

Java-Success.com

Prepare to fast-track, choose & go places with 800+ Java & Big Data Q&As with lots of code & diagrams.

[Home](#) [Why? ▾](#) [300+ Java FAQs ▾](#) [300+ Big Data FAQs ▾](#) [Courses ▾](#)

[👤 Membership ▾](#) [Your Career ▾](#)

[Home](#) > [bigdata-success.com](#) > [Tutorials - Big Data](#) > [TUT - Cloudera on Docker](#) > 31:

Docker Tutorial: Apache Spark streaming in Scala with Apache Kafka on Cloudera quickstart

31: Docker Tutorial: Apache Spark streaming in Scala with Apache Kafka on Cloudera quickstart

 Posted on [July 8, 2019](#)

This extends [27: Docker Tutorial: Apache Kafka with Java API on Cloudera quickstart](#).

Pre-requisite: Java 8 & Kafka installed with the “MyTestTopic” as per the previous tutorials.

Install scala & sbt

300+ Java Interview FAQs

300+ Java FAQs



16+ Java Key Areas Q&As



150+ Java Architect FAQs



80+ Java Code Quality Q&As



150+ Java Coding Q&As



300+ Big Data Interview FAQs

300+ Big Data FAQs



Tutorials - Big Data



TUT -  Starting Big Data

TUT - Starting Spark & Scala

We will install Scala 2.10.7 from “<https://www.scala-lang.org/download/2.10.7.html>”, which uses Java 8.

Step 1: Install Scala via curl.

```
1 [root@quickstart /]# curl -O -L https://downloads.
2
```

Copy to “/opt” folder and untar.

```
1 [root@quickstart /]# cp scala-2.10.7.tgz /opt
2 [root@quickstart /]# cd /opt
3 [root@quickstart opt]# tar xzf scala-2.10.7.tgz
4 [root@quickstart opt]# rm -f scala-2.10.7.tgz
5 [root@quickstart opt]# ls -ltr
6 total 16
7 drwxr-xr-x 5 cloudera cloudera 4096 Aug 1
8 drwxr-xr-x 4 cloudera cloudera 4096 Aug 2
9 drwxr-xr-x 4 cloudera-scm cloudera-scm 4096 Apr 6
10 drwxrwxr-x 9 1001 1001 4096 Nov 3
11 [root@quickstart opt]#
12
```

Step 2: Install sbt via curl.

Go to <https://www.scala-lang.org/download/2.10.7.html> and click on “Download SBT”.

```
1 [root@quickstart opt]# curl -O -L https://piccolo.
2
```

Untar the tgz file.

```
1 [root@quickstart opt]# tar xzf sbt-1.2.8.tgz
2 [root@quickstart opt]# rm -f sbt-1.2.8.tgz
3 [root@quickstart opt]# ls -ltr
4 total 20
5 drwxr-xr-x 5 cloudera cloudera 4096 Aug 1
```

TUT - Starting with Python

TUT - Kafka

TUT - Pig

TUT - Apache Storm

TUT - Spark Scala on Zeppelin

TUT - Cloudera

TUT - Cloudera on Docker

TUT - File Formats

TUT - Spark on Docker

TUT - Flume

TUT - Hadoop (HDFS)

TUT - HBase (NoSQL)

TUT - Hive (SQL)

TUT - Hadoop & Spark

TUT - MapReduce

TUT - Spark and Scala

TUT - Spark & Java

TUT - PySpark on Databricks

TUT - Zookeeper

800+ Java Interview Q&As

300+ Core Java Q&As



300+ Enterprise Java Q&As



150+ Java Frameworks Q&As



120+ Companion Tech Q&As



Tutorials - Enterprise Java



```
6 drwxr-xr-x 4 cloudera cloudera 4096 Aug 21
7 drwxr-xr-x 4 cloudera-scm cloudera-scm 4096 Apr 6
8 drwxrwxr-x 9 1001 1001 4096 Nov 3
9 drwxrwxr-x 5 1000 1000 4096 Dec 30
10 [root@quickstart opt]#
11
```

Update ~/.bashrc

Step 3: If you are a root user update the ~/.bashrc file, and if you are a any other user update the ~/.bash_profile so that “scala” and “sbt” command can be run from any folder.

```
1 [root@quickstart opt]# vi ~/.bashrc
2
1 # .bashrc
2
3 # User specific aliases and functions
4
5 alias rm='rm -i'
6 alias cp='cp -i'
7 alias mv='mv -i'
8
9 # Source global definitions
10 if [ -f /etc/bashrc ]; then
11     . /etc/bashrc
12 fi
13
14 export JAVA_HOME=/usr/lib/jvm/jre-1.8.0-openjdk.x86_64
15 export PATH=$JAVA_HOME/bin:$PATH
16
17 SCALA_HOME=/opt/scala-2.10.7
18 SBT_HOME=/opt/sbt
19
20 export PATH=$PATH:$SCALA_HOME/bin:$SBT_HOME/bin
21
```

Activate:

```
1 [root@quickstart opt]# source ~/.bashrc
2
```

Check scala command prompt:

```
1 [root@quickstart opt]# scala
2 Welcome to Scala version 2.10.7 (OpenJDK 64-Bit Se
3 Type in expressions to have them evaluated.
4 Type :help for more information.
5
6 scala>
7
```

Check sbt command prompt:

```
1 [root@quickstart opt]# sbt
2 [info] Updated file /opt/project/build.properties:
3 [info] Loading project definition from /opt/project
4 [info] Updating ProjectRef(uri("file:/opt/project/
5 [info] Done updating.
6 [info] Set current project to opt (in build file:/o
7 [info] sbt server started at local:///root/.sbt/1.0
8 sbt:opt>
9
```

Create Scala project structure

Step 4: Unlike maven archetype:generate, sbt does not create the basic project structure. We can create the sbt project structure with a shell script.

```
1 [root@quickstart ~]# cd ~
2 [root@quickstart ~]# pwd
3 /root
4 [root@quickstart ~]# mkdir projects
5 [root@quickstart ~]# cd projects/
6 [root@quickstart projects]# mkdir my-app
7 [root@quickstart projects]# cd my-app
8 [root@quickstart my-app]# vi mkdirs4sbt.sh
9
```

The “mkdirs4sbt.sh”

```
1 #!/bin/sh
```

```
2 mkdir -p src/{main,test}/{java,resources,scala}
3 mkdir lib project target
4
5 # create an initial build.sbt file
6 echo 'name := "my-app"
7 version := "1.0"
8 scalaVersion := "2.10.7"' > build.sbt
9
1 [root@quickstart my-app]# chmod 755 mkdirs4sbt.sh
2 [root@quickstart my-app]# ./mkdirs4sbt.sh
3
1 [root@quickstart my-app]# tree
2 .
3 |— build.sbt
4 |— lib
5 |— mkdirs4sbt.sh
6 |— project
7 |— src
8 |   |— main
9 |   |   |— java
10 |   |   |— resources
11 |   |   |— scala
12 |   |— test
13 |   |   |— java
14 |   |   |— resources
15 |   |   |— scala
16 |— target
17
18 12 directories, 2 files
19 [root@quickstart my-app]#
20
```

Add spark dependency in build.sbt

Step 5: To write Spark code spark-core api library is required.

```
1 [root@quickstart my-app]# vi build.sbt
2
```

```
1 name := "my-app"
2 version := "1.0"
3 scalaVersion := "2.10.7"
4
5 libraryDependencies += "org.apache.spark" %% "spark-core"
6 libraryDependencies += "org.apache.spark" %% "spark-streaming"
```

```
7 libraryDependencies += "org.apache.spark" %% "spark-streaming-kafka"
8
```

Create the Spark job in Scala

Step 6: Create the package “com.mycompany.app”.

```
1 [root@quickstart my-app]# mkdir -p src/main/scala/com/mycompany/app
```

Step 7: Create “SimpleSparkStreaming.scala”.

```
1 [root@quickstart my-app]# vi src/main/scala/com/mycompany/app/SimpleSparkStreaming.scala
```

```
1 package com.mycompany.app
2
3 import kafka.serializer.StringDecoder
4 import kafka.serializer.DefaultDecoder
5 import org.apache.spark._
6 import org.apache.spark.streaming._
7 import org.apache.spark.streaming.kafka.KafkaUtils
8
9 object SimpleSparkStreaming {
10   def main(args: Array[String]) {
11     val conf = new SparkConf().setAppName("SimpleSparkStreaming")
12     val ssc = new StreamingContext(conf, Seconds(1))
13
14     // zookeeper 2181, kafka 9092
15     val kafkaParams = Map[String, String]("metadata.max.request.size", "1048576")
16     val kafkaTopics = Set("MyTestTopic")
17
18     val directKafkaStream = KafkaUtils.createDirectStream(
19       ssc, kafkaParams, kafkaTopics
20     )
21
22
23     directKafkaStream.foreachRDD { rdd =>
24       rdd.foreach { content =>
25         // code to handle the string here
26         println(content)
27       }
28     }
29
30
31     ssc.start()
32     ssc.awaitTermination()
33
34   }
```

```
35 }  
36  
37
```

Compile & Package with sbt

Step 8: Package it with “sbt”

```
1 [root@quickstart my-app]# sbt package  
2 ....  
3
```

```
1 [root@quickstart my-app]# ls -ltr target/scala-2.10  
2 total 12  
3 drwxr-xr-x 5 root root 4096 Jun  8 09:06 resolution  
4 drwxr-xr-x 3 root root 4096 Jun  8 09:06 classes  
5 -rw-r--r-- 1 root root 2334 Jun  8 09:06 my-app_2.10-1.0.jar  
6 [root@quickstart my-app]#  
7
```

spark-submit to run the spark job

Step 9: Run the Spark job in the jar file via Spark-submit command.

Local client mode

```
1 [root@quickstart my-app]# spark-submit \  
2 --class com.mycompany.app.SimpleSparkStreaming \  
3 --master local \  
4 --deploy-mode client \  
5 target/scala-2.10/my-app_2.10-1.0.jar  
6
```

Publish messages to Kafka topic MyTestTopic

Open a new terminal, and login to the running container:

```
1 $ docker ps
2
3 CONTAINER ID          IMAGE                COMMAND
4 1c781f94a1c2         gdancik/cloudera    "/usr/bin/c
5
1 $ docker exec -it 1c781f94a1c2 /bin/bash
2
3 [root@quickstart /]#
4
```

publish messages to the topic:

```
1 [root@quickstart /]# /home/kafka/kafka/bin/kafka-co
2
1 >Sending a test message
2 >to the Kafka topic
3 >
4
```

Spark streaming console output

```
1 .....
2 (null, Sending a test message)
3 .....
4 (null, to the Kafka topic)
5
```

Where “null” is the key and “Sending a test message” is the value.

If you want to only print the values use “content._2”:

```
1 .....
2   directKafkaStream.foreachRDD { rdd =>
3     rdd.foreach { content =>
4       // code to handle the string here
5       println(content._2)
6     }
7
8   }
9
```


Word count example

Here is the Scala code:

```
1 package com.mycompany.app
2
3 import kafka.serializer.StringDecoder
4 import kafka.serializer.DefaultDecoder
5 import org.apache.spark._
6 import org.apache.spark.streaming._
7 import org.apache.spark.streaming.kafka.KafkaUtils
8
9 object SimpleSparkStreaming {
10   def main(args: Array[String]) {
11     val conf = new SparkConf().setAppName("Simple
12     val ssc = new StreamingContext(conf, Seconds(1))
13
14     // zookeeper 2181, kafka 9092
15     val kafkaParams = Map[String, String]("metadata" -> "localhost:2181")
16     val kafkaTopics = Set("MyTestTopic")
17
18     val directKafkaStream = KafkaUtils.createDirectStream(
19       ssc, kafkaParams, kafkaTopics
20     )
21
22
23     directKafkaStream.map(record => record._2)
24                       .flatMap(line => line.split(" "))
25                       .map(word => (word, 1))
26                       .reduceByKey(_+_ )
27                       .print()
28
29
30
31     ssc.start()
32     ssc.awaitTermination()
33
34   }
35 }
36
```

Package it with sbt:

```
1 [root@quickstart my-app]# sbt package
2
```

Run the Spark job:

```
1 [root@quickstart my-app]# spark-submit \
2 --class com.mycompany.app.SimpleSparkStreaming \
3 --master local \
4 --deploy-mode client \
5 target/scala-2.10/my-app_2.10-1.0.jar
6
```

Publish a message to the topic from a different terminal as illustrated before:

```
1 >a big brown fox jumped over a brown fence
2
```

The output on the Spark streaming console will be:

```
1 -----
2 Time: 1562508765000 ms
3 -----
4 (fox,1)
5 (a,2)
6 (big,1)
7 (fence,1)
8 (over,1)
9 (brown,2)
10 (jumped,1)
11
```

◀ 30: Docker Tutorial: Apache Spark streaming in Python 3 with Apache Kafka on Cloudera quickstart

01: Getting started with Jenkins on Docker tutorial ▶

Disclaimer

The contents in this Java-Success are copyrighted and from EmpoweringTech pty Ltd. The EmpoweringTech pty Ltd has the right to correct or enhance the current content without any prior notice. These are general advice only, and one needs to take his/her own circumstances into consideration. The EmpoweringTech pty Ltd will not be held liable for any damages caused or alleged to be caused either directly or

indirectly by these materials and resources. Any trademarked names or labels used in this blog remain the property of their respective trademark owners. Links to external sites do not imply endorsement of the linked-to sites. [Privacy Policy](#).