# BÁO CÁO LAB5 LƯU TRỮ XỬ LÝ DỮ LIỆU LỚN

Nhóm: Squad Game

Thành viên: Lại Ngọc Thăng Long 20183581

Nguyễn Đình Dũng 20183506

Nguyễn Thành Long 20183586

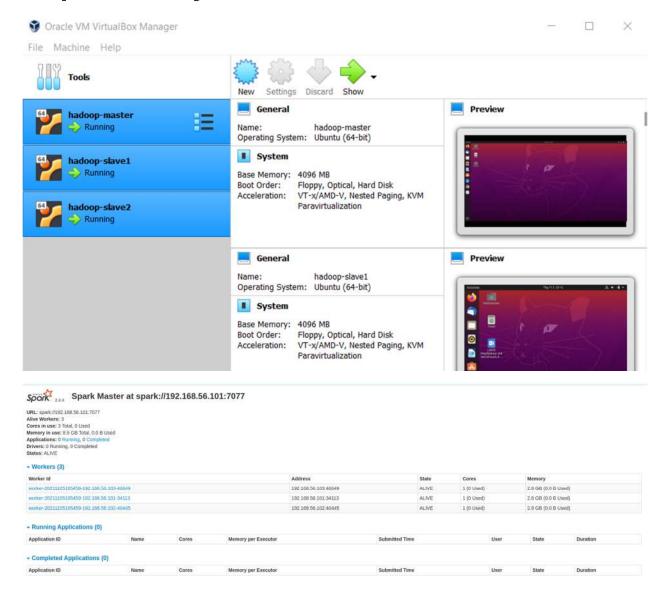
Nguyễn Khương Duy 20183513

**BÀI LÀM:** 

# I. Cài đặt Spark Cluster

1st Step: Cài đặt Spark Cluster trên 3 máy

### Step 1 — 3 máy Master, Slave1, Slave2



Ta thấy Alive Workers: 3

# II, Create a WordCount Application

- Move to the SPARK HOME folder:

\$ cd \$SPARK\_HOME

- Create a directory to save the source code: \$ mkdir -p examples/socket-stream/src/main/scala

- Create a file name SocketStream.scala in \$SPARK\_HOME/examples/socket-stream/src/main/scala folder:

```
hadoopuser@hadoop-master:/usr/local/spark/examples/socket-stream/src/main/scala

SocketStream.scala

bbject SocketStream {
    def main(args: Array[String]) {
        val conf = new SparkConf().setAppName("Socket-Stream")
        // Create a StreamingContext with a 1-second batch size from a SparkConf
        val ssc = new StreamingContext(conf, Seconds(1))
        // Create a DStream using data received after connecting to port 7777 on the

        // local machine
        val lines = ssc.socketTextStream("localhost", 7777)
        // Filter our DStream for lines with "error"
        val errorLines = lines.filter(_.contains("error"))

// Print out the lines with errors
        errorLines.print()

// Start our streaming context and wait for it to "finish"
        ssc.start()

// Walt for the job to finish
        ssc.awaitTermination()

}
```

- Create file build.sbt in \$SPARK HOME/examples/socket-stream:

```
hadoopuser@hadoop-master:/usr/local/spark/examples/socket-stream$ sudo nano buil
d.sbt
[sudo] password for hadoopuser:
hadoopuser@hadoop-master:/usr/local/spark/examples/socket-stream$
```

```
hadoopuser@hadoop-master: /usr/local/spark/examples/w... Q = __

GNU nano 4.8 build.sbt Mc

name := "socket-stream"

version := "0.0.1"

scalaVersion := "2.11.12"

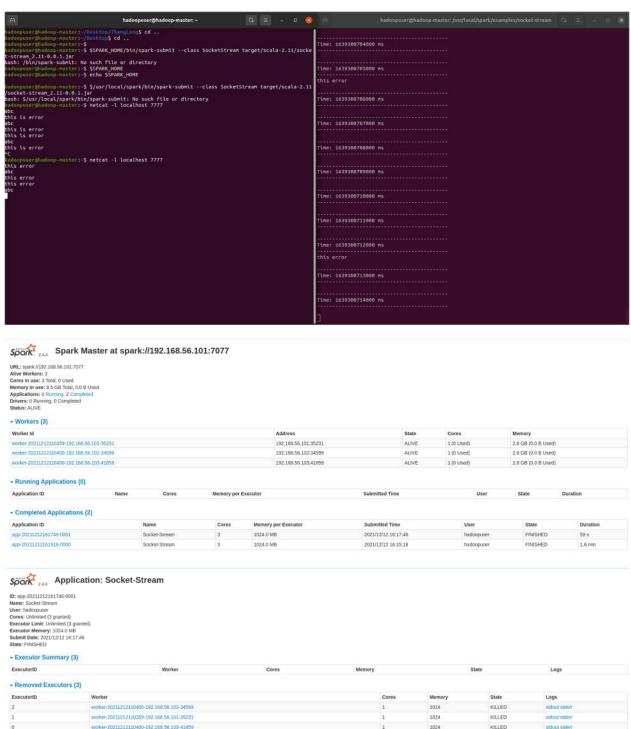
// additional libraries
libraryDependencies ++= Seq(
"org.apache.spark" %% "spark-core" % "2.4.1" % "provided",
"org.apache.spark" %% "spark-streaming" % "2.4.1"
)
```

### - Build application: \$ sbt clean package

```
hadoopuser@hadoop-master:/usr/local/spark/examples/socket-stream$ sbt clean pack
age
[info] Updated file /usr/local/spark/examples/socket-stream/project/build.proper
ties: set sbt.version to 1.5.6
[info] welcome to sbt 1.5.6 (Private Build Java 1.8.0_292)
[info] loading project definition from /usr/local/spark/examples/socket-stream/p
[info] loading settings for project socket-stream from build.sbt ...
[info] set current project to socket-stream (in build file:/usr/local/spark/exam
ples/socket-stream/)
[success] Total time: 0 s, completed Dec 12, 2021 4:13:30 PM
[info] compiling 1 Scala source to /usr/local/spark/examples/socket-stream/targe
t/scala-2.11/classes ...
https://repo1.maven.org/maven2/org/scala-sbt/compiler-bridge_2.11/1.5.7/compile...
  100.0% [########] 2.7 KiB (1.6 KiB / s)
https://repo1.maven.org/maven2/org/scala-sbt/util-interface/1.5.0/util-interfac...
 100.0% [########] 2.5 KiB (6.4 KiB / s)
https://repo1.maven.org/maven2/org/scala-sbt/compiler-bridge_2.11/1.5.7/compile...
 100.0% [########] 52.2 KiB (87.5 KiB / s)
[info] Non-compiled module 'compiler-bridge_2.11' for Scala 2.11.12. Compiling..
[info] Compilation completed in 22.78s.
```

## - Submit and run in Spark:

# \$ \$SPARK\_HOME/bin/spark-submit --class SocketStream --master: spark://192.168.56.101:7077 target/scala-2.11/socket-stream\_2.11-0.0.1.jar



# III, Create a Log Analyzer

- Move to the SPARK\_HOME folder:

```
$ cd $SPARK_HOME
```

- Create a directory to save the source code:

\$ mkdir -p examples/logs-analyzer/src/main/scala

```
hadoopuser@hadoop-master:/usr/local/spark/examples/socket-stream$ cd $SPARK_HOM E
hadoopuser@hadoop-master:/usr/local/spark$ mkdir -p examples/logs-analyzer/src/
main/scala
hadoopuser@hadoop-master:/usr/local/spark$
```

- File log.txt:

```
Open
1 64.242.88.10 - - [07/Mar/2004:16:05:49 -0800] "GET /twiki/bin/edit/Main/Double_bounce_sender?
  topicparent=Main.ConfigurationVariables HTTP/1.1" 401 12846
2 64.242.88.10 - - [07/Mar/2004:16:06:51 -0800] "GET /twiki/bin/rdiff/TWiki/NewUserTemplate?-
  rev1=1.3&rev2=1.2 HTTP/1.1" 200 4523
3 64.242.88.10 - - [07/Mar/2004:16:10:02 -0800] "GET /mailman/listinfo/hsdivision HTTP/1.1"
  200 6291
4 64.242.88.10 - - [07/Mar/2004:16:11:58 -0800] "GET /twiki/bin/view/TWiki/WikiSyntax HTTP/-
  1.1" 200 7352
5 64.242.88.10 - - [07/Mar/2004:16:20:55 -0800] "GET /twiki/bin/view/Main/DCCAndPostFix HTTP/-
  1.1" 200 5253
664.242.88.10 - - [07/Mar/2004:16:23:12 -0800] "GET /twiki/bin/oops/TWiki/AppendixFileSystem?-
  template=oopsmore&param1=1.12&param2=1.12 HTTP/1.1" 200 11382
7 64.242.88.10 - - [07/Mar/2004:16:24:16 -0800] "GET /twiki/bin/view/Main/PeterThoeny HTTP/-
  1.1" 200 4924
8 64.242.88.10 - - [07/Mar/2004:16:29:16 -0800] "GET /twiki/bin/edit/Main/Header_checks?-
  topicparent=Main.ConfigurationVariables HTTP/1.1" 401 12851
9 64.242.88.10 - - [07/Mar/2004:16:30:29 -0800] "GET /twiki/bin/attach/Main/OfficeLocations
  HTTP/1.1" 401 12851
10 64.242.88.10 - - [07/Mar/2004:16:31:48 -0800] "GET /twiki/bin/view/TWiki/-
  WebTopicEditTemplate HTTP/1.1" 200 3732
11 64.242.88.10 - - [07/Mar/2004:16:32:50 -0800] "GET /twiki/bin/view/Main/WebChanges HTTP/1.1"
  200 40520
12 64.242.88.10 - - [07/Mar/2004:16:33:53 -0800] "GET /twiki/bin/edit/Main/-
  Smtpd_etrn_restrictions?topicparent=Main.ConfigurationVariables HTTP/1.1" 401 12851
13 64.242.88.10 - - [07/Mar/2004:16:35:19 -0800] "GET /mailman/listinfo/business HTTP/1.1" 200
14 64.242.88.10 - - [07/Mar/2004:16:36:22 -0800] "GET /twiki/bin/rdiff/Main/WebIndex?-
  rev1=1.2&rev2=1.1 HTTP/1.1" 200 46373
15 64.242.88.10 - - [07/Mar/2004:16:37:27 -0800] "GET /twiki/bin/view/TWiki/DontNotify HTTP/-
  1.1" 200 4140
16 64.242.88.10 - - [07/Mar/2004:16:39:24 -0800] "GET /twiki/bin/view/Main/TokyoOffice HTTP/-
  1.1" 200 3853
17 64.242.88.10 - - [07/Mar/2004:16:43:54 -0800] "GET /twiki/bin/view/Main/MikeMannix HTTP/1.1"
  200 3686
18 64.242.88.10 - - [07/Mar/2004:16:45:56 -0800] "GET /twiki/bin/attach/Main/PostfixCommands
  HTTP/1.1" 401 12846
10 64 2/2 88 18 . . [87/Mar/2884:16:47:12 _8888] "CET /cobote tvt HTTD/1 1" 288 68
```

- Create a file name ApacheAccessLog.scala in \$SPARK\_HOME/examples/logs-analyzer/src/main/scala folder:

```
1 v case class ApacheAccessLog(
     ipAddress: String,
      clientIdentd: String,
      userId: String,
     dateTime: String,
     method: String,
     endpoint: String,
      protocal: String,
      responseCode: Int,
      contentSize: Long){
13 ∨ object ApacheAccessLog {
    val PATTERN = """^(\S+) (\S+) \[([\w:/]+\s[+\-]\d{4})\] "(\S+) (\S+) (\S+)" (\d{3}) (\d+)""".r
        def parseLogLine(log: String): ApacheAccessLog = {
      log match {
        case PATTERN(ipAddress,
        clientIdentd,
       userId,
dateTime,
       method,
       endpoint,
      protocol,
        responseCode,
       contentSize) => new ApacheAccessLog(ipAddress,
             clientIdentd,
              userId,
             dateTime,
             method,
            endpoint,
protocol,
responseCode.toInt,
contentSize.toLong)
          case => throw new RuntimeException(s"""Cannot parse log line $log""")
```

# - Create a file name LogAnalyzerStreaming.scala in \$SPARK\_HOME/examples/logsanalyzer/src/main/scala folder:

```
### Accession of the procession of the procession of the printin("On accession of the printin("Content Sizes Args K, Min. %s, Man. %s".format(contentSizes.map(..., m, ",")")""")

### Accession of the printin("S""Pasponse code counts: %[responseCodeToCount.mkString("[", ",", "]"))"")

### Accession of the printin("S""Pasponse code counts: %[responseCodeToCount.mkString("[", ",", "]"))"")

### Accession of the printin("S"" Padddresses > 10 imag. ("inplication of the printin("S"") Padddresses > 10 imag. ("inplication of the printin("S"
```

- Create file build.sbt in \$SPARK\_HOME/examples/logs-analyzer

```
build.sbt
1    name := "log-analyzer"
2
3    version := "0.0.1"
4
5    scalaVersion := "2.11.12"
6    // additional libraries
7    libraryDependencies ++= Seq(
8    "org.apache.spark" %% "spark-core" % "2.4.1" % "provided",
9    "org.apache.spark" %% "spark-streaming" % "2.4.1"
10    )
11    |
```

- Create the shell script named "stream.sh" that emulates network stream by periodically

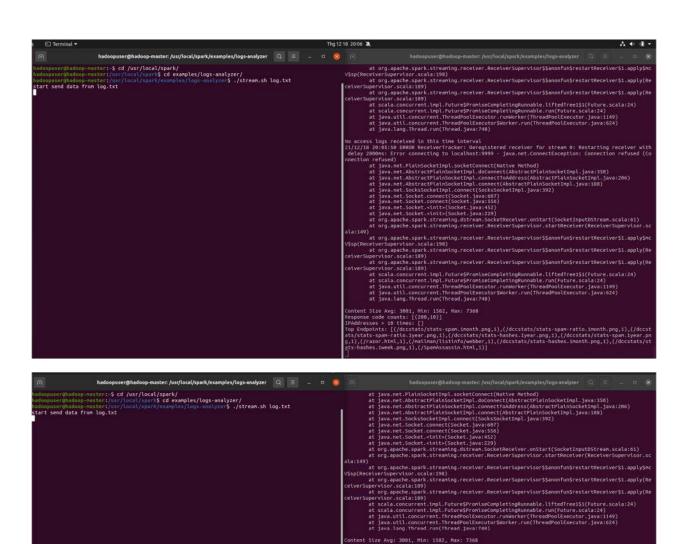
# sending portions of the sample log file to a network socket:

```
$ stream.sh
     set -o nounset
     set -o errexit
     test $# -eq 1 || ( echo "Incorrect number of arguments" ; exit 1 )
     file="$1"
     echo "start send data from $file"
    network port=9999
    lines in batch=100
     interval sec=10
    n lines=$(cat $file | wc -l)
     cursor=10
     while test $cursor -le $n lines
       tail -n $cursor $file | head -n $lines_in_batch | nc -l $network_port
         cursor=$(($cursor + $lines in batch))
         sleep $interval sec
25
```

- Build application: \$ sbt clean package

```
hadoopuser@hadoop-master:/usr/local/spark/examples/logs-analyzer$ sbt clean pack age
[info] welcome to sbt 1.5.6 (Private Build Java 1.8.0_292)
[info] loading project definition from /usr/local/spark/examples/logs-analyzer/p roject
[info] loading settings for project logs-analyzer from build.sbt ...
[info] set current project to log-analyzer (in build file:/usr/local/spark/examples/logs-analyzer/)
[success] Total time: 0 s, completed Dec 18, 2021 7:07:09 PM
[info] compiling 2 Scala sources to /usr/local/spark/examples/logs-analyzer/targ et/scala-2.11/classes ...
[success] Total time: 22 s, completed Dec 18, 2021 7:07:31 PM
hadoopuser@hadoop-master:/usr/local/spark/examples/logs-analyzer$
```

- Submit and run in Spark: \$\$SPARK\_HOME/bin/spark-submit --class "LogAnalyzerStreaming" --master: spark://192.168.56.101:7077 target/scala-2.11/log-analyzer\_2.11-0.0.1.jar



at java-lang Thread-cun(Thread-java-128)

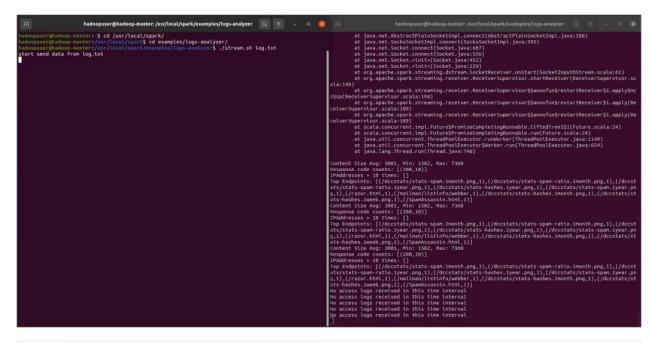
at java-lang Thread-cun(Thread-java-128)

Content Size Avg: 2001, Min: 1582, Max: 7308

Response Yook counts: (200,100)

IPAddresses > 10 times: []

Top Endpoints: [/dccstat/stats-spen.inonth.png.1), (/dccstats/stats-spen-ratio.imonth.png.1), (/dccstats/stats-spen-ratio.imonth.png.1)



#### Spork 244 Spark Master at spark://192.168.56.101:7077

URL: spark://192.168.56.101:7077 Alive Workers: 3 Cores in use: 3 Total, 3 Used Memory in use: 85 GB Total, 3.0 GB Used Applications: 1 Rizzning, 6 Completed Drivers: 0 Rizzning, 0 Completed Status: ALIVE

#### - Workers (3)

Worker Id	Address	State	Cores	Memory
worker-20211212110359-192.168.56.101-39231	192.168.56.101:35231	ALIVE	1 (1 Used)	2.8 GB (1024.0 MB Used)
worker-20211212110400-192-168-56-102-34599	192.168.56.102:34599	ALIVE	1 (1 Used)	2.8 GB (1024.0 MB Used)
worker-20211212110400-192.168.56.103-41859	192 168 56 103 41859	ALIVE	1 (1 Used)	2.8 GB (1024.0 MB Used)

#### - Running Applications (1)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
app-20211212175057-0006 (00)	Log Analyzer Streaming in Scala	3	1024.0 MB	2021/12/12 17:50:57	hadoopuser	RUNNING	4.7 min

#### - Completed Applications (6)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
врр-20211212175022-0005	Log Analyzer Streaming in Scala	3	1024.0 MB	2021/12/12 17:50:22	hadoopuser	FINISHED	19 s
app-20211212174923-0004	Log Analyzer Streaming in Scala	3	1024.0 MB	2021/12/12 17:49:23	hadoopuser	FINISHED	51 s
арр-20211212173540-0003	Log Analyzer Streaming in Scala	3	1024.0 MB	2021/12/12 17:35:40	hadoopuser	FINISHED	11 min
арр-20211212172036-0002	Log Analyzer Streaming in Scala	3	1024.0 MB	2021/12/12 17:20:36	hadoopuser	FINISHED	8.3 min
арр-20211212161746-0001	Socket-Stream	3	1024.0 MB	2021/12/12 16:17:46	hadoopuser	FINISHED	59 s
app-20211212161518-0000	Socket-Stream	3	1024.0 MB	2021/12/12 16:15:18	hadoopuser	FINISHED	1.6 min



ExecutorID

#### Spork 24.4 Application: Log Analyzer Streaming in Scala

D: app-20211212175022-0005
Name: Log Analyzer Streaming in Scala
User: hadoopuser
Cores: Unlimited (3 granted)
Executor Limit: Unlimited (3 granted)
Executor Memory: 1024-0 MB
Sulbmit Date: 2021/12/12 17:50:22
State: FINISHED

#### \* Executor Summary (3)

Lincolne		Soles monthly				Logs	cogs	
- Removed Execut	tors (3)							
ExecutorID	Worker		Cores	Memory	State	Logs		
2	worker-20211212110400-192.168.56.102-34599		1	1024	KILLED	stdout stderr		
1	worker-20211212110359-192 168.56 101-35231		1	1024	KILLED	stdout stderr		
0	worker-20211212110400-192.188.56.103-41859		1	3024	KILLED	stdout stdern		



#### Spark Master at spark://192.168.56.101:7077

URL: speck.#192.108.56.101:7077
Allve Workers: 3
Cores in use: 3 Total, 0 Used
Memory in use: 8.5 GB Total, 0.0 B Used
Applications: 0 Running, 1 Completed
Drivers: 0 Parming, 0 Completed
Status: ALIVE

#### - Workers (3)

Worker Id	Address	State	Cores	Memory
Worker-20211218185741-102-168-56-101-33787	192-168-56-101-33787	ALIVE	1 (D Lisard)	2 S GR (0.0 R Used)
worker-20211218185742-192.108.56.102-36657	192 168 56 102 36657	ALIVE	I (0 Used)	2.8 GB (0.0 B Used)
worker-20211218185742-192.168.56.103-41619	192 168 56 103 41619	ALIVE	1 (0 Used)	2.8 GB (0.0 B Used)

- Running Applications (0)							
Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration

#### - Completed Applications (1)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
app-20211218191810-0000	Log Analyzer Streaming in Scala	3	1024.0 MB	2021/12/18 19:18:10	hadoopuser	FINISHED	49 min



#### Spork 244 Application: Log Analyzer Streaming in Scala

Dic app-20211218191810-0000
Name: Log Analyzer Streaming in Scala
User: hadsoguser
Cores: Unlimited (3 granted)
Executor Limit: Unlimited (3 granted)
Executor Memory: 1022-0 MB
Submit Date: 202112718 19:18:10
State: FINIS

#### - Executor Summary (4)

ExecutorID	Worker	Cores	Memory	State	Logs
	worker-20211218185741-192.168.58.101-33787	1	1024	EXITED	stdout stderr
of remindrational community					

#### - Removed Executors (3)

ExecutorID	Worker	Cores	Memory	State	Logs
1	worker-20211218185742-192.168.56.103-41619	1	1024	KILLED	stdout stderr
3	worker-20211218185741-192.168:56.101-33787	1	1024	KILLED	stdout stderr
0	worker-20211218185742-192-168.56-102-36657	1	1024	KILLED	stdout stderr