# 1. Для части по SQOOP

**Провести импорт таблицы из вашего сервера БД в Hadoop с использованием SQOOP в любых двух вариантах из перечисленных ниже.**
**a. в Hive-таблицу (--hive-import)**
**b. в HDFS в формате avro (--as-avrodatafile)**
**c. в HDFS в формате sequencefile (--as-sequencefile)**
**Если у вас нет своего сервера то можно использовать тот Postgres, который я показал на лекции. Пароль expoter_pass**

Посмотрим при помощи SQOOP содержимое в PosgreSQL.



**Sqoop Help** не работает, через ssh перехожу на node3.novalocal

```
        at org.apache.sqoop.Sqoop.runTool(Sqoop.java:224)
        at org.apache.sqoop.Sqoop.runTool(Sqoop.java:243)
        at org.apache.sqoop.Sqoop.main(Sqoop.java:252)
Exception in thread "main" java.lang.RuntimeException: javax.xml.parsers.ParserConfigurationException: Feature 'http://apache.org/xml/features/xinclude' is not recogniz
ed.
        at org.apache.hadoop.conf.Configuration.loadResource(Configuration.java:2820)
        at org.apache.hadoop.conf.Configuration.loadResources(Configuration.java:2653)
        at org.apache.hadoop.conf.Configuration.getProps(Configuration.java:2559)
        at org.apache.hadoop.conf.Configuration.get(Configuration.java:1078)
        at org.apache.sqoop.tool.SqoopTool.loadPluginsFromConfDir(SqoopTool.java:170)
        at org.apache.sqoop.tool.SqoopTool.loadPlugins(SqoopTool.java:140)
        at org.apache.sqoop.Sqoop.runTool(Sqoop.java:224)
        at org.apache.sqoop.Sqoop.runTool(Sqoop.java:243)
        at org.apache.sqoop.Sqoop.main(Sqoop.java:252)
Caused by: javax.xml.parsers.ParserConfigurationException: Feature 'http://apache.org/xml/features/xinclude' is not recognized.
        at org.apache.xerces.jaxp.DocumentBuilderFactoryImpl.newDocumentBuilder(Unknown Source)
        at org.apache.hadoop.conf.Configuration.loadResource(Configuration.java:2694)
        ... 8 more
[student4_10@manager ~]$ ssh node3.novalocal
Last login: Sat Jun 20 18:27:03 2020 from manager.novalocal
[student4_10@node3 ~]$ sqoop help
Warning: /opt/cloudera/parcels/CDH-5.16.2-1.cdh5.16.2.p0.8/bin/../lib/sqoop/../accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
20/06/21 17:26:38 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-cdh5.16.2
usage: sqoop COMMAND [ARGS]

Available commands:
  codegen            Generate code to interact with database records
  create-hive-table  Import a table definition into Hive
  eval               Evaluate a SQL statement and display the results
  export             Export an HDFS directory to a database table
  help               List available commands
  import             Import a table from a database to HDFS
  import-all-tables  Import tables from a database to HDFS
  import-mainframe   Import datasets from a mainframe server to HDFS
  job                Work with saved jobs
  list-databases     List available databases on a server
  list-tables        List available tables in a database
  merge              Merge results of incremental imports
  metastore          Run a standalone Sqoop metastore
  version            Display version information

See 'sqoop help COMMAND' for information on a specific command.
[student4_10@node3 ~]$
```

работает !

Проверим таблицы в базе pg_db



```
[student4_10@node3 ~]$ sqoop list-databases --connect jdbc:postgresql://node3.novalocal/pg_db --username exporter --password exporter_pass
Warning: /opt/cloudera/parcels/CDH-5.16.2-1.cdh5.16.2.p0.8/bin/../lib/sqoop/../accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
20/06/20 19:57:50 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-cdh5.16.2
20/06/20 19:57:50 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
20/06/20 19:57:50 INFO manager.SqlManager: Using default fetchSize of 1000
template1
template0
postgres
pg_db
[student4_10@node3 ~]$ sqoop list-tables --connect jdbc:postgresql://node3.novalocal/pg_db --username exporter --password exporter_pass
Warning: /opt/cloudera/parcels/CDH-5.16.2-1.cdh5.16.2.p0.8/bin/../lib/sqoop/../accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
20/06/20 19:58:06 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-cdh5.16.2
20/06/20 19:58:06 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
20/06/20 19:58:06 INFO manager.SqlManager: Using default fetchSize of 1000
character
character_work
paragraph
sales_large
wordform
work
chapter
[student4_10@node3 ~]$
```

sqoop list-tables --connect jdbc:postgresql://node3.novalocal/pg_db --username exporter --password exporter_pass

Скопируем Таблицу Work в локальную папку



```
[student4_10@node3 ~]$ sqoop import --connect jdbc:postgresql://node3.novalocal/pg_db --username exporter --password exporter_pass --table work --target-dir /user/stude
nt4_10/hw_5/work --as-avrodatafile
Warning: /opt/cloudera/parcels/CDH-5.16.2-1.cdh5.16.2.p0.8/bin/../lib/sqoop/../accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
20/06/21 18:00:45 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-cdh5.16.2
20/06/21 18:00:45 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
20/06/21 18:00:45 INFO manager.SqlManager: Using default fetchSize of 1000
20/06/21 18:00:45 INFO tool.CodeGenTool: Beginning code generation
20/06/21 18:00:45 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM "work" AS t LIMIT 1
20/06/21 18:00:45 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce
Note: /tmp/sqoop-student4_10/compile/618b15e8a75495be294c889b3e6f1766/work.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
20/06/21 18:00:47 INFO orm.CompilationManager: Writing jar file: /tmp/sqoop-student4_10/compile/618b15e8a75495be294c889b3e6f1766/work.jar
20/06/21 18:00:47 WARN manager.PostgresqlManager: It looks like you are importing from postgresql.
20/06/21 18:00:47 WARN manager.PostgresqlManager: This transfer can be faster! Use the --direct
20/06/21 18:00:47 WARN manager.PostgresqlManager: option to exercise a postgresql-specific fast path.
20/06/21 18:00:47 INFO mapreduce.ImportJobBase: Beginning import of work
20/06/21 18:00:48 INFO Configuration.deprecation: mapred.jar is deprecated. Instead, use mapreduce.job.jar
20/06/21 18:00:49 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM "work" AS t LIMIT 1
20/06/21 18:00:49 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM "work" AS t LIMIT 1
20/06/21 18:00:49 INFO mapreduce.DataDrivenImportJob: Writing Avro schema file: /tmp/sqoop-student4_10/compile/618b15e8a75495be294c889b3e6f1766/work.avsc
20/06/21 18:00:49 INFO Configuration.deprecation: mapred.map.tasks is deprecated. Instead, use mapreduce.job.maps
```

Проверим папку /user/student4_10



```
drwx------   - hdfs        supergroup          0 2020-06-01 16:11 /user/hdfs
drwxr-xr-x   - mapred      supergroup          0 2019-11-18 19:33 /user/history
drwxr-xr-x   - hive        hive                0 2019-11-18 19:57 /user/hive
drwxrwxr-x   - hue         hue                 0 2019-12-08 22:25 /user/hue
drwxr-xr-x   - instructor  instructor          0 2020-03-17 19:35 /user/instructor
drwxr-x--x   - spark       spark               0 2020-01-19 20:20 /user/spark
drwxr-xr-x   - studen4_6   studen4_6           0 2020-05-23 20:27 /user/studen4_6
drwxr-xr-x   - student     student             0 2019-12-02 15:02 /user/student
drwxr-xr-x   - student4_1  student4_1          0 2020-06-10 13:44 /user/student4_1
drwxr-xr-x   - student4_10 student4_10         0 2020-06-21 18:01 /user/student4_10
drwxr-xr-x   - student4_11 student4_11         0 2020-06-07 06:05 /user/student4_11
drwxr-xr-x   - student4_12 student4_12         0 2020-06-21 15:49 /user/student4_12
drwxr-xr-x   - student4_13 student4_13         0 2020-06-16 11:35 /user/student4_13
drwxr-xr-x   - student4_14 student4_14         0 2020-06-16 08:47 /user/student4_14
drwxr-xr-x   - student4_15 student4_15         0 2020-06-15 20:09 /user/student4_15
drwxr-xr-x   - student4_16 student4_16         0 2020-06-07 06:08 /user/student4_16
drwxr-xr-x   - student4_17 student4_17         0 2020-06-07 06:10 /user/student4_17
drwxr-xr-x   - student4_18 student4_18         0 2020-06-07 06:10 /user/student4_18
drwxr-xr-x   - student4_19 student4_19         0 2020-06-07 06:10 /user/student4_19
drwxr-xr-x   - student4_2  student4_2          0 2020-06-06 20:47 /user/student4_2
drwxr-xr-x   - student4_20 student4_20         0 2020-06-07 06:11 /user/student4_20
drwxr-xr-x   - student4_3  student4_3          0 2020-06-21 12:22 /user/student4_3
drwxr-xr-x   - student4_4  student4_4          0 2020-06-14 03:36 /user/student4_4
drwxr-xr-x   - student4_5  student4_5          0 2020-06-08 15:26 /user/student4_5
drwxr-xr-x   - student4_6  student4_6          0 2020-06-07 13:34 /user/student4_6
drwxr-xr-x   - student4_7  student4_7          0 2020-05-31 12:55 /user/student4_7
drwxr-xr-x   - student4_8  student4_8          0 2020-05-23 20:27 /user/student4_8
drwxr-xr-x   - student4_9  student4_9          0 2020-06-05 12:49 /user/student4_9
drwxr-xr-x   - sudent4_3   sudent4_3           0 2020-05-19 19:19 /user/sudent4_3
[student4_10@node3 ~]$ hdfs dfs -ls /user/student4_10/
Found 4 items
drwx------   - student4_10 student4_10         0 2020-05-27 11:00 /user/student4_10/.Trash
drwx------   - student4_10 student4_10         0 2020-06-21 18:01 /user/student4_10/.staging
drwxr-xr-x   - student4_10 student4_10         0 2020-06-05 17:02 /user/student4_10/Datasets
drwxr-xr-x   - student4_10 student4_10         0 2020-06-21 18:01 /user/student4_10/hw_5
[student4_10@node3 ~]$ hdfs dfs -ls /user/student4_10/hw_5/
Found 1 items
drwxr-xr-x   - student4_10 student4_10         0 2020-06-21 18:01 /user/student4_10/hw_5/work
[student4_10@node3 ~]$
```

Скопируем схему структуры таблицы с локальной директории через команду COPYFROMLOCAL

Создадим Таблицу с путями для схемы work.avsc и файла таблицы work

Копируем Paragraph Таблицу в формате Parquet, но сначало создадтм таблицу так как в Parquet нельзя импортировать таким же способ как avro таблицу , без готовой структуры таблицы в базе данных, так как как не импортируется файл со схемой таблицы

Поэтому через Sqoop проверим схему команды

```
sqoop import --m 1 --connect jdbc:postgresql://node3.novalocal/pg_db --username exporter --password
exporter_pass --query "SELECT column_name, DATA_TYPE FROM INFORMATION_SCHEMA.Columns WHERE
table_name='paragraph' AND \$CONDITIONS" --target-dir '/user/student4_10/hw_5_1/work/'
```

```
sqoop import --m 1 --connect jdbc:postgresql://node3.novalocal/pg_db --username exporter --password
exporter_pass --query "SELECT column_name, DATA_TYPE FROM INFORMATION_SCHEMA.Columns WHERE
table_name='paragraph' AND \$CONDITIONS" --target-dir '/user/student4_10/hw_5_1/paragraph/'
```

Check what you got imported from pg_database (table paragraph)

## Create Table in Parquet



Импортируем данные в таблицу

sqoop import --connect jdbc:postgresql://node3.novalocal/pg_db --username exporter --password exporter_pass --table paragraph --hive-import --hive-database student4_10 --hive-table paragraph --as-parquetfile

```
phonetictext,text
stemtext,text
paragraphtype,character varying
section,integer
chapter,integer
charcount,integer
wordcount,integer
[student4_10@node3 ~]$ sqoop import --connect jdbc:postgresql://node3.novalocal/pg_db --username exporter --password exporter_pass --t
able paragraph --hive-import --hive-database student4_10 --hive-table paragraph --as-parquetfile
Warning: /opt/cloudera/parcels/CDH-5.16.2-1.cdh5.16.2.p0.8/bin/../lib/sqoop/../accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
20/06/22 00:44:24 INFO sqoop.Sqoop: Running Sqoop version: 1.4.6-cdh5.16.2
20/06/22 00:44:24 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
20/06/22 00:44:24 INFO tool.BaseSqoopTool: Using Hive-specific delimiters for output. You can override
20/06/22 00:44:24 INFO tool.BaseSqoopTool: delimiters with --fields-terminated-by, etc.
20/06/22 00:44:24 INFO manager.SqlManager: Using default fetchSize of 1000
20/06/22 00:44:24 INFO tool.CodeGenTool: Beginning code generation
20/06/22 00:44:24 INFO tool.CodeGenTool: Will generate java class as codegen_paragraph
20/06/22 00:44:24 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM "paragraph" AS t LIMIT 1
20/06/22 00:44:24 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce
Note: /tmp/sqoop-student4_10/compile/c225d08da52ea0139b23e05152e1e9fd/codegen_paragraph.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
20/06/22 00:44:26 INFO orm.CompilationManager: Writing jar file: /tmp/sqoop-student4_10/compile/c225d08da52ea0139b23e05152e1e9fd/codeg
en_paragraph.jar
20/06/22 00:44:26 WARN manager.PostgresqlManager: It looks like you are importing from postgresql.
20/06/22 00:44:26 WARN manager.PostgresqlManager: This transfer can be faster! Use the --direct
20/06/22 00:44:26 WARN manager.PostgresqlManager: option to exercise a postgresql-specific fast path.
20/06/22 00:44:26 INFO mapreduce.ImportJobBase: Beginning import of paragraph
20/06/22 00:44:27 INFO Configuration.deprecation: mapred.jar is deprecated. Instead, use mapreduce.job.jar
20/06/22 00:44:28 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM "paragraph" AS t LIMIT 1
20/06/22 00:44:28 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM "paragraph" AS t LIMIT 1
20/06/22 00:44:29 INFO hive.metastore: Trying to connect to metastore with URI thrift://manager.novalocal:9083
20/06/22 00:44:30 INFO hive.metastore: Opened a connection to metastore, current connections: 1
20/06/22 00:44:30 INFO hive.metastore: Connected to metastore.
20/06/22 00:44:30 WARN mapreduce.DataDrivenImportJob: Target Hive table 'paragraph' exists! Sqoop will append data into the existing H
ive table. Consider using --hive-overwrite, if you do NOT intend to do appending.
20/06/22 00:44:32 INFO Configuration.deprecation: mapred.map.tasks is deprecated. Instead, use mapreduce.job.maps
20/06/22 00:44:32 INFO client.RMProxy: Connecting to ResourceManager at manager.novalocal/89.208.221.132:8032
```

Делаем
select * from student4_10.paragraph limit 10;



| | paragraph.workid | paragraph.paragraphid | paragraph.paragraphnum | paragraph.charid | paragraph.plaintext |
|---|---|---|---|---|---|
| 1 | henry4p1 | 639729 | 1537 | henry5 | Let's see what they be: read them. |
| 2 | henry4p1 | 639730 | 1538 | peto | [Reads] Item, A capon,. . 2s. 2d. [p]Item, Sa |
| 3 | henry4p1 | 639731 | 1543 | henry5 | O monstrous! but one half-penny-worth of |
| 4 | henry4p1 | 639732 | 1553 | xxx | [Exeunt] |
| 5 | henry4p1 | 639733 | 1554 | peto | Good morrow, good my lord. |
| 6 | henry4p1 | 639734 | 1557 | xxx | [Enter HOTSPUR, WORCESTER, MORTIMEF |
| 7 | henry4p1 | 639735 | 1558 | mortimer | These promises are fair, the parties sure, [p |
| 8 | henry4p1 | 639736 | 1560 | hotspur | Lord Mortimer, and cousin Glendower, [p]W |
| 9 | henry4p1 | 639737 | 1564 | glendower | No, here it is. [p]Sit, cousin Percy; sit, good |
| 10 | henry4p1 | 639738 | 1569 | hotspur | And you in hell, as oft as he hears Owen Gl |

Проверим папку paragraph

Точно таким способом можно выполнить работу и с Avro файлом

**1.Создать Flume-агент с именем, соответствующим имени своего пользователя (например Flume4_20)**
**2. Создать любой Flume поток используя Flume сервис соответствующего номера.**
**• Тип источника источник – exec**
**• Тип канала – memory**
**• Тип слива – hdfs**
3. **Убедиться что данные поступают в слив.**
4. **Создать поверх данных в hdfs таблицу через которую можно просмотреть полученные данные.**
5. **[Продвинутый вариант] Сделать то-же самое используя несколько сливов в разные места, например в HDFS и в HIve одновременно**
6. **[Продвинутый** вариант] Повторить стандартный пример с выборкой сообщений из Twitter.

**1.Создать Flume-агент с именем, соответствующим имени своего пользователя (например Flume4_20)**
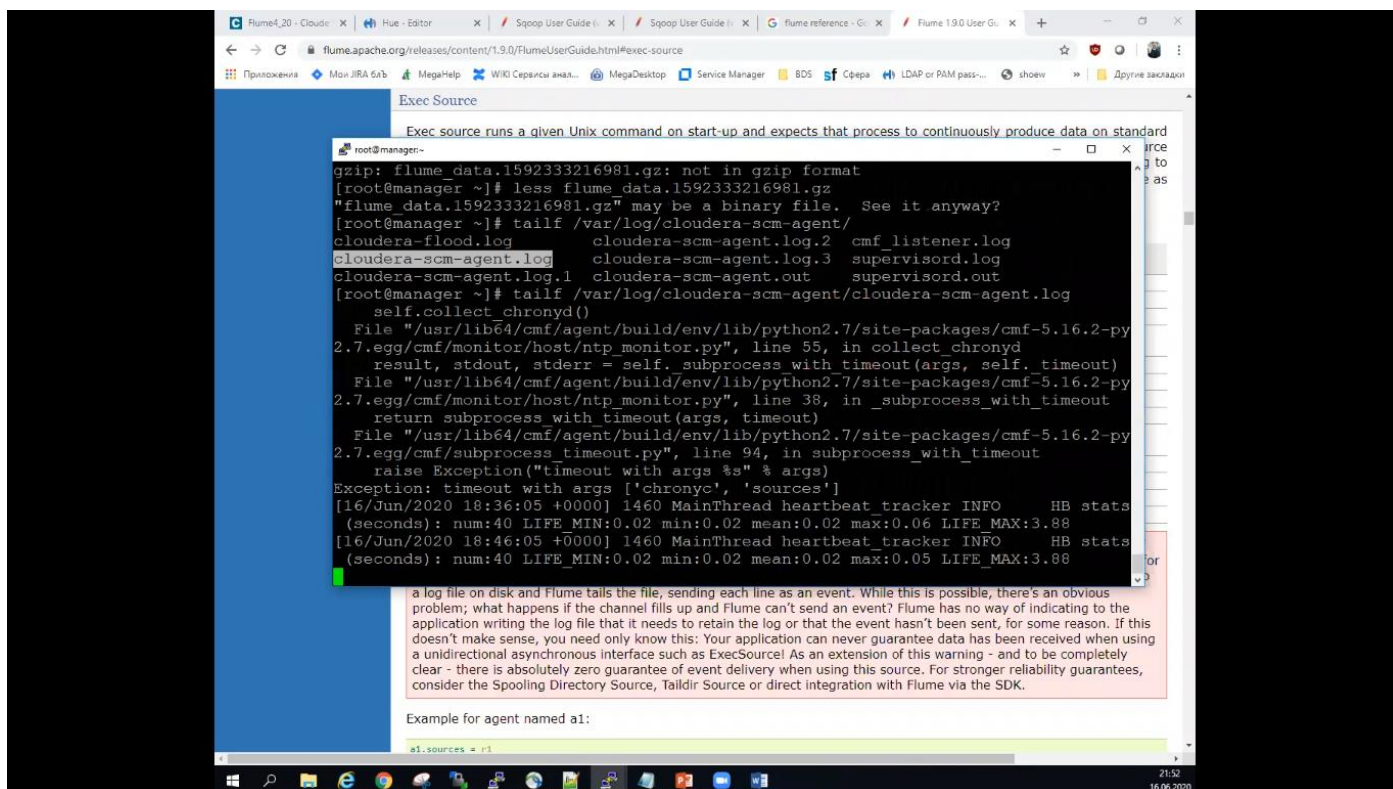
BAD HEALTH!

## 2. Создать любой Flume поток используя Flume сервис соответствующего номера.
• **Тип источника источник – exec**
• **Тип канала – memory**
• **Тип слива – hdfs**

**tailif /var/log/cloudera-scm-agent/cloudera-scm-agent.log**

**работает с этим лог файлом**

**Писать буду в Flume папку !**



**Код для configuration file**

```
# Naming the components on the current agent
Flume4_10.sources = ExecSource
Flume4_10.channels = MemChannel
Flume4_10.sinks = HdfsSink


# Describing/Configuring the source
Flume4_10.sources.ExecSource.type = exec
Flume4_10.sources.ExecSource.command = tailf /var/log/cloudera-scm-agent/cloudera-scm-agent.log
Flume4_10.sources.ExecSource.interceptors = TimestampInterceptor
Flume4_10.sources.ExecSource.interceptors.TimestampInterceptor.type = timestamp


# Describing/Configuring the HDFS sink
Flume4_10.sinks.HdfsSink.type = hdfs
Flume4_10.sinks.LoggerSink.hdfs.path= /flume/flume4_10/%y-%m-%d/
Flume4_10.sinks.HdfsSink.hdfs.filePrefix = events

# Describing/Configuring the channel
Flume4_10.channels.MemChannel.type = memory
Flume4_10.channels.MemChannel.capacity = 10000
Flume4_10.channels.MemChannel.transactionCapacity = 10

# Bind the source and sink to the channel
Flume4_10.sources.ExecSource.channels = MemChannel
Flume4_10.sinks.HdfsSink.channel = MemChannel
```

```
Configuration File          Agent Default Group  ↩

                            # Naming the components on the current agent
                            Flume4_10.sources = ExecSource
                            Flume4_10.channels = MemChannel
                            Flume4_10.sinks = HdfsSink


                            # Describing/Configuring the source
                            Flume4_10.sources.ExecSource.type = exec
                            Flume4_10.sources.ExecSource.command = tailf /var/log/cloudera-scm-agent/cloudera-scm-agent.log
                            Flume4_10.sources.ExecSource.interceptors = TimestampInterceptor
                            Flume4_10.sources.ExecSource.interceptors.TimestampInterceptor.type = timestamp

                            # Describing/Configuring the HDFS sink
                            Flume4_10.sinks.HdfsSink.type = hdfs
                            Flume4_10.sinks.LoggerSink.hdfs.path= /flume/flume4_10/%y-%m-%d/
                            Flume4_10.sinks.HdfsSink.hdfs.filePrefix = events

                            # Describing/Configuring the channel
                            Flume4_10.channels.MemChannel.type = memory
                            Flume4_10.channels.MemChannel.capacity = 10000
                            Flume4_10.channels.MemChannel.transactionCapacity = 10

                            # Bind the source and sink to the channel
                            Flume4_10.sources.ExecSource.channels = MemChannel
                            Flume4_10.sinks.HdfsSink.channel = MemChannel
```

Configuration changes have been saved successfully.

Save Changes



**BAD HEALTH !!!!**

**cloudera** MANAGER   Clusters   Hosts   Diagnostics   Charts   Audits   Administration   Search   Support   student4

○ Flume4_10 (GeekBrains Cluster)

**Start Command**

Status ○ Running   Context Flume4_10   Jun 21, 4:16:55 PM

▼ Completed 0 of 1 step(s).

● Show All Steps   ○ Show Only Failed Steps   ○ Show Running Steps

> ○ Starting 1 roles on service                                Jun 21, 4:16:55 PM   Abort
      0/1 start commands completed.

Abort   Close

**Health History**

> ⊘ 2:49:23 PM          Agent Health Disabled   Show
> ⊘ 2:49:18 PM          Agent Health Unknown    Show
> ⊕ May 3 11:17:50 PM   Agent Health Disabled   Show
> ⊕ May 3 11:17:28 PM   Agent Health Unknown    Show
> ⊘ Apr 29 10:16 PM     Agent Health Disabled   Show
> ⊘ Apr 29 9:57 PM      Agent Health Good       Show

bad health 0 — concerning health 0 — disabled health 100
good health 0 — unknown health 0
Alerts 0 — Critical Events 0 — Important Events 0

---

**cloudera** MANAGER   Clusters   Hosts   Diagnostics   Charts   Audits   Administration   Search   Support   student4

○ Flume4_10 (GeekBrains Cluster)

**Start Command**

Status ● Failed   Context Flume4_10   Jun 21, 4:16:55 PM   2.5m

Failed to start service.

▼ Completed 1 of 1 step(s).

○ Show All Steps   ● Show Only Failed Steps   ○ Show Running Steps

▼ ● Starting 1 roles on service                               Jun 21, 4:16:55 PM   2.5m
      Service did not start successfully; not all of the required roles started: only 0/1 roles started. Reasons : Service has only 0 Agent roles running instead of minimum required 1.

   ▼ ● Execute command Start this Agent on role Agent (node2)   Agent (node2)   Jun 21, 4:16:55 PM   2.5m
         Command aborted because of exception: Command timed-out after 150 seconds

      ▼ ● Start a role                                         Agent (node2)   Jun 21, 4:16:55 PM   2.5m
            Role is starting.
            Role Log

Close

> ⊘ Apr 29 10:16 PM     Agent Health Disabled   Show
> ⊕ Apr 29 9:57 PM      Agent Health Good       Show

---

**Я принял решение запустить с Flume 4_2 так как у него GOOD HEALTH**

# Naming the components on the current agent
Flume4_2.sources = ExecSource
Flume4_2.channels = MemChannel
Flume4_2.sinks = HdfsSink


# Describing/Configuring the source
Flume4_2.sources.ExecSource.type = exec
Flume4_2.sources.ExecSource.command = tailf /var/log/cloudera-scm-agent/cloudera-scm-agent.log
Flume4_2.sources.ExecSource.interceptors = TimestampInterceptor
Flume4_2.sources.ExecSource.interceptors.TimestampInterceptor.type = timestamp

# Describing/Configuring the HDFS sink
Flume4_2.sinks.HdfsSink.type = hdfs
Flume4_2.sinks.LoggerSink.hdfs.path= /flume/Flume4_2/%y-%m-%d/
Flume4_2.sinks.HdfsSink.hdfs.filePrefix = events

# Describing/Configuring the channel
Flume4_2.channels.MemChannel.type = memory
Flume4_2.channels.MemChannel.capacity = 10000
Flume4_2.channels.MemChannel.transactionCapacity = 10

# Bind the source and sink to the channel
Flume4_2.sources.ExecSource.channels = MemChannel
Flume4_2.sinks.HdfsSink.channel = MemChannel

**Save changes and restart**

**нету ничего**

**Попробую flume-7 в configuration file**

**Flume-7**



**Не работает**