LARGE SCALE DATA PROCESSING – DIGITAL ASSIGNMENT 2

CASE STUDY: TWITTER DATA STREAMING USING APACHE FLUME

[BEFORE the Twitter data streaming implementation, we understand the concepts of Flume]

**Flume: Introduction**

Apache Flume is a service for streaming logs into Hadoop. Flume is distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of streaming data into the Hadoop Distributed File System (HDFS). It has a simple and flexible architecture based on streaming data flows; and is robust and fault tolerant with tuneable reliability mechanisms for failover and recovery. YARN coordinates data ingest from Apache Flume and other services that deliver raw data into a Hadoop cluster.
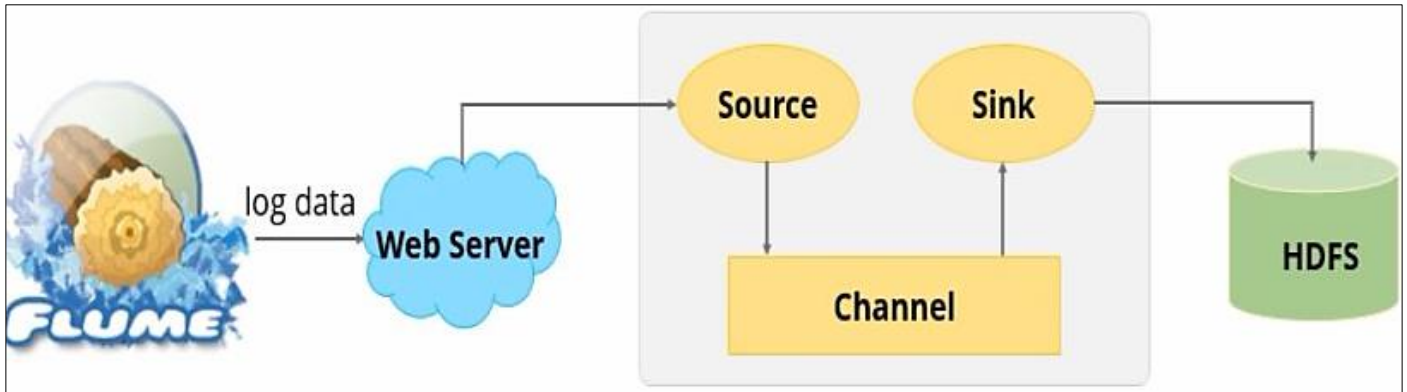
**Flume: Features**

Flume lets Hadoop users ingest high-volume streaming data into HDFS for storage. Specifically, Flume allows users to:

| FEATURES | DESCRIPTION |
|---|---|
| Stream data | Ingest streaming data from multiple sources into Hadoop for storage and analysis. |
| Insulate systems | Buffer storage platform from transient spikes, when the rate of incoming data exceeds the rate at which data can be written to the destination. |
| Guarantee data delivery | Flume NG uses channel-based transactions to guarantee reliable message delivery. When a message moves from one agent to another, two transactions are started, one on the agent that delivers the event and the other on the agent that receives the event. This ensures guaranteed delivery semantics. |
| Scale horizontally | To ingest new data streams and additional volume as needed. |

Enterprises use Flume's powerful streaming capabilities to land data from high-throughput streams in the HDFS. Typical sources of these streams are application logs, sensor and machine data, geo-location data and social media. These different types of data can be landed in Hadoop for future analysis using interactive queries in Apache Hive. Or they can feed business dashboards served ongoing data by Apache HBase.

Example usecase: Flume is used to log manufacturing operations. When one run of product comes off the line, it generates a log file about that run. Even if this occurs hundreds or thousands of times per day, the large volume log file data can stream through Flume into a tool for same-day analysis with Apache Storm or months or years of production runs can be stored in HDFS and analysed by a quality assurance engineer using Apache Hive.
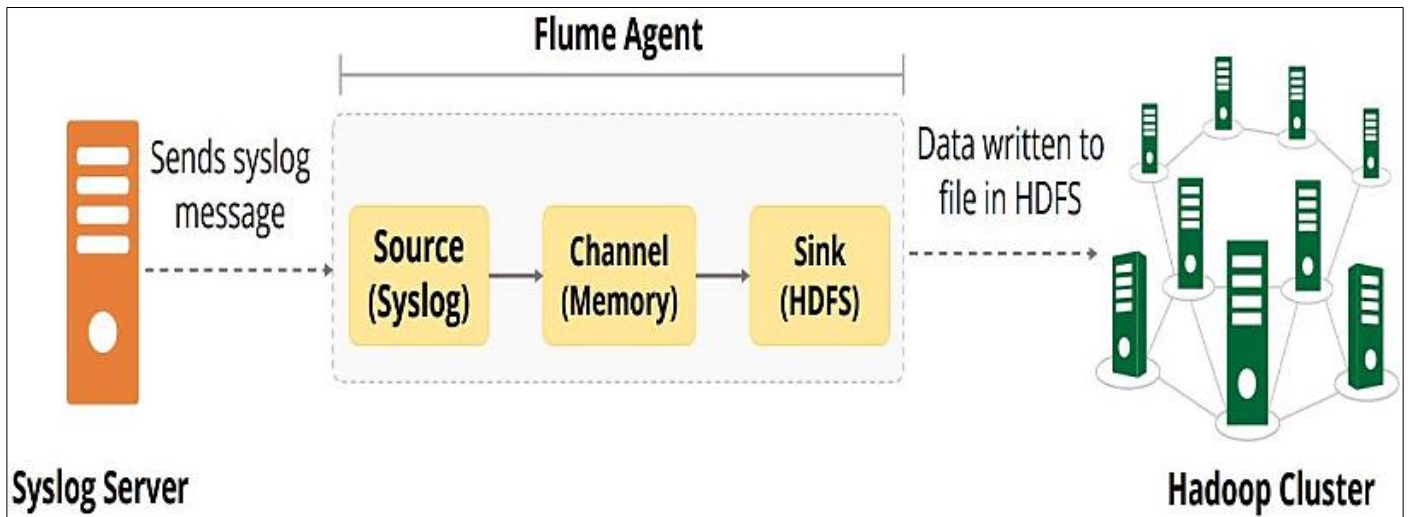
**Flume: Working and operations**

Flume's high-level architecture is built on a streamlined codebase that is easy to use and extend. The project is highly reliable, without the risk of data loss. Flume also supports dynamic reconfiguration without the need for a restart, which reduces downtime for its agents. The following components make up Apache Flume:

| COMPONENT | DEFINITION |
|-----------|------------|
| Event | A singular unit of data that is transported by Flume (typically a single log entry). |
| Source | The entity through which data enters into Flume. Sources either actively poll for data or passively wait for data to be delivered to them. A variety of sources allow data to be collected, such as log4j logs and syslogs. |
| Sink | The entity that delivers the data to the destination. A variety of sinks allow data to be streamed to a range of destinations. One example is the HDFS sink that writes events to HDFS. |
| Channel | The conduit between the Source and the Sink. Sources ingest events into the channel and the sinks drain the channel. |
| Agent | Any physical Java virtual machine running Flume. It is a collection of sources, sinks and channels. |
| Client | The entity that produces and transmits the Event to the Source operating within the Agent. |

Flume components interact in the following way:

1. A flow in Flume starts from the **Client**.
2. The **Client** transmits the **Event** to a **Source** operating within the **Agent**.
3. The **Source** receiving this **Event** then delivers it to one or more **Channels**.
4. One or more **Sinks** operating within the same **Agent** drains these **Channels**.
5. **Channels** decouple the ingestion rate from drain rate using the familiar producer-consumer model of data exchange.
6. When spikes in client side activity cause data to be generated faster than can be handled by the provisioned destination capacity can handle, the **Channel** size increases. This allows sources to continue normal operation for the duration of the spike.
7. The **Sink** of one **Agent** can be chained to the **Source** of another **Agent**. This chaining enables the creation of complex data flow topologies.

How is Flume easily scalable?

*Because Flume's distributed architecture requires no central coordination point. Each agent runs independently of others with no inherent single point of failure, and Flume can easily scale horizontally.*

**FLUME CASE STUDY: TWITTER DATA STREAMING**

**STEP 1:**

Setting up Apache Flume: We download and install **Flume 1.9**. We also set the appropriate JAVA_HOME path and the FLUME_HOME path in the .bashrc file.

To test if Flume has been installed correctly, we try executing ./flume-ng. On successful installation, we get the following output.

```
dhruvsgarg@dhruvsgarg:~$ ls
 anaconda3    Downloads                              Hadoop        NLP_try.ipynb    Public          Template
 Desktop      federated-averaging-tutorials-master   Music         nltk_data        summary.py      Videos
 Documents    Flume                                  NLPProject    Pictures         'TCS Internship'
dhruvsgarg@dhruvsgarg:~$ cd Flume/
dhruvsgarg@dhruvsgarg:~/Flume$ ls
bin         conf        doap_Flume.rdf  flume.conf  LICENSE   README.md       tools
CHANGELOG   DEVNOTES    docs            lib         NOTICE    RELEASE-NOTES
dhruvsgarg@dhruvsgarg:~/Flume$ cd bin/
dhruvsgarg@dhruvsgarg:~/Flume/bin$ ./flume-ng
Error: Unknown or unspecified command ''

Usage: ./flume-ng <command> [options]...

commands:
  help                    display this help text
  agent                   run a Flume agent
  avro-client             run an avro Flume client
  version                 show Flume version info

global options:
  --conf,-c <conf>        use configs in <conf> directory
  --classpath,-C <cp>     append to the classpath
  --dryrun,-d             do not actually start Flume, just print the command
  --plugins-path <dirs>   colon-separated list of plugins.d directories. See the
```

**STEP 2:**

Setting up Hadoop version 3.2: We download and install **Hadoop 3.1.2.** We also set the appropriate JAVA_HOME path and HADOOP_HOME path along with YARN_HOME, etc. We must also manually add content to the various .xml files to set the appropriate configuration settings.

```
dhruvsgarg@dhruvsgarg:~$ hadoop version
Hadoop 3.1.2
Source code repository https://github.com/apache/hadoop.git -r 1019dde65bcf12e05ef48ac71e84550d589e5d9a
Compiled by sunilg on 2019-01-29T01:39Z
Compiled with protoc 2.5.0
From source with checksum 64b8bdd4ca6e77cce75a93eb09ab2a9
This command was run using /home/dhruvsgarg/Hadoop/share/hadoop/common/hadoop-common-3.1.2.jar
dhruvsgarg@dhruvsgarg:~$ cd $FLUME_HOME
dhruvsgarg@dhruvsgarg:~/Flume$ ls
bin  CHANGELOG  conf  DEVNOTES  doap_Flume.rdf  docs  lib  LICENSE  NOTICE  README.md  RELEASE-NOTES  tools
```

**STEP 3:**

Creating a Twitter application: After logging into Twitter, we create a new application using https://apps.twitter.com. Here I have created an application named **dhruvFirstFlumeApp**.



**STEP 4:**

After creating the application, we can **access the keys and tokens** of the application.

There are 4 important application settings:

1. Consumer Key (API Key)
2. Consumer Secret (API Secret)
3. Access Token
4. Access Token Secret.

**STEP 5:**

We create a **flume.conf file** in Flume's root directory as shown below. As discussed in the Flume Architecture, we configure our Source, Sink and Channel. In this application, the <u>Source is Twitter</u>, from where we are streaming our data and our <u>Sink is HDFS</u>, where we are writing the data.



**Source properties:** In source configuration we are passing the Twitter source type as org.apache.flume.source.twitter.TwitterSource. Then, we are passing all the four tokens which we received from Twitter. At last in source configuration we are passing the keywords on which we are going to fetch the tweets.

```
13    TwitterAgent.sinks.HDFS.channel = MemChannel
14    TwitterAgent.sinks.HDFS.type = hdfs
15    TwitterAgent.sinks.HDFS.path = hdfs://localhost:9000/user/flume/tweets/
16    TwitterAgent.sinks.HDFS.fileType = DataStream
17    TwitterAgent.sinks.HDFS.writeFormat = Text
18    TwitterAgent.sinks.HDFS.batchSize = 1000
19    TwitterAgent.sinks.HDFS.rollSize = 0
20    TwitterAgent.sinks.HDFS.rollCount = 10000
21    TwitterAgent.sinks.HDFS.rollInterval = 600
```

**Sink properties:** In the sink configuration, we set the type, path, file type, batch size and roll parameters.

```
22
23    TwitterAgent.channels.MemChannel.type = memory
24    TwitterAgent.channels.MemChannel.capacity = 10000
25    TwitterAgent.channels.MemChannel.transactionCapacity = 100
26
```

**Channel properties:** In the channel properties, we set the memory, capacity and the capacity of each transaction.

## STEP 6:

After setting up all the required tools, we are **set for execution**. We perform the Twitter Data Streaming execution using the following command:

$FLUME_HOME/bin/flume-ng agent --conf ./conf/ -f $FLUME_HOME/flume.conf

```
dhruvsgarg@dhruvsgarg:~/Flume$ $FLUME_HOME/bin/flume-ng agent --conf ./conf/ -f $FLUME_HOME/flume.conf
Info: Sourcing environment configuration script /home/dhruvsgarg/Flume/conf/flume-env.sh
Info: Including Hadoop libraries found via (/home/dhruvsgarg/Hadoop/bin/hadoop) for HDFS access
```

We stop the execution by entering Ctrl+C after some time.

## STEP 7:

After some time, the **execution completes** and the data has been streamed into HDFS from Twitter.



## RESULT

Thus Flume was successfully used for Twitter Data Streaming and the output was stored in HDFS.

**OTHER APPLICATIONS OF APACHE FLUME**

- Flume is used in <u>e-commerce companies</u> to analyse the customer behaviour of different regions.
- It is used <u>to feed huge log data generated</u> by application servers into HDFS at a higher speed.
- Flume is used in <u>online analytics</u> application hence it is a backbone for real-time event processing.
- Flume efficiently <u>feed log data files from multiple web servers</u> into a centralized store i.e HDFS or HBase.
- Flume can <u>process data in-flight</u> using interceptors. These can be very useful for data masking or filtering.
- Flume is <u>distributed in nature</u> i.e agents can be installed on many machine.
- Flume is able to <u>collect real-time and batch mode data</u> from multiple servers and immediately feed into Hadoop.
- Flume is also used to <u>import huge volumes of event data</u> (instantaneous values) produced <u>by social networking sites</u> like Facebook and Twitter, and e-commerce websites like Amazon and Flipkart.