Spark. Oreilly
- James A. Scott. Getting started with Apache Spark. MapR Technologies

Slides:

- Amir H. Payberah. Scalable Stream Processing – Spark Streaming and Flink
- Matteo Nardelli. Spark Streaming: Hands on Session
- DataBricks. Spark Streaming
- DataBricks: Spark Streaming: Best Practices

28

28