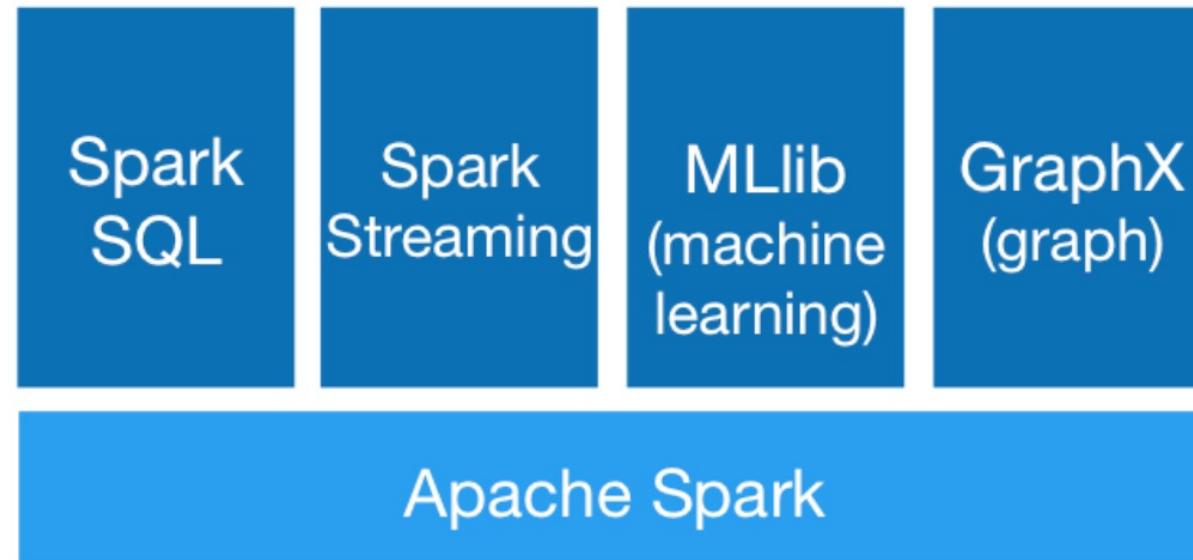


# Spark Streaming Introduction

# Spark Streaming makes it easy to build scalable fault-tolerant streaming applications



# Spark Streaming makes it easy to build scalable fault-tolerant streaming applications

## Ease of Use

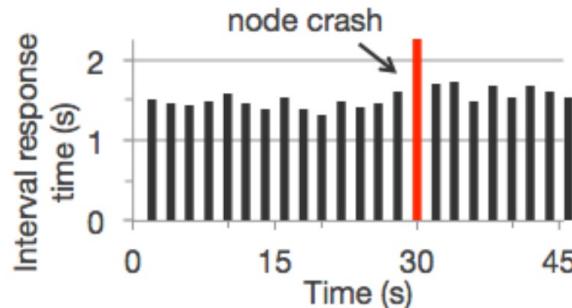
Build applications through high-level operators.

```
TwitterUtils.createStream(...)  
  .filter(_.getText.contains("Spark"))  
  .countByWindow(Seconds(5))
```

Counting tweets on a sliding window

## Fault Tolerance

Stateful exactly-once semantics out of the box.



## Spark Integration

Combine streaming with batch and interactive queries.

```
stream.join(historicCounts).filter {  
  case (word, (curCount, oldCount)) =>  
    curCount > oldCount  
}
```

Find words with higher frequency than historic data

# Alternatives



Process unbounded streams of data in real-time manner and used for real-time analytics, online machine learning, continuous computation, distributed RPC, ETL, and more.



# Alternatives



Apache Flink® is an open-source stream processing framework for distributed, high-performing, always-available, and accurate data streaming applications.



# Spark Streaming overview



- Spark Streaming can read data from HDFS, Flume, Kafka, Twitter, ZeroMQ and TCP sockets.
- You can also define your own custom data sources.

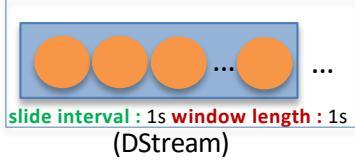
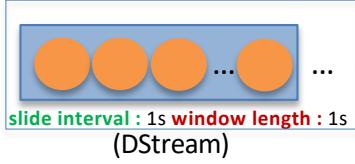
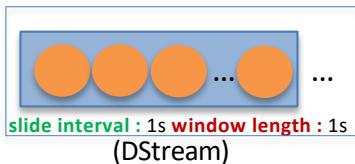
# Spark Streaming API Programming model



Spark Streaming API for Python: `pyspark.streaming`

**DStream** object represents the streaming data

**DStream** type provides transformation and actions like RDD



`.window(4, 1)`

```
1  from pyspark import SparkContext
2  from pyspark.streaming import StreamingContext
3
4
5  # 1. Create StreamingContext
6  ssc = StreamingContext(1)
7
8  # 2. Connect to the data source
9  lines = ssc.socketTextStream("devenv", 9999)
10
11 # 3. Analysis with transformans and actions
12 words = lines.flatMap(f1)
13 pairs = words.map(f2)
14 word_counts = pairs.reduceByKey(f3)
15
16 word_counts.pprint(30)
17
18 # 4. Start the computation
19 ssc.start()
20 ssc.awaitTermination()
```

# Spark Streaming Run-time Execution Flow

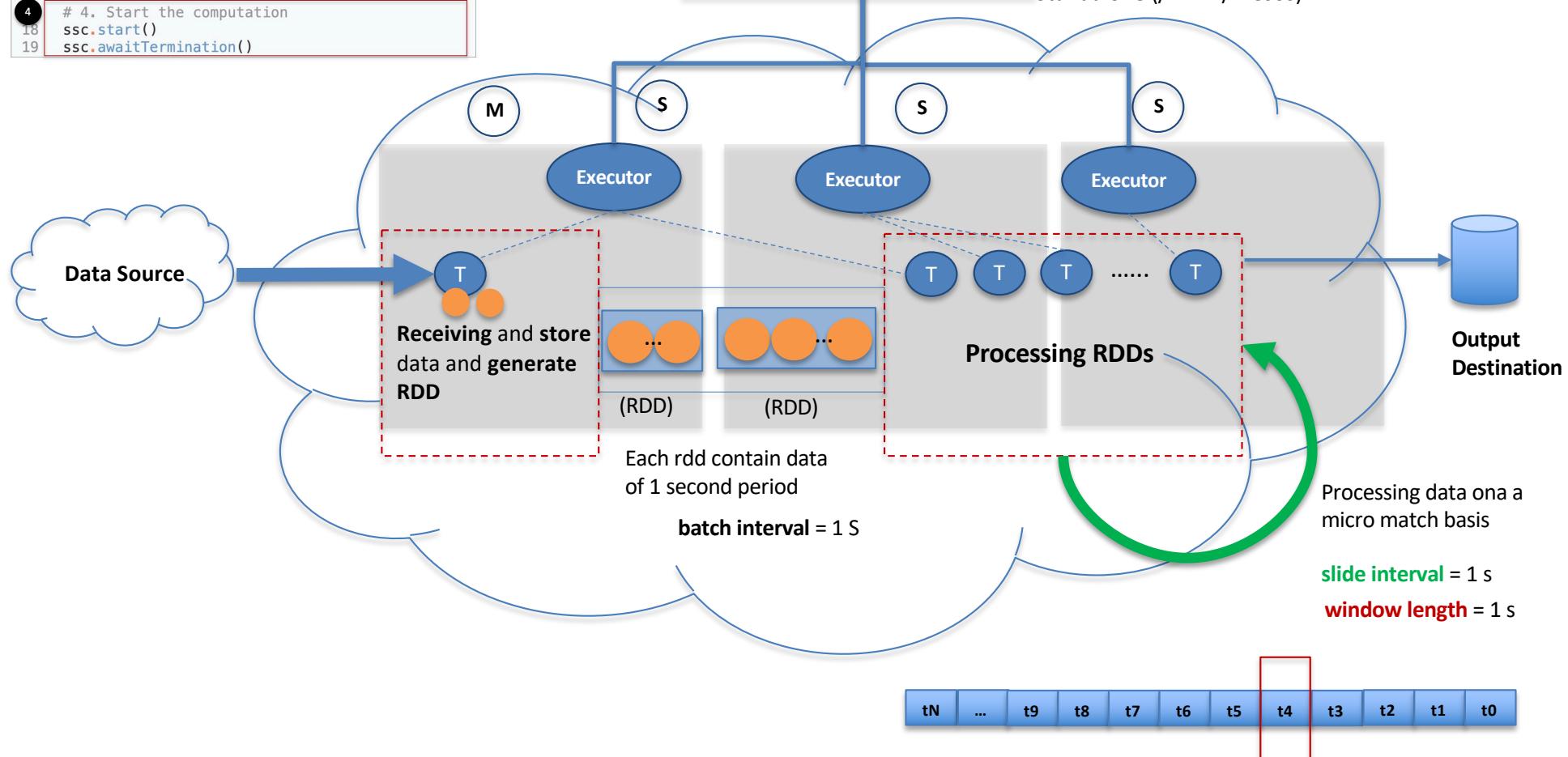
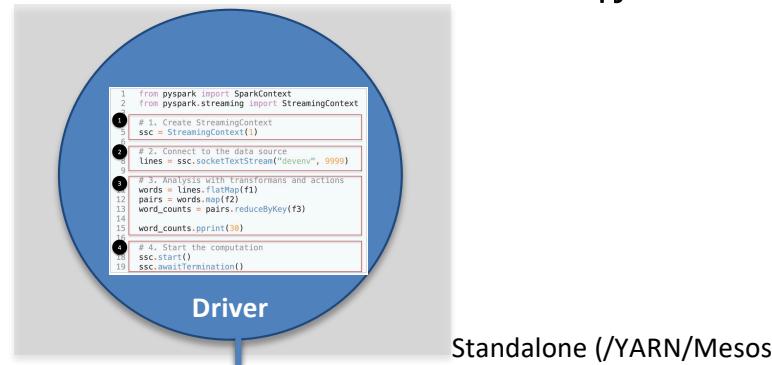
😊 Spark Streaming API for Python: `pyspark.sql`

```

1  from pyspark import SparkContext
2  from pyspark.streaming import StreamingContext
3
4  # 1. Create StreamingContext
5  ssc = StreamingContext(1)
6
7  # 2. Connect to the data source
8  lines = ssc.socketTextStream("devenv", 9999)
9
10 # 3. Analysis with transforms and actions
11 words = lines.flatMap(f1)
12 pairs = words.map(f2)
13 word_counts = pairs.reduceByKey(f3)
14
15 word_counts.pprint(30)
16
17 # 4. Start the computation
18 ssc.start()
19 ssc.awaitTermination()

```

`spark-submit -master ... network_wordcount.py`



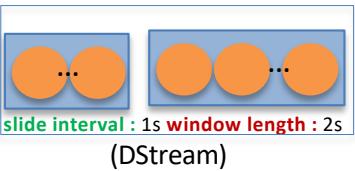
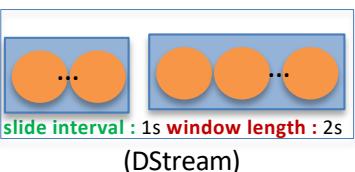
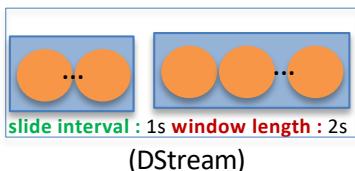
# Spark Streaming API Programming model



Spark Streaming API for Python: `pyspark.streaming`

**DStream** object represents the streaming data

**DStream** type provides transformation and actions like RDD



```
1  from pyspark import SparkContext
2  from pyspark.streaming import StreamingContext
3
4
5  # 1. Create StreamingContext
6  ssc = StreamingContext(1)
7
8  # 2. Connect to the data source
9  lines = ssc.socketTextStream("devenv", 9999).window(2, 1)
10
11 # 3. Analysis with transformans and actions
12 words = lines.flatMap(f1)
13 pairs = words.map(f2)
14 word_counts = pairs.reduceByKey(f3)
15
16 word_counts.pprint(30)
17
18 # 4. Start the computation
19 ssc.start()
20 ssc.awaitTermination()
```

# Spark Streaming Run-time Execution Flow

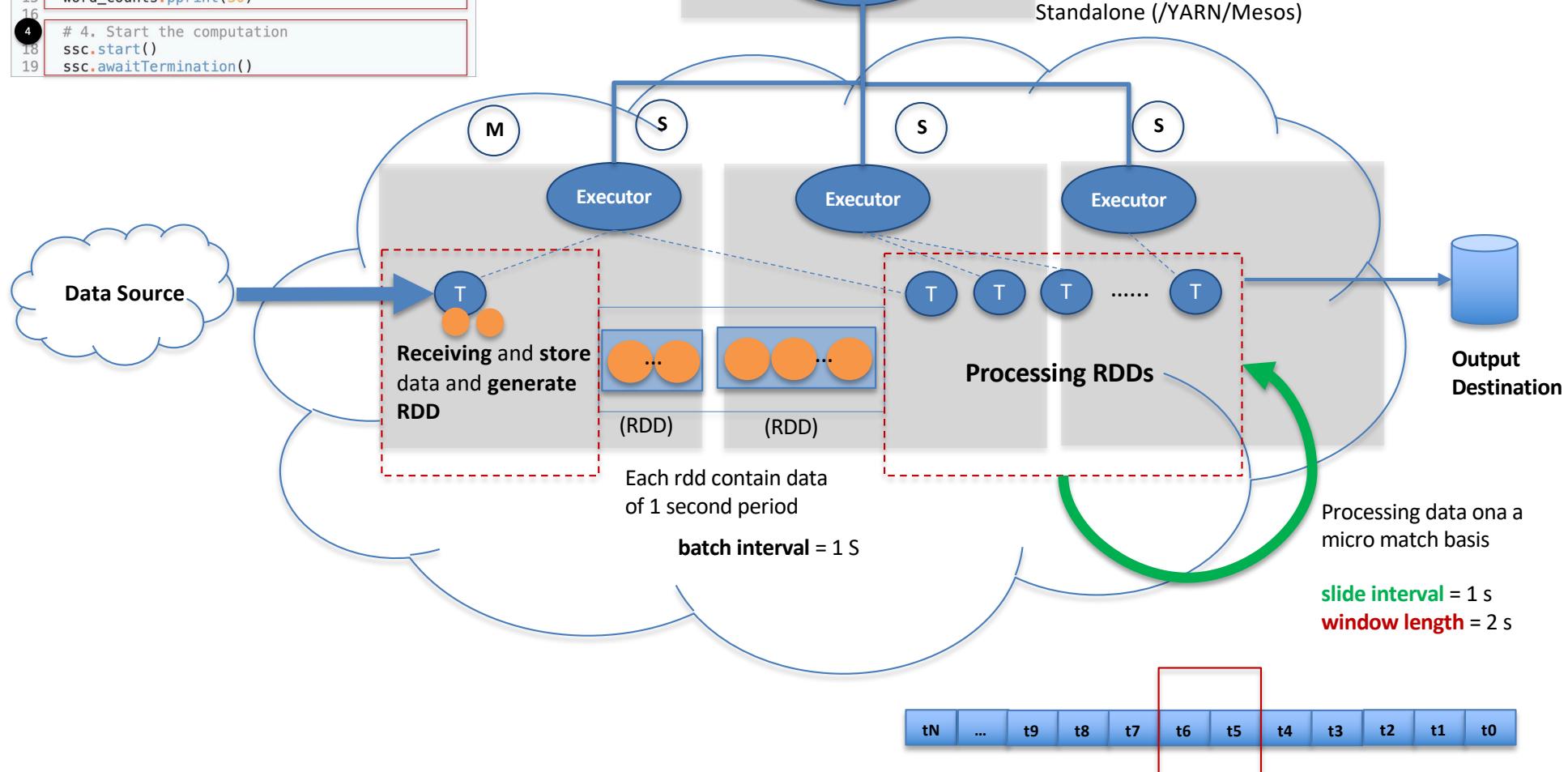
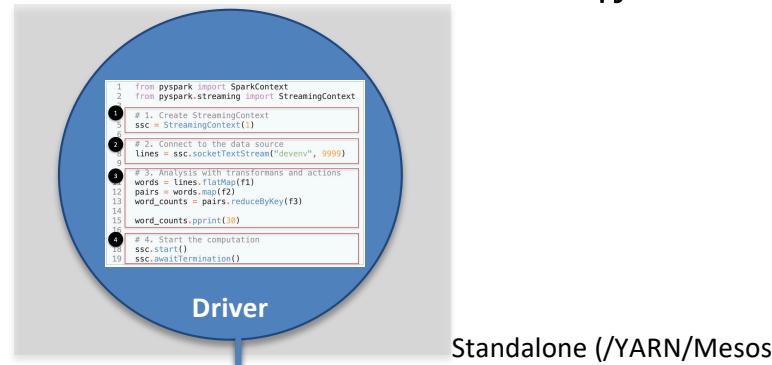
😊 Spark Streaming API for Python: `pyspark.sql`

```

1  from pyspark import SparkContext
2  from pyspark.streaming import StreamingContext
3
4  # 1. Create StreamingContext
5  ssc = StreamingContext(1)
6
7  # 2. Connect to the data source
8  lines = ssc.socketTextStream("devenv", 9999)
9
10 # 3. Analysis with transforms and actions
11 words = lines.flatMap(f1)
12 pairs = words.map(f2)
13 word_counts = pairs.reduceByKey(f3)
14
15 word_counts.pprint(30)
16
17 # 4. Start the computation
18 ssc.start()
19 ssc.awaitTermination()

```

`spark-submit -master ... network_wordcount.py`



# Dive-in Demo: Network Wordcount

```
1  from pyspark import SparkContext
2  from pyspark.streaming import StreamingContext
3
4  if __name__ == "__main__":
5      # Create a local StreamingContext with two working thread and batch interval of 1 second
6      sc = SparkContext()
7      ssc = StreamingContext(sc, 1)
8
9      # Create a DStream that will connect to hostname:port, like localhost:9999
10     lines = ssc.socketTextStream("localhost", 9999)
11
12     # Split each line into words
13     words = lines.flatMap(lambda line: line.split(" "))
14
15     # Count each word in each batch
16     pairs = words.map(lambda word: (word, 1))
17     word_counts = pairs.reduceByKey(lambda x, y: x + y)
18
19     # Print the first ten elements of each RDD generated in this DStream to the console
20     word_counts.pprint()
21
22     ssc.start()          # Start the computation
23     ssc.awaitTermination() # Wait for the computation to terminate
```

# Dive-in Demo: Network Wordcount

- Launch NC server

```
[cloudera@quickstart Desktop]$ nc -lk 9999  
a b c  
a b
```

- Launch the spark streaming job

# StreamingContext

- The main entry point of all Spark Streaming functionality.
- Batch interval
  - Batch interval decides the batch computation interval and the size the each RDD in the input DStreams.
  - Should be based on the latency requirements and available cluster resources

```
from pyspark import SparkContext
from pyspark.streaming import StreamingContext

sc = SparkContext(master, appName)
ssc = StreamingContext(sc, 1)
```



# StreamingContext

- StreamingContext is also used when input DStreams are created.
  - Basic sources
    - Sources directly available in the StreamingContext API. e.g. file systems, socket connections, ...
  - Advanced sources
    - Sources like Kafka, Flume, Kinesis, Twitter, etc. are available through extra utility classes.

# DStream

- A Dstream (Discretized Stream) represents a continuous stream of data.
- Internally, a DStream is represented by a continuous series of RDDs.



# DStreams and Receivers

- Every input DStream is associated with a Receiver object. (except file stream)
  - Input DStreams are DStreams representing the stream of input data received from streaming sources.
- The receiver receives the data from a source and stores it in Spark's memory for processing.

# DStreams and Receivers

- Allocate enough cores to both run the receivers and process the received data.
  - i.e. the # of cores allocated to the Spark Streaming application must be more than the number of receivers. Otherwise the system will receive data, but not be able to process it.
- Creating multiple input Dstreams create multiple receivers.

# File Streams

- File Streams
  - for reading data from files on any file system compatible with the HDFS API (HDFS, S3,...)  
(for general hadoop files)

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

`streamingContext.textFileStream(dataDirectory)` (for text)

- It monitors the directory and process any files created in that directory
  - Files written in nested directories not supported
  - The files must have the same data format.
  - The files must be created by atomically *moving* or *renaming* them into the data directory. (e.g. hadoop fs -mv)
  - Once moved, the files must not be changed.

# WordCount from a text folder

```
1  from pyspark import SparkContext
2  from pyspark.streaming import StreamingContext
3
4  if __name__ == "__main__":
5      # Create a local StreamingContext with two working thread and batch interval of 1 second
6      sc = SparkContext()
7      ssc = StreamingContext(sc, 1)
8
9      # Create a DStream that will monitor the textfolder
10     lines = ssc.textFileStream("hdfs://localhost/user/cloudera/spark_streaming_101/textfolder")
11
12     # Split each line into words
13     words = lines.flatMap(lambda line: line.split(" "))
14
15     # Count each word in each batch
16     pairs = words.map(lambda word: (word, 1))
17     word_counts = pairs.reduceByKey(lambda x, y: x + y)
18
19     # Print the first ten elements of each RDD generated in this DStream to the console
20     word_counts.pprint()
21
22     ssc.start()          # Start the computation
23     ssc.awaitTermination() # Wait for the computation to terminate
```

# Spark Streaming + Kafka Integration

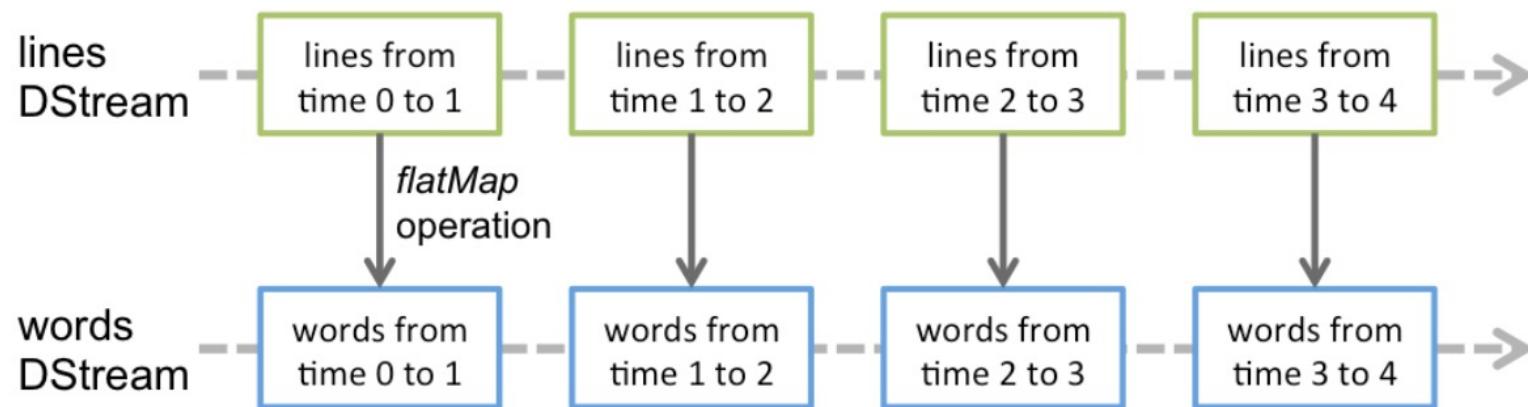
## Receiver-based Approach

```
# Create a DStream from Kafka "page_views" topic and apply window to the stream
page_view_stream = KafkaUtils.createStream(ssc, "localhost:2181", "consumer-group", {"page_views_logs_stream": 1})
lines = page_view_stream.map(lambda (k,v): v).window(60,10)
```

- Use a Receiver to receive the data. The Receiver is implemented using the Kafka high-level consumer API.
- Require external Kafka receiver libraries to use the API.

# Dstream operations

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



These underlying RDD transformations are computed by the Spark engine.

# Common Transformation Methods on DStream

<b>map(func)</b>	Return a new DStream by passing each element of the source DStream through a function <i>func</i> .
<b>flatMap(func)</b>	Similar to map, but each input item can be mapped to 0 or more output items.
<b>filter(func)</b>	Return a new DStream by selecting only the records of the source DStream on which <i>func</i> returns true.
<b>repartition(numPartitions)</b>	Changes the level of parallelism in this DStream by creating more or fewer partitions.
<b>union(otherStream)</b>	Return a new DStream that contains the union of the elements in the source DStream and <i>otherDStream</i> .
<b>count()</b>	Return a new DStream of single-element RDDs by counting the number of elements in each RDD of the source DStream.
<b>reduce(func)</b>	Return a new DStream of single-element RDDs by aggregating the elements in each RDD of the source DStream using a function <i>func</i> (which takes two arguments and returns one). The function should be associative and commutative so that it can be computed in parallel.
<b>countByValue()</b>	When called on a DStream of elements of type K, return a new DStream of (K, Long) pairs where the value of each key is its frequency in each RDD of the source DStream.
<b>reduceByKey(func, [numTasks])</b>	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. <b>Note:</b> By default, this uses Spark's default number of parallel tasks (2 for local mode, and in cluster mode the number is determined by the config property <code>spark.default.parallelism</code> ) to do the grouping. You can pass an optional numTasks argument to set a different number of tasks.

# Common Transformation Methods on DStream

## **groupByKey(*numPartitions=None*)**

Return a new DStream by applying groupByKey on each RDD.

## **mapValues(*f*)**

Return a new DStream by applying a map function to the value of each key-value pairs in this DStream without changing the key.

## **transform(*func*)**

Return a new DStream by applying a RDD-to-RDD function to every RDD of the source DStream. This can be used to do arbitrary RDD operations on the DStream.

# Common Transformation Methods on DStream

<b>window(windowLength, slideInterval)</b>	Return a new DStream which is computed based on windowed batches of the source DStream.
<b>countByWindow(windowLength, slideInterval)</b>	Return a sliding window count of elements in the stream.
<b>reduceByWindow(func, windowLength, slideInterval)</b>	Return a new single-element stream, created by aggregating elements in the stream over a sliding interval using <i>func</i> . The function should be associative and commutative so that it can be computed correctly in parallel.
<b>reduceByKeyAndWindow(func, windowLength, slideInterval, [numTasks])</b>	When 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 the given reduce function <i>func</i> over batches in a sliding window. <b>Note:</b> By default, this uses Spark's default number of parallel tasks (2 for local mode, and in cluster mode the number is determined by the config property <code>spark.default.parallelism</code> ) to do the grouping. You can pass an optional <i>numTasks</i> argument to set a different number of tasks.
<b>reduceByKeyAndWindow(func, invFunc, windowLength, slideInterval, [numTasks])</b>	A more efficient version of the above <code>reduceByKeyAndWindow()</code> where the reduce value of each window is calculated incrementally using the reduce values of the previous window. This is done by reducing the new data that enters the sliding window, and “inverse reducing” the old data that leaves the window. An example would be that of “adding” and “subtracting” counts of keys as the window slides. However, it is applicable only to “invertible reduce functions”, that is, those reduce functions which have a corresponding “inverse reduce” function (taken as parameter <i>invFunc</i> ). Like in <code>reduceByKeyAndWindow</code> , the number of reduce tasks is configurable through an optional argument. Note that <a href="#">checkpointing</a> must be enabled for using this operation.
<b>countByValueAndWindow(windowLength, slideInterval, [numTasks])</b>	When called on a DStream of (K, V) pairs, returns a new DStream of (K, Long) pairs where the value of each key is its frequency within a sliding window. Like in <code>reduceByKeyAndWindow</code> , the number of reduce tasks is configurable through an optional argument.

# Common Transformation Methods on DStream

**groupByKeyAndWindow**(*windowDuration*, *slideDuration*, *numPartitions=None*)

Return a new DStream by applying *groupByKey* over a sliding window. Similar to *DStream.groupByKey()*, but applies it over a sliding window.

- Parameters:**
- **windowDuration** – width of the window; must be a multiple of this DStream's batching interval
  - **slideDuration** – sliding interval of the window (i.e., the interval after which the new DStream will generate RDDs); must be a multiple of this DStream's batching interval
  - **numPartitions** – Number of partitions of each RDD in the new DStream.

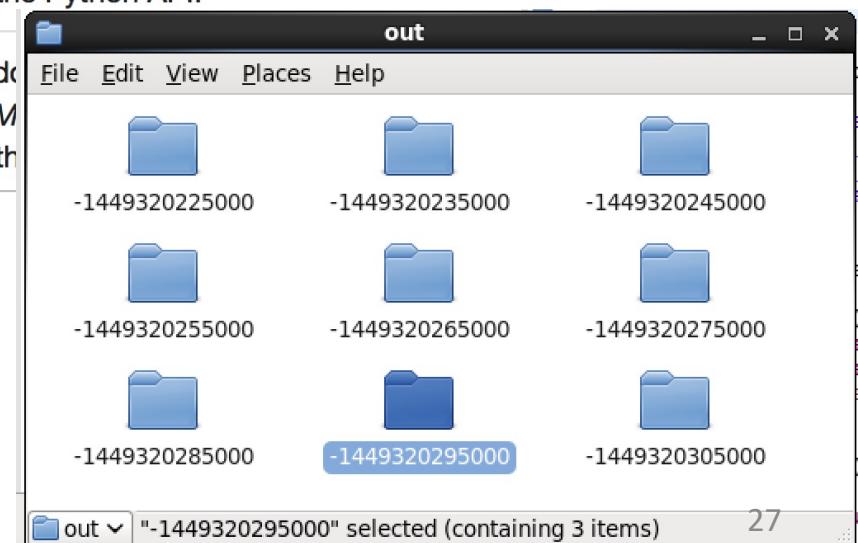
# Window Network WordCount

```
1  from pyspark import SparkContext
2  from pyspark.streaming import StreamingContext
3
4  if __name__ == "__main__":
5      # Create a local StreamingContext with two working thread and batch interval of 1 second
6      sc = SparkContext()
7      ssc = StreamingContext(sc, 1)
8
9      # Create a DStream that will connect to hostname:port, like localhost:9999
10     lines = ssc.socketTextStream("localhost", 9999)
11     lines_windowed = lines.window(20, 10) # process last 20 seconds of data, every 10 seconds
12
13     # Split each line into words
14     words = lines_windowed.flatMap(lambda line: line.split(" "))
15
16     # Count each word in each batch
17     pairs = words.map(lambda word: (word, 1))
18     word_counts = pairs.reduceByKey(lambda x, y: x + y)
19
20     # Print the first ten elements of each RDD generated in this DStream to the console
21     word_counts.pprint()
22
23     ssc.start()          # Start the computation
24     ssc.awaitTermination() # Wait for the computation to terminate
```

# Output Operations on DStreams

Output operations trigger the actual execution of all the DStream transformations (similar to actions for RDDs)

<b>print()</b>	Prints the first ten elements of every batch of data in a DStream on the driver node running the streaming application. This is useful for development and debugging. <b>Python API</b> This is called <code>pprint()</code> in the Python API.
<b>saveAsTextFiles(prefix, [suffix])</b>	Save this DStream's contents as text files. The file name at each batch interval is generated based on <i>prefix</i> and <i>suffix</i> : " <i>prefix-TIME_IN_MS[.suffix]</i> ".
<b>saveAsObjectFiles(prefix, [suffix])</b>	Save this DStream's contents as SequenceFiles of serialized Java objects. The file name at each batch interval is generated based on <i>prefix</i> and <i>suffix</i> : " <i>prefix-TIME_IN_MS[.suffix]</i> ". <b>Python API</b> This is not available in the Python API.
<b>saveAsHadoopFiles(prefix, [suffix])</b>	Save this DStream's contents as Hadoop files based on <i>prefix</i> and <i>suffix</i> : " <i>prefix-TIM... <b>Python API</b> This is not available in th</i>



# Output Operations on Dstreams

`foreachRDD(func)` allows data to be sent out to external systems

`foreachRDD(func)`

The most generic output operator that applies a function, *func*, to each RDD generated from the stream. This function should push the data in each RDD to an external system, such as saving the RDD to files, or writing it over the network to a database. Note that the function *func* is executed in the driver process running the streaming application, and will usually have RDD actions in it that will force the computation of the streaming RDDs.

# Design Patterns for using foreachRDD

This does not work!

```
def sendRecord(rdd):
    connection = createNewConnection() # executed at the driver
    rdd.foreach(lambda record: connection.send(record))
    connection.close()

dstream.foreachRDD(sendRecord)
```

Poor solution

```
def sendRecord(record):
    connection = createNewConnection()
    connection.send(record)
    connection.close()

dstream.foreachRDD(lambda rdd: rdd.foreach(sendRecord))
```

# Design Patterns for using foreachRDD

Better solution

```
def sendPartition(iter):
    connection = createNewConnection()
    for record in iter:
        connection.send(record)
    connection.close()

dstream.foreachRDD(lambda rdd: rdd.foreachPartition(sendPartition))
```

# Persistence

- Use `persist()` method on a DStream will automatically persist every RDD of that DStream in memory.
  - This is useful if the data in the DStream will be computed multiple times (e.g., multiple operations on the same data)
- DStreams generated by window-based and `updateStateByKey` operations are automatically persisted in memory, without calling `persist()`

Storage Level	Meaning
MEMORY_ONLY	Store RDD as deserialized Java objects in the JVM. If the RDD does not fit in memory, some partitions will not be cached and will be recomputed on the fly each time they're needed. This is the default level.
MEMORY_AND_DISK	Store RDD as deserialized Java objects in the JVM. If the RDD does not fit in memory, store the partitions that don't fit on disk, and read them from there when they're needed.
MEMORY_ONLY_SER (Java and Scala)	Store RDD as <i>serialized</i> Java objects (one byte array per partition). This is generally more space-efficient than deserialized objects, especially when using a <a href="#">fast serializer</a> , but more CPU-intensive to read.
MEMORY_AND_DISK_SER (Java and Scala)	Similar to MEMORY_ONLY_SER, but spill partitions that don't fit in memory to disk instead of recomputing them on the fly each time they're needed.
DISK_ONLY	Store the RDD partitions only on disk.
MEMORY_ONLY_2, MEMORY_AND_DISK_2, etc.	Same as the levels above, but replicate each partition on two cluster nodes.
OFF_HEAP (experimental)	Similar to MEMORY_ONLY_SER, but store the data in <a href="#">off-heap memory</a> . This requires off-heap memory to be enabled.

# Custom Receiver

- If none suits you, you can create your own ones.
- A custom receiver must extend this abstract class *Receiver* by implementing two methods

```
class CustomReceiver(host: String, port: Int)
  extends Receiver[String](StorageLevel.MEMORY_AND_DISK_2) with Logging {

  def onStart() {
    // Start the thread that receives data over a connection
    new Thread("Socket Receiver") {
      override def run() { receive() }
    }.start()
  }

  def onStop() {
    // There is nothing much to do as the thread calling receive()
    // is designed to stop by itself if isStopped() returns false
  }
}
```

# Custom Receiver (Reference)

- You can create your own ones if none suits your custom receiver.
- A custom receiver must extend this abstract class *Receiver* by implementing two methods.
  - `onStart()`:
    - Things to do to start receiving data like creating receiving threads.
    - `onStart()` must not block indefinitely.
    - Inside the implementation, once the data is received, call `Receiver's store(data)` method to save the data to Spark.
  - `onStop()`:
    - Things to do to stop receiving data.

# A custom receiver example

A custom receiver that receives a stream of text over a socket. It treats \n delimited lines in the text stream as records and stores them with Spark

```
class CustomReceiver(host: String, port: Int)
  extends Receiver[String](StorageLevel.MEMORY_AND_DISK_2) with Logging {

  def onStart() {
    // Start the thread that receives data over a connection
    new Thread("Socket Receiver") {
      override def run() { receive() }
    }.start()
  }

  def onStop() {
    // There is nothing much to do as the thread calling receive()
    // is designed to stop by itself if isStopped() returns false
  }
}
```

# A custom receiver example

```
/** Create a socket connection and receive data until receiver is stopped */
private def receive() {
    var socket: Socket = null
    var userInput: String = null
    try {
        // Connect to host:port
        socket = new Socket(host, port)

        // Until stopped or connection broken continue reading
        val reader = new BufferedReader(
            new InputStreamReader(socket.getInputStream(), StandardCharsets.UTF_8))
        userInput = reader.readLine()
        while(!isStopped && userInput != null) {
            store(userInput)
            userInput = reader.readLine()
        }
        reader.close()
        socket.close()

        // Restart in an attempt to connect again when server is active again
        restart("Trying to connect again")
    } catch {
        case e: java.net.ConnectException =>
            // restart if could not connect to server
            restart("Error connecting to " + host + ":" + port, e)
        case t: Throwable =>
            // restart if there is any other error
            restart("Error receiving data", t)
    }
}
```

Data can be stored inside Spark by calling store(data)

# A custom receiver example

- Using the custom receiver in a Spark Streaming application

```
// Assuming ssc is the StreamingContext
val customReceiverStream = ssc.receiverStream(new CustomReceiver(host, port))
val words = lines.flatMap(_.split(" "))
...
```