

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Presentation 4 / 75

Recap

Challenges with Real-Time Data

Apache Spark Streaming makes it easier to process real-time data. However, there are challenges:

- It's difficult to handle data in a timely manner.
- Some data is incomplete or inconsistent. This can potentially reduce the quality of your data and processing.
- Data is often unstructured.
- High volume of data.
- Rapidly changing data.
- Data is often noisy.

Why is Kafka Needed?

Let us see how using big data pipelines makes a system more complicated.

Kafka Cluster

Kafka is a distributed, reliable, and fault-tolerant service for efficiently collecting, aggregating, and moving large amounts of log data from various sources to a distributed database.

Why Flume?

Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of log data from various sources to a distributed database.

What is Flume?

Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of log data from various sources to a distributed database.

Flume Architecture : Deep Dive

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:44 IN 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Topics

- Spark Streaming Features
- Spark Workflow
- Streaming Context
- Dstream
- Caching
- Accumulators, Broadcast Variables and Checkpoints
- Streaming Word Count
- Stateful Operators

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:44 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Objectives

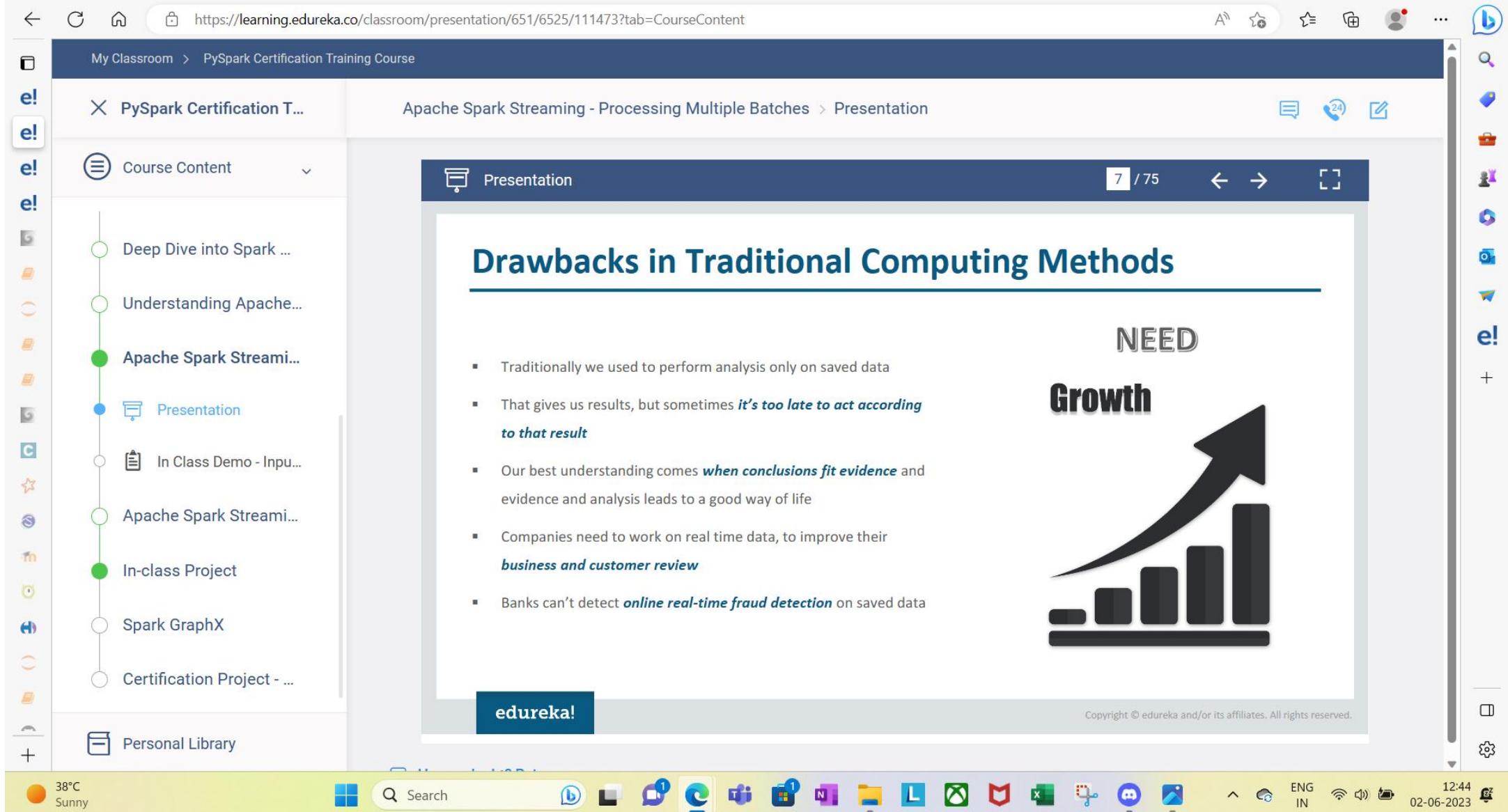
After completing this module, you should be able to:

- Analyze the Drawbacks in Existing Computing Methods
- Describe why Streaming is Necessary
- Explain what is Spark Streaming
- Describe Spark Streaming Features
- Describe Spark Streaming Workflow
- Understand how Uber uses Streaming Data
- Explain Streaming Context and DStreams
- Describe Transformations on DStreams
- Execute a Streaming WordCount Program



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.



X PySpark Certification T...

Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

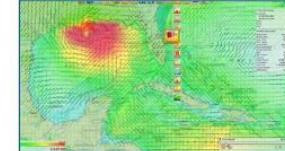
Certification Project - ...

Personal Library

Solution : Spark Streaming



| | | |
|--------|--------|-------|
| 402.13 | +6.84 | +49 |
| 34.18 | -1.95 | -15 |
| 58.84 | +7.82 | +372 |
| 887.32 | +9.88 | +3.96 |
| 573.54 | +14.28 | +254 |
| 532.89 | -11.32 | -213 |
| 401.76 | +9.15 | +15 |



Spark Streaming is used to stream real-time data from various sources like *Twitter*, *Stock Market* and *Geographical Systems* and perform *powerful analytics* to help businesses

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

- Deep Dive into Spark ...
 - Understanding Apache...
 - Apache Spark Streami...
 -  Presentation
 -  In Class Demo - Inpu...
 - Apache Spark Streami...
 - In-class Project
 - Spark GraphX
 - Certification Project - ...

Apache Spark Streaming - Processing Multiple Batches > Presentation



Let us understand what is Streaming

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

X PySpark Certification T...

Apache Spark Streaming - Processing Multiple Batches > Presentation



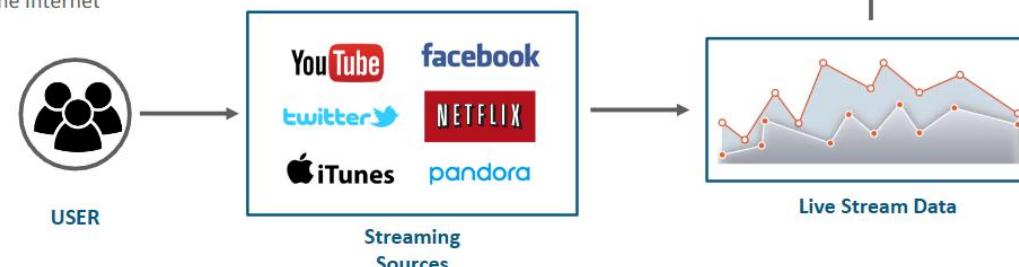
 Course Content

- Deep Dive into Spark ...
 - Understanding Apache...
 - Apache Spark Streami...**
 -  **Presentation**
 -  **In Class Demo - Inpu...**
 - Apache Spark Streami...
 - In-class Project**
 - Spark GraphX
 - Certification Project - ...

What is Streaming?

- **Data Streaming** is a technique for transferring data so that it can be processed as a steady and continuous stream
 - A **data stream** is an unbounded sequence of data arriving continuously
 - Streaming technologies are becoming increasingly important with the growth of the Internet

"Without stream processing there's no big data and no Internet of Things" – Dana Sandu, SQL stream™



My Classroom - My Learning Environment

Apache Spark Streaming - Processing Multiple Batches > Presentation

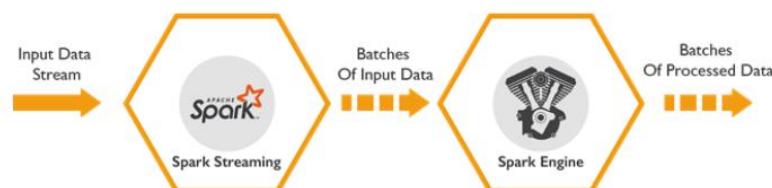


 Course Content

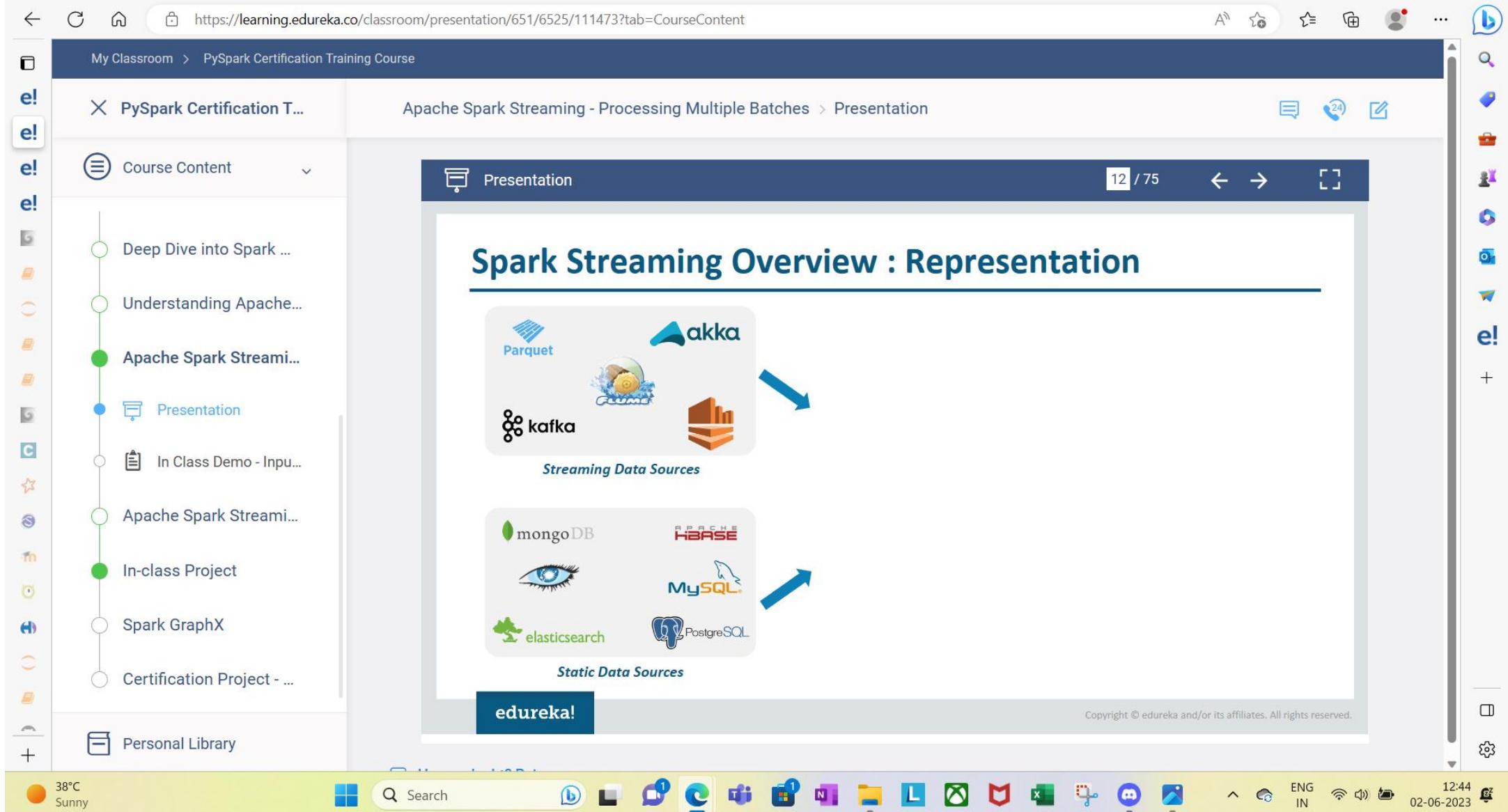
- Deep Dive into Spark ...
 - Understanding Apache...
 - Apache Spark Streami...**
 -  **Presentation**
 -  In Class Demo - Inpu...
 - Apache Spark Streami...
 - In-class Project**
 - Spark GraphX
 - Certification Project - ...

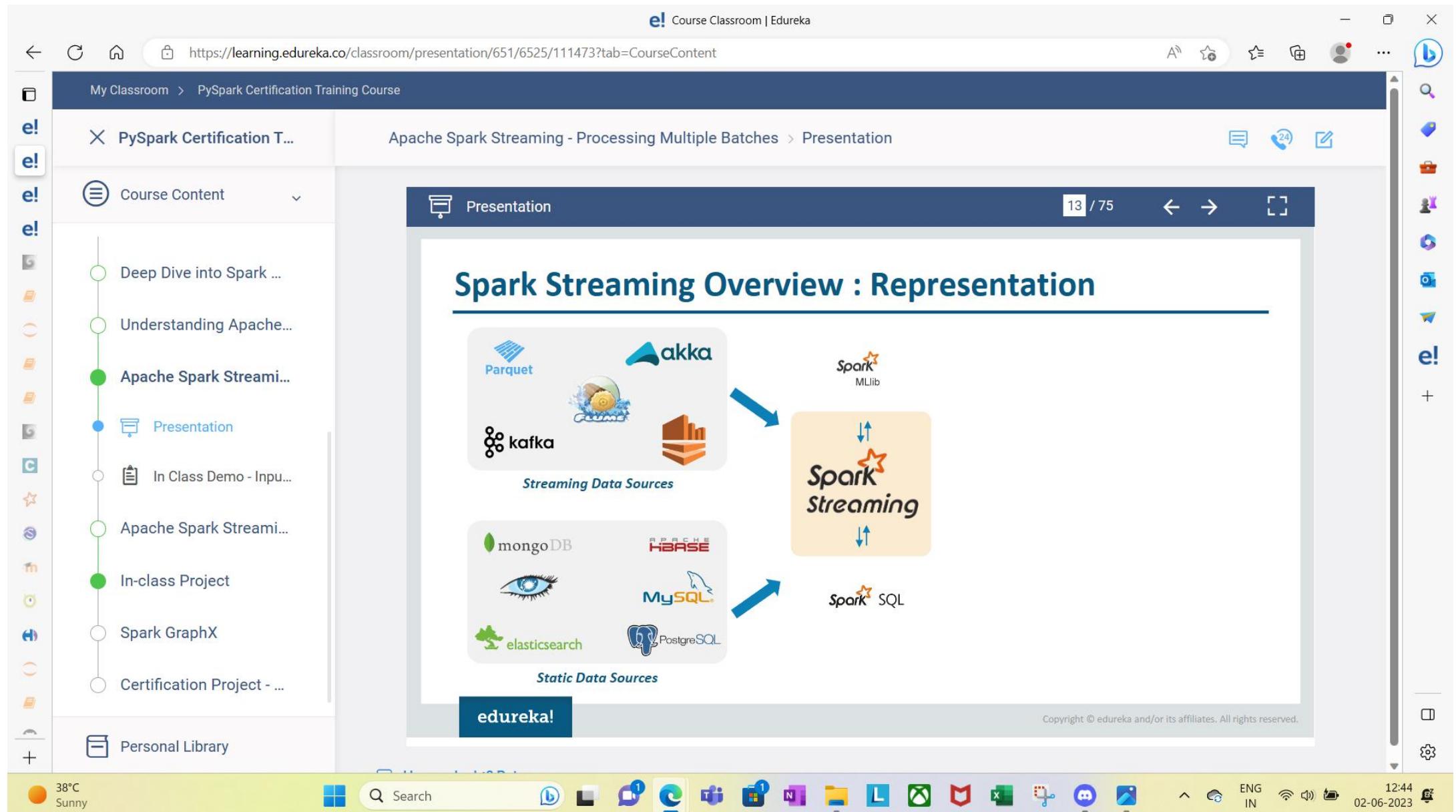
Spark Streaming Overview

- **Spark Streaming** is used for processing *real-time streaming data*
 - It is a useful **addition** to the **core Spark API**
 - Spark Streaming enables **high-throughput** and **fault-tolerant** stream processing of live data streams
 - The fundamental stream unit is **DStream** which is basically a series of RDDs to process the real-time data



edureka!





My Classroom > PySpark Certification Training Course

X PySpark Certification T...

e! Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library



https://learning.edureka.co/course/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:44 02-06-2023

https://learning.edureka.co/classroom/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Companies using Streaming Data

16 / 75

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:44 02-06-2023

https://learning.edureka.co/course/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:44 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...**
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...
- Personal Library

Spark Streaming in Spark Ecosystem

The diagram illustrates the Spark Ecosystem architecture. At the bottom is a thick blue bar labeled "Spark Core Engine". Above it are five teal-colored boxes, each representing a different component:

- Spark SQL (SQL)**: Used for structured data. Can run unmodified hive queries on existing Hadoop deployment.
- Spark Streaming (Streaming)**: Enables analytical and interactive apps for live streaming data.
- MLlib (Machine Learning)**: Machine learning libraries being built on top of Spark.
- GraphX (Graph Computation)**: Graph Computation engine (Similar to Giraph). Combines data-parallel and graph-parallel concepts.
- SparkR (R on Spark)**: Package for R language to enable R-users to leverage Spark power from R shell.

A large blue arrow points downwards from the components to a callout box at the bottom: "The core engine for entire Spark framework. Provides utilities and architecture for other components".

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:44 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

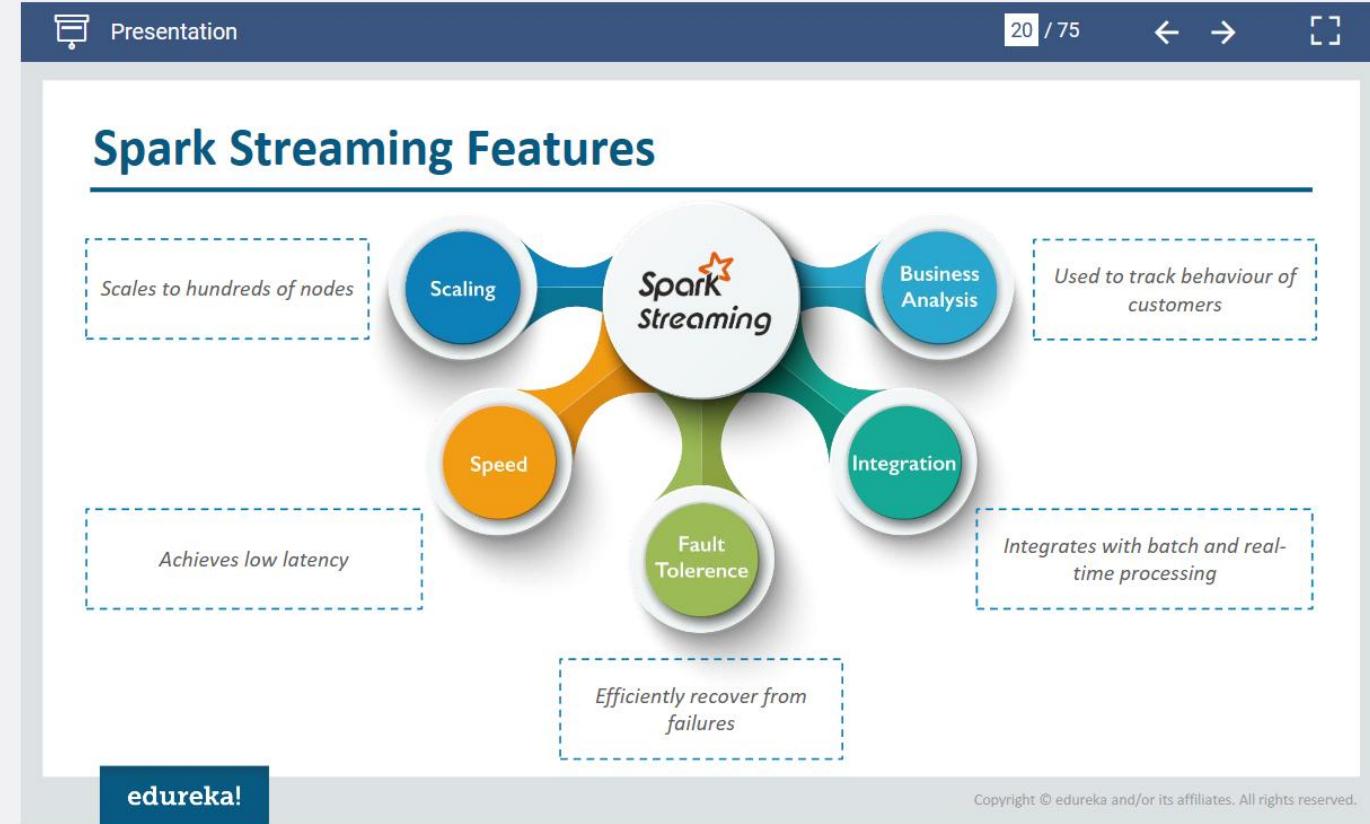
12:44 02-06-2023



Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library



https://learning.edureka.co/course/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:44 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...
- Personal Library

Spark Streaming Workflow

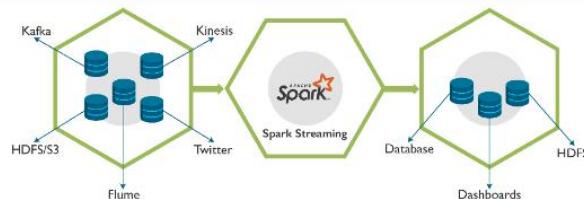


Figure: Data from a variety of sources to various storage systems

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:44 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Input...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...
- Personal Library

Apache Spark Streaming - Processing Multiple Batches > Presentation

Presentation 23 / 75 ← →

Spark Streaming Workflow

Figure: Data from a variety of sources to various storage systems

Figure: Incoming streams of data divided into batches

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:44 02-06-2023



Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Input...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Spark Streaming Workflow

Figure: Data from a variety of sources to various storage systems

Figure: Incoming streams of data divided into batches

Figure: Input data stream divided into discrete chunks of data

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.



Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Input...

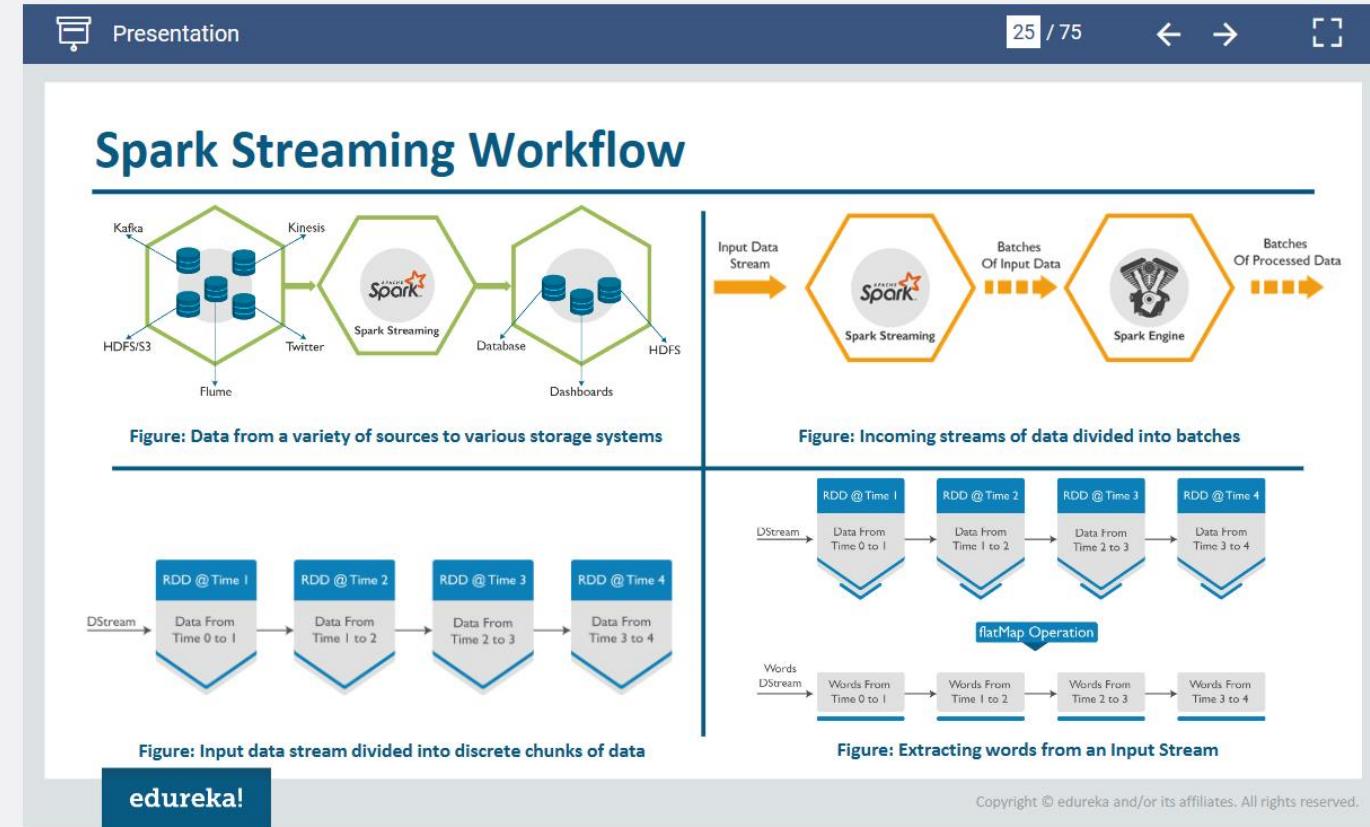
Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library



Copyright © edureka and/or its affiliates. All rights reserved.



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

- Deep Dive into Spark ...
 - Understanding Apache...
 - Apache Spark Streami...
 -  Presentation
 -  In Class Demo - Inpu...
 - Apache Spark Streami...
 - In-class Project
 - Spark GraphX
 - Certification Project - ...

Apache Spark Streaming - Processing Multiple Batches > Presentation



Let us have a look at how Streaming is done at Uber

edureka!

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

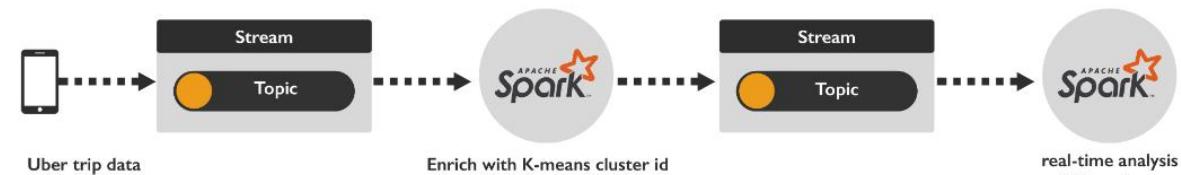
Certification Project - ...

Personal Library



Uber Streaming : Process

- Uber trip data is published to a MapR Streams topic using the Kafka API
- A Spark streaming application subscribed to the first topic:
 - *Ingests a stream of uber trip events*
 - *Identifies the location cluster corresponding to the latitude and longitude of the uber trip*
 - *Adds the cluster location to the event and publishes the results in JSON format to another topic*
- A Spark streaming application subscribed to the second topic:
 - *Analyses the uber trip location clusters that are popular by date and time*



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

X PySpark Certification T...

Apache Spark Streaming - Processing Multiple Batches > Presentation



Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Presentation

28 / 75

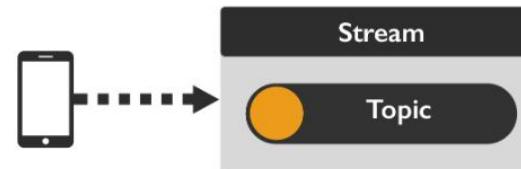


- ## Getting Data
- The example data set is Uber trip data
 - The incoming data is in CSV format, an example is shown below , with the header:

date/time, latitude, longitude, base

2014-08-01 00:00:00,40.729,-73.9422,B02598

Uber trip data



2014-08-01 00:00:00,40.729,-73.9422,B02598

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

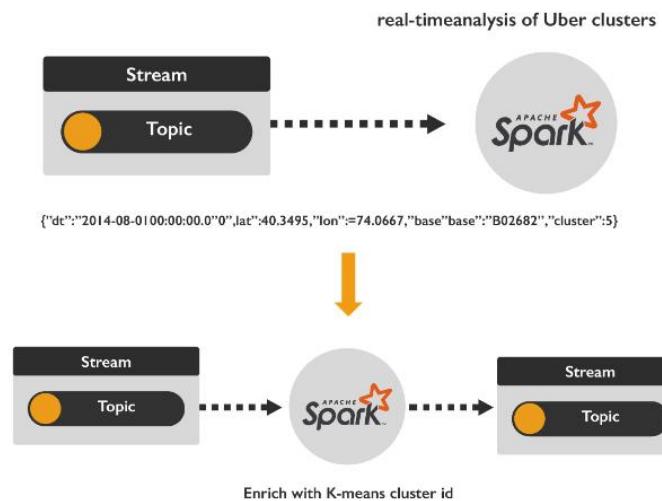
Spark GraphX

Certification Project - ...

Personal Library

Workflow

The enriched Data Records are in JSON format. An example line is shown below :



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

https://learning.edureka.co/classroom/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

L

12:45 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Input...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

```
graph LR; SC((1) Streaming Context) --- D((2) DStream); D --- C((3) Caching); D --- ABBC((4) Accumulators, Broadcast Variables and Checkpoints); D --- IT((2.1) Input DStream); D --- DT((2.2) DStream Transformations); D --- OD((2.3) Output DStream);
```

31 / 75

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Input...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

```
graph LR; A([Steaming Context 1]) --- B[DStream 2]; B --- C1[Input DStream 2.1]; B --- C2[DStream Transformations 2.2]; B --- C3[Output DStream 2.3]; B --- D[Caching 3]; D --- E[Accumulators, Broadcast Variables and Checkpoints 4]
```

38°C Sunny

Search

12:45 02-06-2023



Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Streaming Context

- It is the main **entry point** for Spark functionality
- Spark provides a number of default implementations of sources like **Twitter**, **Akka Actor** and **ZeroMQ** that are accessible from the context
- Whatever we do in Spark Streaming has to start from creating an instance of **StreamingContext**



- Consumes a stream of data in Spark
- Registers an **InputDStream** to produce a **Receiver** object
- It provides methods used to create DStreams from various input sources

Copyright © edureka and/or its affiliates. All rights reserved.

X PySpark Certification T...

Apache Spark Streaming - Processing Multiple Batches > Presentation



e! Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Presentation

34 / 75



Streaming Context : Initialization

- Streaming Context can be created in 3 ways:
 - By providing a *Spark master URL* and *an appName*, or
 - from a *org.apache.spark.SparkConf* configuration, or
 - from an *existing org.apache.spark.SparkContext*
- The associated SparkContext can be accessed using *context.sparkContext*. After creating and transforming DStreams, the streaming computation can be started and stopped using *context.start()* and *context.stop()*, respectively
- Let's see a simple initialization of StreamingContext :

```
#  Spark
from pyspark import SparkContext
#  Spark Streaming
from pyspark.streaming import StreamingContext
sc = SparkContext(appName="PythonSparkStreaming")
ssc = StreamingContext(sc, 60)
```

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

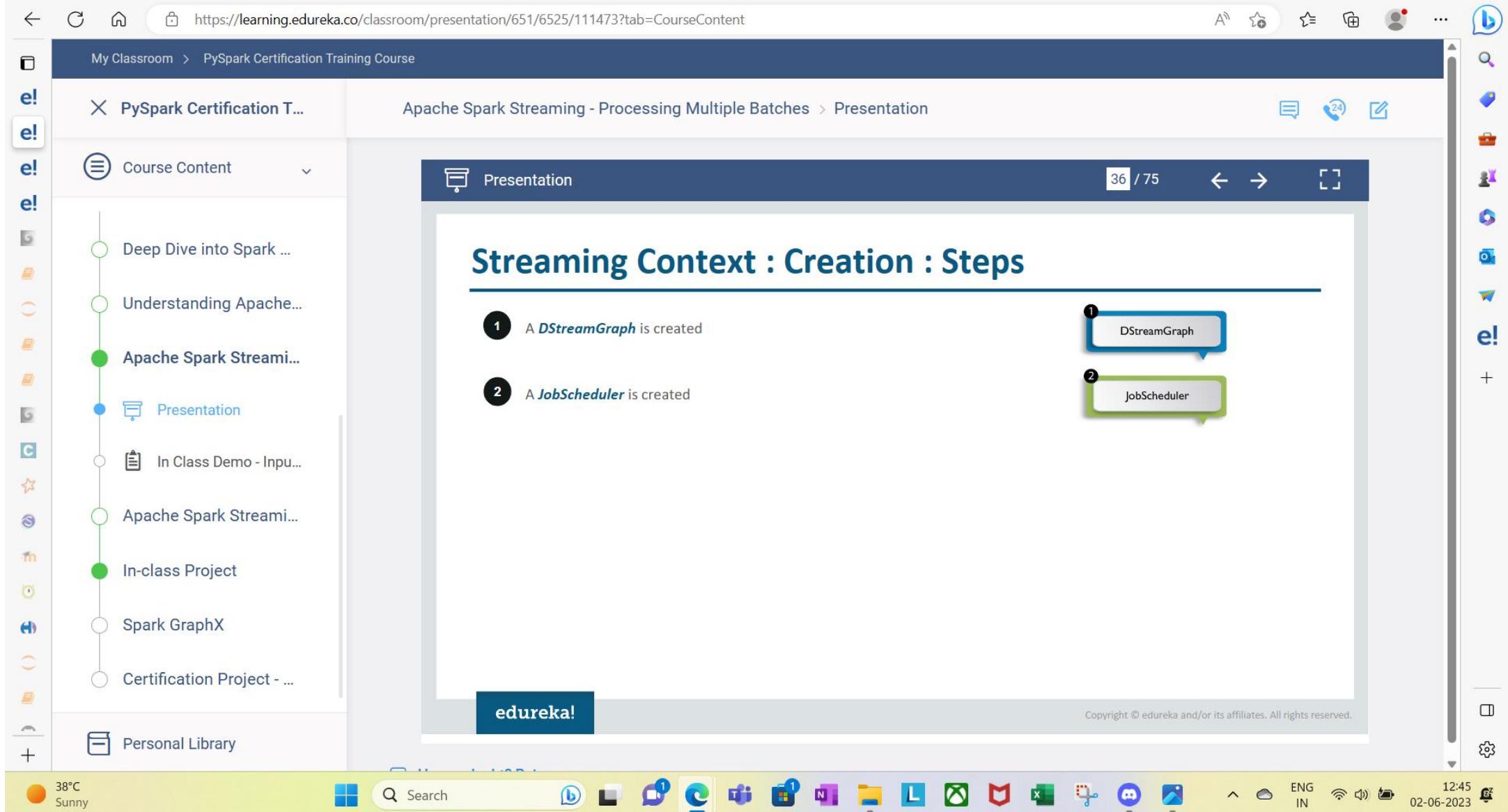
38°C Sunny

Search

12:45 02-06-2023

1 A *DStreamGraph* is created

1 DStreamGraph



My Classroom > PySpark Certification Training Course



X PySpark Certification T...

Apache Spark Streaming - Processing Multiple Batches > Presentation



Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library



Presentation

37 / 75



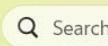
Streaming Context : Creation : Steps

- 1 A *DStreamGraph* is created
- 2 A *JobScheduler* is created
- 3 A *StreamingJobProgressListener* is created



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.



ENG

IN



12:45

02-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

e! Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Apache Spark Streaming - Processing Multiple Batches > Presentation

Presentation 38 / 75 ← →

Streaming Context : Creation : Steps

- 1 A *DStreamGraph* is created
- 2 A *JobScheduler* is created
- 3 A *StreamingJobProgressListener* is created
- 4 Streaming tab in web UI is created (when *spark.ui.enabled* is enabled)

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

e! Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Apache Spark Streaming - Processing Multiple Batches > Presentation

Presentation 39 / 75 ← →

Streaming Context : Creation : Steps

- 1 A *DStreamGraph* is created
- 2 A *JobScheduler* is created
- 3 A *StreamingJobProgressListener* is created
- 4 Streaming tab in web UI is created (when `spark.ui.enabled` is enabled)
- 5 A *StreamingSource* is instantiated

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Streaming Context : Creation : Steps

- 1 A *DStreamGraph* is created
- 2 A *JobScheduler* is created
- 3 A *StreamingJobProgressListener* is created
- 4 Streaming tab in web UI is created (when `spark.ui.enabled` is enabled)
- 5 A *StreamingSource* is instantiated
- 6 At this point, StreamingContext enters *INITIALIZED* state

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023

https://learning.edureka.co/classroom/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Input...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

```
graph LR; SC((1 Steaming Context)) --- D((2 DStream)); D --- ID((2.1 Input DStream)); D --- DT((2.2 DStream Transformations)); D --- OD((2.3 Output DStream)); C((3 Caching)) --- D; AV((4 Accumulators, Broadcast Variables and Checkpoints)) --- D;
```

42 / 75

38°C Sunny 12:45 02-06-2023

https://learning.edureka.co/classroom/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...
- Personal Library

DStream

- Discretized Stream* (DStream) is the basic abstraction provided by Spark Streaming
- It is a *continuous stream* of data

The diagram illustrates the concept of a DStream. It shows a horizontal arrow labeled "DStream" pointing to four vertical arrows labeled "RDD @ Time 1", "RDD @ Time 2", "RDD @ Time 3", and "RDD @ Time 4". Each vertical arrow points to a trapezoidal shape representing a data chunk. The first trapezoid is labeled "Data From Time 0 to 1", the second "Data From Time 1 to 2", the third "Data From Time 2 to 3", and the fourth "Data From Time 3 to 4". Arrows connect the trapezoids sequentially from left to right.

Figure: Input data stream divided into discrete chunks of data

- It is received from *source* or from a *processed data stream* generated by transforming the input stream
- Internally, a DStream is represented by a continuous *series of RDDs* and each *RDD* contains data from a *certain interval*

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

DStream Operation

- Any operation applied on a DStream translates to operations on the underlying RDDs
- For example, in the example of converting a stream of lines to words, the `flatMap` operation is applied on each RDD in the lines DStream to generate the RDDs of the words **DStream**

flatMap Operation

Figure: Extracting words from an Input Stream

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Input...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

```
graph LR; SC((1) Streaming Context) --> D((2) DStream); D --> C((3) Caching); D --> ABBC((4) Accumulators, Broadcast Variables and Checkpoints); D --> ID((2.1) Input DStream); D --> DT((2.2) DStream Transformations); D --> ODS((2.3) Output DStream);
```

45 / 75

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023



Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Input DStreams

- *Input DStreams* are DStreams representing the stream of *input data* received from streaming *sources*

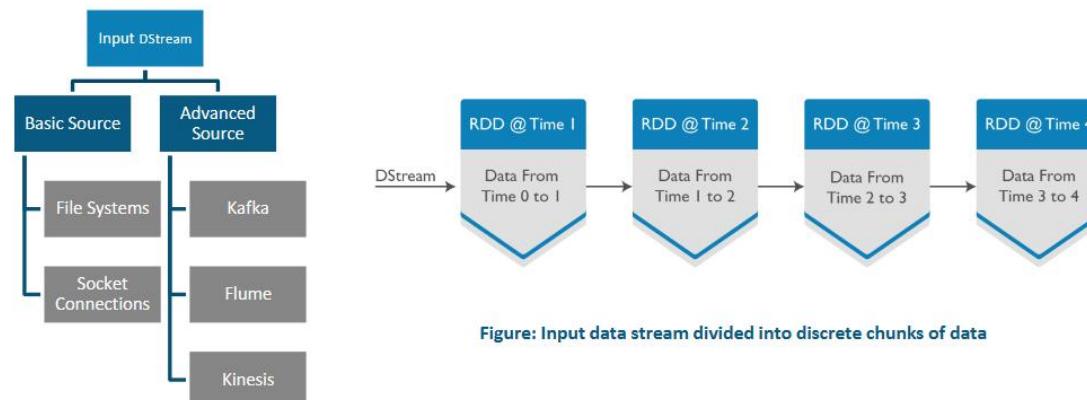


Figure: Input data stream divided into discrete chunks of data

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

Apache Spark Streaming - Processing Multiple Batches > Presentation



 Course Content

- Deep Dive into Spark ...
 - Understanding Apache...
 - Apache Spark Streami...
 -  Presentation
 -  In Class Demo - Inpu...
 - Apache Spark Streami...
 - In-class Project
 - Spark GraphX
 - Certification Project - ...

Receiver

- Every input DStream is associated with a **Receiver** object which receives the data from a **source** and stores it in **Spark's memory** for processing

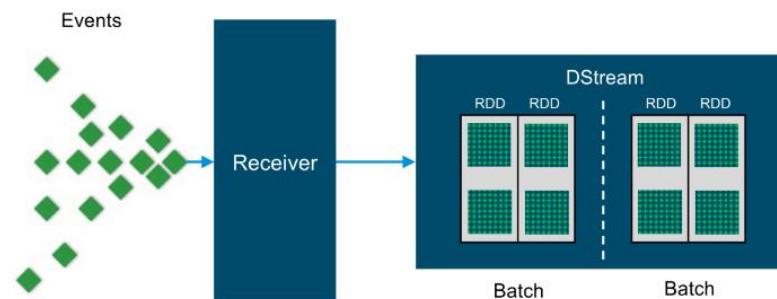


Figure: The Receiver sends data onto the DStream where each Batch contains RDDs

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Input...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

```
graph LR; SC((1 Steaming Context)) --- D((2 DStream)); D --- IC((2.1 Input DStream)); D --- DT((2.2 DStream Transformations)); D --- OD((2.3 Output DStream)); C((3 Caching)) --- AV((4 Accumulators, Broadcast Variables and Checkpoints))
```

48 / 75

38°C Sunny 12:45 02-06-2023

Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Presentation

49 / 75

Transformations
on DStreamsMost Popular Spark
Streaming Transformations

- 1 map
- 2 flatMap
- 3 filter
- 4 reduce
- 5 groupBy

- **Transformations** allow the data from the **input DStream** to be **modified** similar to RDDs
- DStreams support many of the transformations available on normal Spark RDDs

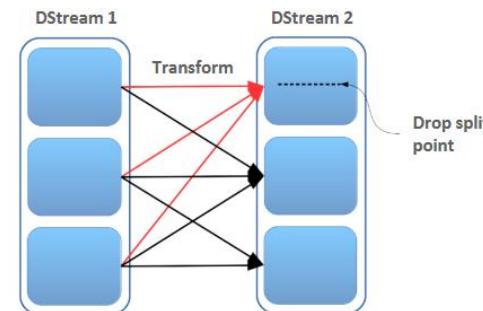


Figure: DStream Transformations

Copyright © edureka and/or its affiliates. All rights reserved.

Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Input...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Presentation

50 / 75



Transformations on DStreams

1 map

2 flatMap

3 filter

4 reduce

5 groupBy

- **map(func)**

- **map(func)** returns a new **DStream** by passing each element of the source **DStream** through a function **func**



Figure: Input DStream being converted through `map(func)`

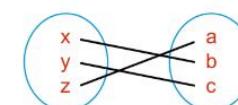


Figure: Map Function

Copyright ©, edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Transformations on DStreams

- 1 map
- 2 flatMap
- 3 filter
- 4 reduce
- 5 groupBy

flatMap(func)

flatMap(func) is similar to map(func) but each input item can be mapped to 0 or more output items and returns a new DStream by passing each source element through a function func

Input Data Stream → flatMap → Batches Of Input Data → Node

Figure: Input DStream being converted through flatMap(func)

Figure: flatMap Function

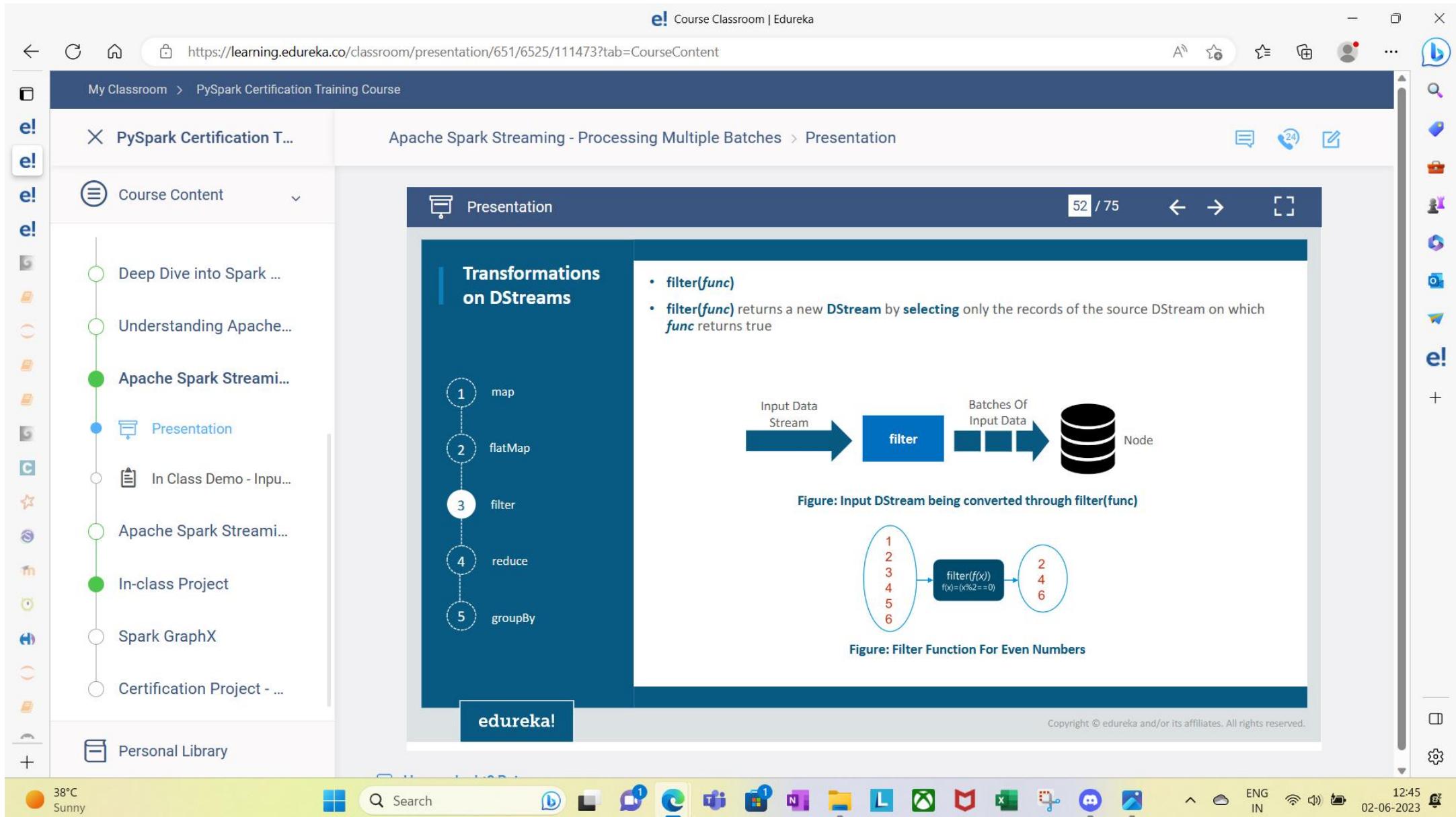
edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023



My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Transformations on DStreams

- map
- flatMap
- filter
- reduce
- groupBy

• `reduce(func)`

• `reduce(func)` returns a new DStream of **single-element RDDs** by **aggregating** the elements in each RDD of the source DStream using a function `func`

Input Data Stream → reduce → Batches Of Input Data → Node

Figure: Input DStream being converted through `reduce(func)`

Figure: Reduce Function To Get Cumulative Sum

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 IN 02-06-2023



Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Presentation

54 / 75

Transformations
on DStreams

map

flatMap

filter

reduce

groupBy

• groupBy(func)

- groupBy(func) returns the new RDD which basically is made up with a key and corresponding list of items of that group



Figure: Input DStream being converted through groupBy(func)

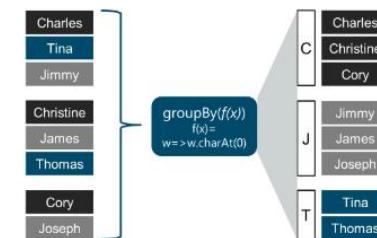
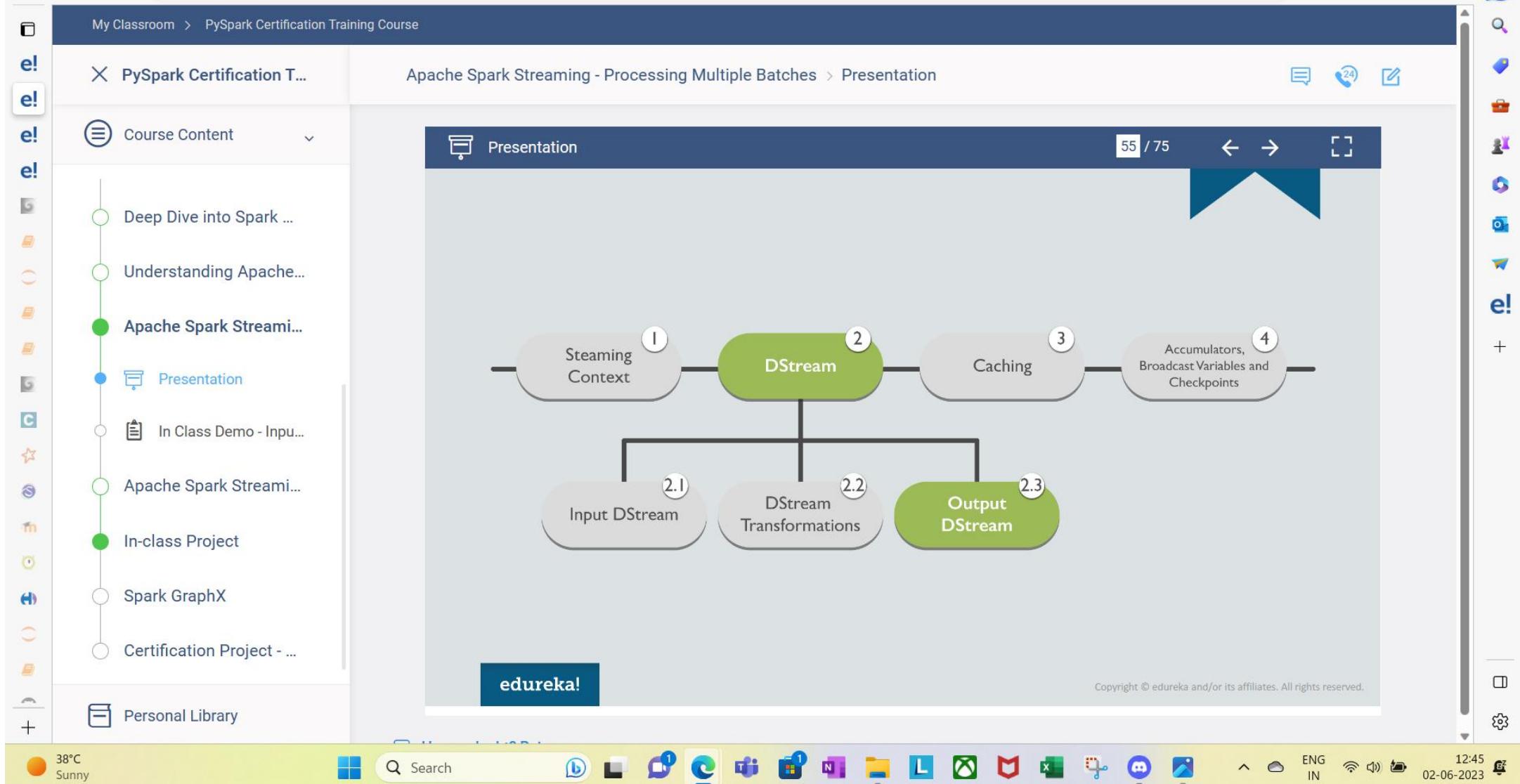


Figure: Grouping By First Letters

Copyright © edureka and/or its affiliates. All rights reserved.



Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Output Operations on DStreams

- **Output operations** allow DStream's data to be pushed out to **external systems** like **databases** or **file systems**
- Output operations **trigger** the **actual execution** of all the DStream transformations

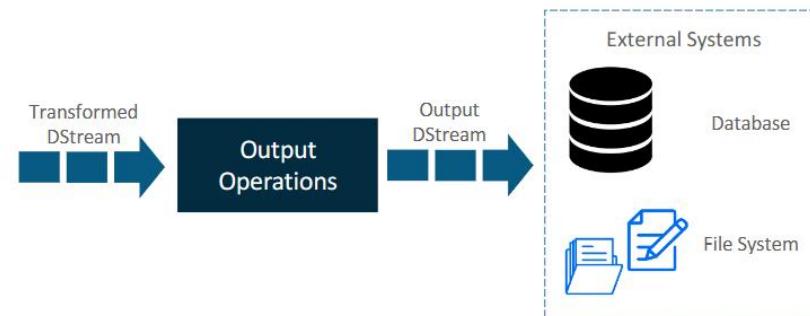


Figure: Output Operations on DStreams

Copyright © edureka and/or its affiliates. All rights reserved.

edureka!

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

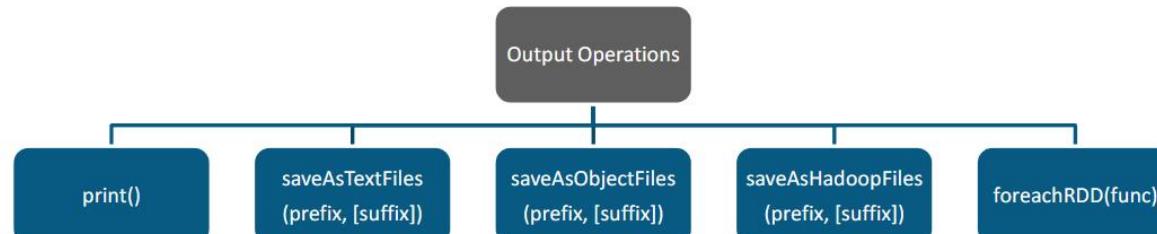
Spark GraphX

Certification Project - ...

Personal Library

Output Operations on DStreams

- Currently, the following output operations are defined:



My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Output Operations Example - `foreachRDD`

- `foreachRDD`
- `dstream.foreachRDD` is a powerful primitive that allows data to be sent out to external systems
- The lazy evaluation achieves the most efficient transfer of data

```
def sendRecord(record):  
    connection = createNewConnection()  
    connection.send(record)  
    connection.close()  
dstream.foreachRDD(lambda rdd: rdd.foreach(sendRecord))
```

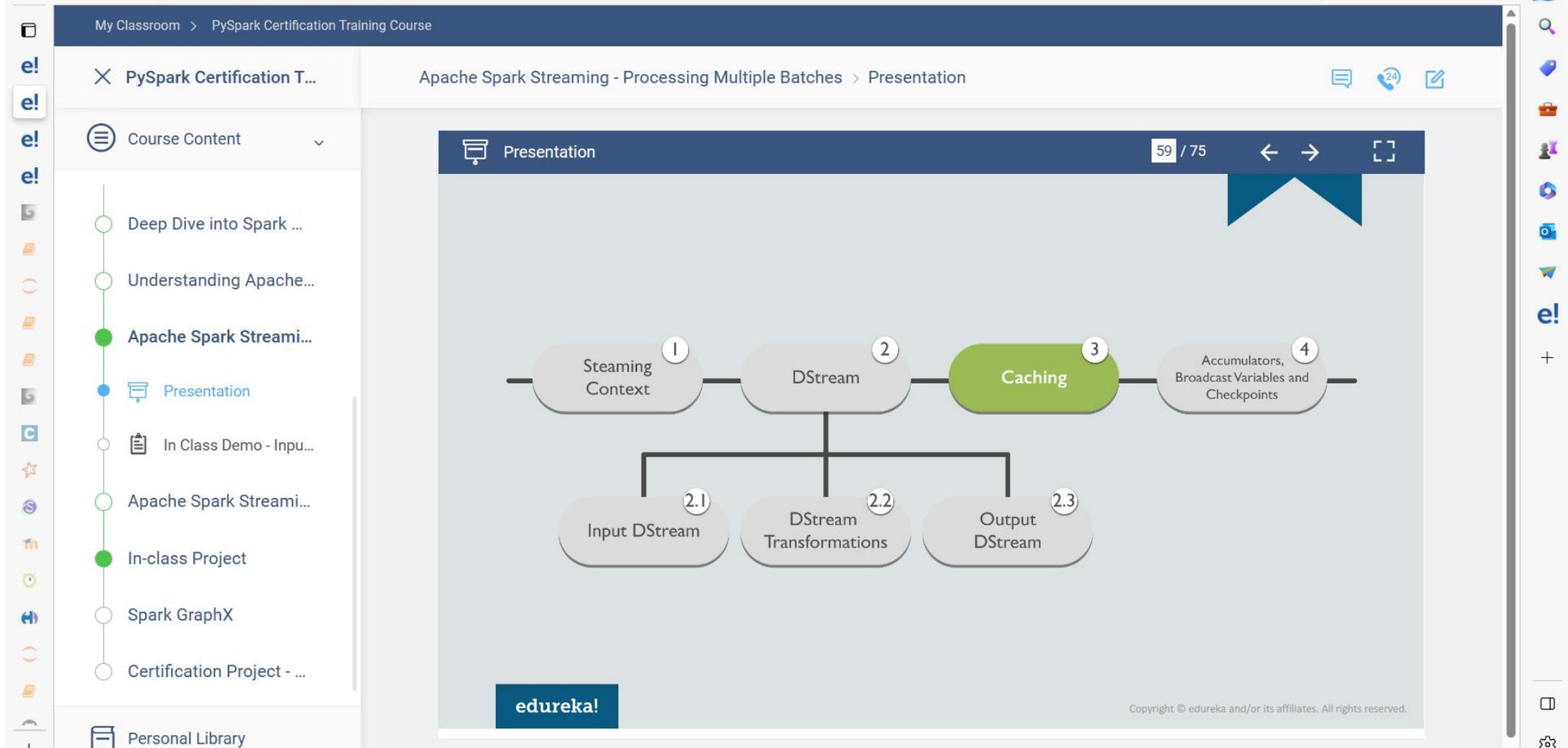
edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023



https://learning.edureka.co/classroom/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

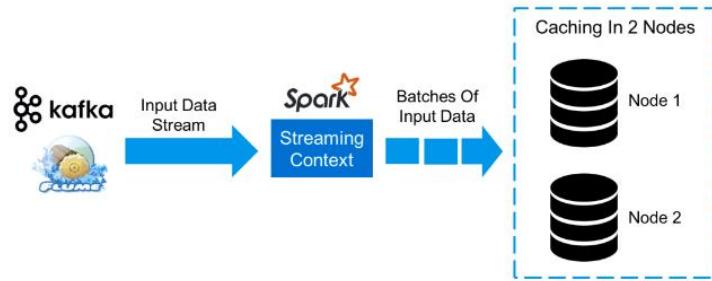
PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Caching/Persistence

- DStreams allow developers to **cache/persist** the stream's data in memory
- This is useful if the data in the DStream will be computed multiple times
- This can be done using the **persist()** method on a DStream
- For input streams that receive data over the network (such as Kafka, Flume, Sockets, etc.), the default persistence level is set to replicate the data to two nodes for fault-tolerance


Figure: Caching Into 2 Nodes

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Input...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

```
graph LR; SC((1) Steaming Context) --- D((2) DStream); D --- C((3) Caching); C --- ABBC((4) Accumulators, Broadcast Variables and Checkpoints); D --- IT((2.1) Input DStream); D --- DT((2.2) DStream Transformations); D --- OT((2.3) Output DStream)
```

61 / 75

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023

Course Content

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Accumulators

- **Accumulators** are variables that are only added through an associative and commutative operation
- They are used to implement **counters** or **sums**
- Tracking accumulators in the UI can be useful for **understanding** the **progress** of running stages
- Spark natively supports **numeric** accumulators. We can create **named** or **unnamed** accumulators

| Accumulators | | | | | | | | |
|--------------|-------|---------|---------|----------------|--------------------|---------------------|----------|-------------|
| Accumulable | Value | | | | | | | |
| counter | 45 | | | | | | | |
| Tasks | | | | | | | | |
| Index | ID | Attempt | Status | Locality Level | Executor ID / Host | Launch Time | Duration | GC Time |
| 0 | 0 | 0 | SUCCESS | PROCESS_LOCAL | driver / localhost | 2016/04/21 10:10:41 | 17 ms | |
| 1 | 1 | 0 | SUCCESS | PROCESS_LOCAL | driver / localhost | 2016/04/21 10:10:41 | 17 ms | counter: 1 |
| 2 | 2 | 0 | SUCCESS | PROCESS_LOCAL | driver / localhost | 2016/04/21 10:10:41 | 17 ms | counter: 2 |
| 3 | 3 | 0 | SUCCESS | PROCESS_LOCAL | driver / localhost | 2016/04/21 10:10:41 | 17 ms | counter: 7 |
| 4 | 4 | 0 | SUCCESS | PROCESS_LOCAL | driver / localhost | 2016/04/21 10:10:41 | 17 ms | counter: 5 |
| 5 | 5 | 0 | SUCCESS | PROCESS_LOCAL | driver / localhost | 2016/04/21 10:10:41 | 17 ms | counter: 6 |
| 6 | 6 | 0 | SUCCESS | PROCESS_LOCAL | driver / localhost | 2016/04/21 10:10:41 | 17 ms | counter: 7 |
| 7 | 7 | 0 | SUCCESS | PROCESS_LOCAL | driver / localhost | 2016/04/21 10:10:41 | 17 ms | counter: 17 |

Copyright © edureka and/or its affiliates. All rights reserved.

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Presentation

In Class Demo - Inpu...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library



Broadcast Variables

- **Broadcast variables** allow the programmer to keep a **read-only variable cached on each machine** rather than shipping a copy of it with tasks
- They can be used to give every node a **copy** of a **large input dataset** in an efficient manner
- **Spark** also attempts to distribute broadcast variables using efficient **broadcast algorithms** to reduce communication cost

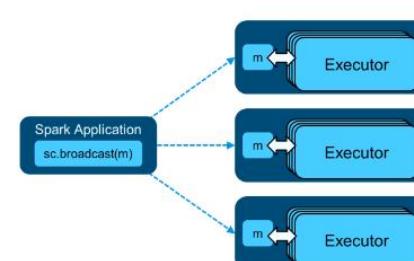


Figure: Broadcasting A Value To Executors

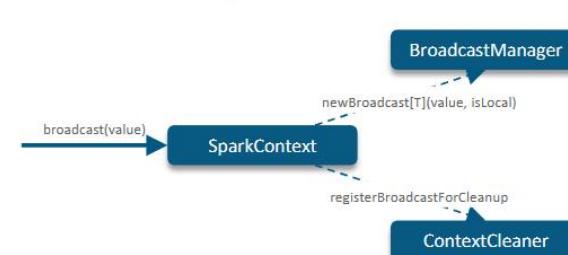
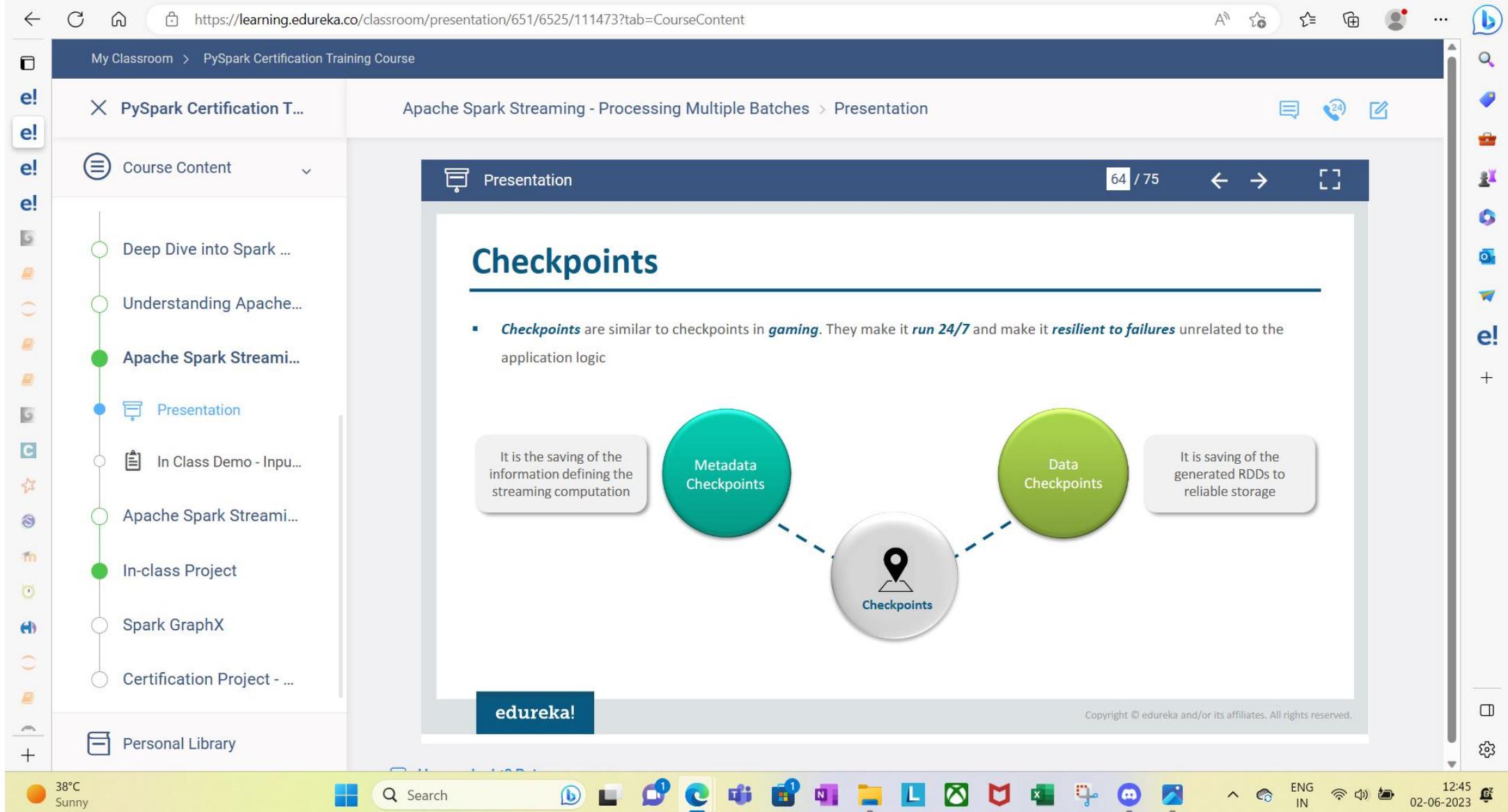


Figure: SparkContext and Broadcasting

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.



https://learning.edureka.co/course/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Streaming WordCount

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

L

12:45 02-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

e! Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Apache Spark Streaming - Processing Multiple Batches > Presentation

Presentation 66 / 75 ← → ⟲ ⟳

Streaming WordCount

```
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
# Create a local StreamingContext with two working thread and batch interval of 1 second
sc = SparkContext("local[2]", "NetworkWordCount")
ssc = StreamingContext(sc, 1)
# Create a DStream that will connect to hostname:port, like localhost:9999
lines = ssc.socketTextStream("localhost", 9999)
#split each line into words
words = lines.flatMap(lambda line: line.split(" "))
```

Importing StreamingContext

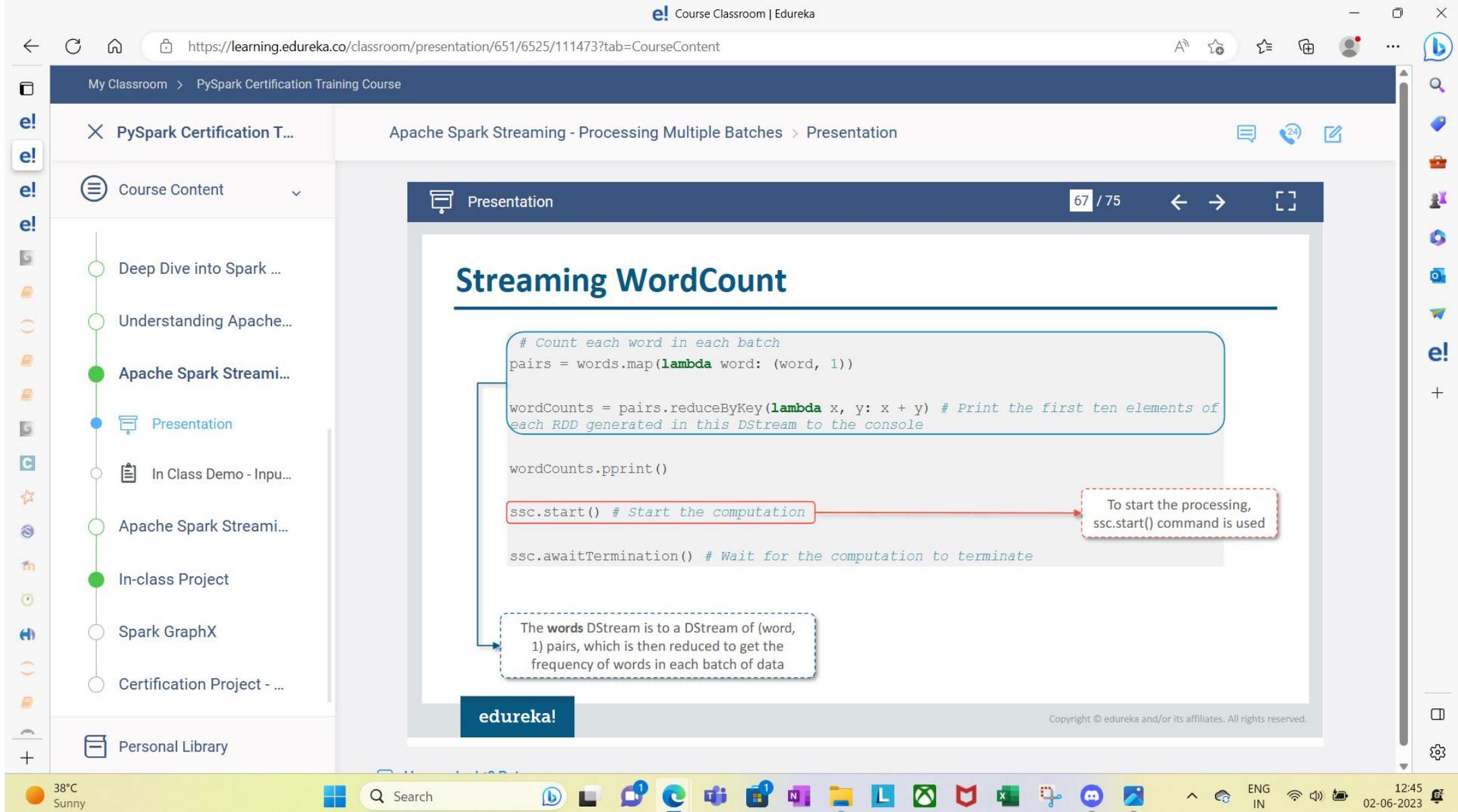
Creating a local StreamingContext with two execution threads, and batch interval of 1 second

The lines DStream represents the stream of data that will be received from the data server

Each line will be split into multiple words and the stream of words is represented as the words DStream

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.



https://learning.edureka.co/classroom/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

Windowed Operators

In short, **windowed operators** allow you to apply transformations over a **sliding window** of data

There are mainly 3 types of operators :

The diagram illustrates the concept of Windowed Operators. At the center is a blue circle labeled "Windowed Operators". Three arrows point from this central circle to three surrounding circles: a green circle labeled "Window", an orange circle labeled "Slice", and a teal circle labeled "reduceByWindow".

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023

https://learning.edureka.co/course/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

L

12:45 02-06-2023

https://learning.edureka.co/classroom/presentation/651/6525/111473?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation**
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

Stateful Operators

- Stateful operators** (like `mapWithState` or `updateStateByKey`) are part of the set of additional operators available on DStreams of *key-value pairs*, i.e. instances of `DStream[(K, V)]`
- They allow you to build stateful stream processing pipelines and are also called cumulative calculations
- By design streaming operators are *stateless* and know nothing about the previous records and hence a state
- If we'd like to react to new records appropriately given the previous records we would have to resort to using **persistent storages** outside Spark Streaming
- Few Examples of Stateful Operators are:
 - `mapWithState`
 - `StateSpec`
 - `updateStateByKey`

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 IN 02-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...
- Personal Library

71 / 75 ← →

Summary

What is Streaming?

- Data streaming is a technique for processing data as it arrives by providing a timely and efficient mechanism.
- It makes it easier to work with real-time data and analyze it in a timely manner.
- Stream processing engines handle stream processing in a distributed environment.

Streaming Sources: YouTube, Facebook, Twitter, LinkedIn, etc. → Processing Stations → Live Stream Data

Spark Streaming Overview : Representation

Spark Streaming in Spark Ecosystem

Spark Streaming Workflow

Streaming WordCount

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 IN 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming - Processing Multiple Batches > Presentation

Course Content

- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Presentation
- In Class Demo - Inpu...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Further Reading

- PySpark Tutorial – Learn Apache Spark using Python
 - <https://www.edureka.co/blog/pyspark-tutorial/>
- Apache Spark with Hadoop
 - <https://www.edureka.co/blog/apache-spark-with-hadoop-why-it-matters/>

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

38°C Sunny

Search

12:45 02-06-2023