



Future Vision

FUTURE VISION BIE

By K B Hemanth Raj

Visit : <https://hemanthrajhemu.github.io>

A Small Contribution Would Support Us.

Dear Viewer,

Future Vision BIE is a free service and so that any Student/Research Personal **Can Access Free of Cost**.

If you would like to say **thanks**, you can make a **small contribution** to the author of this site.

Contribute whatever you feel this is worth to you. This gives **us support** & to bring **Latest Study Material** to you. After the Contribution Fill out this Form (<https://forms.gle/tw3T3bUVpLXL8omX7>). To Receive a **Paid E-Course for Free**, from our End within 7 Working Days.

Regards

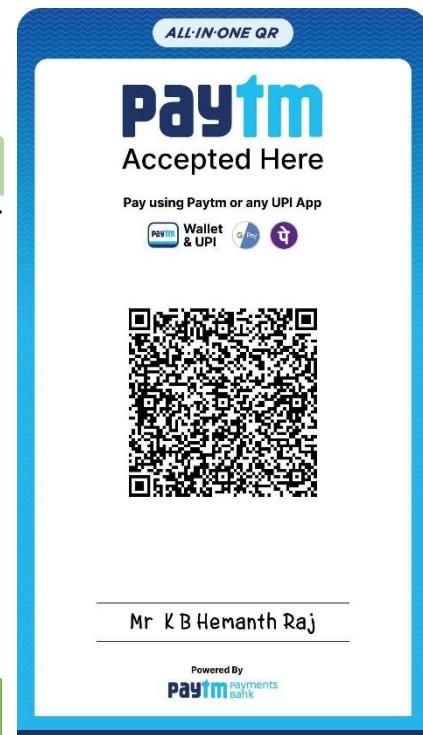
- K B Hemanth Raj (Admin)

Contribution Methods

UPI ID

1. futurevisionbie@oksbi
2. futurevisionbie@paytm

Scan & Pay



Account Transfer

Account Holder's Name: K B Hemanth Raj

Account Number: 39979402438

IFSC Code: SBIN0003982

MICR Code: 560002017

More Info: <https://hemanthrajhemu.github.io/Contribution/>

Gain Access to All Study Materials according to VTU,
CSE – Computer Science Engineering,
ISE – Information Science Engineering,
ECE - Electronics and Communication Engineering & MORE...

Stay Connected... get Updated... ask your queries...

Join Telegram to get Instant Updates: https://bit.ly/VTU_TELEGRAM

Contact: MAIL: futurevisionbie@gmail.com

INSTAGRAM: www.instagram.com/futurevisionbie/

WHATSSAPP SHARE: <https://bit.ly/FVBIESHARE>

Hadoop® 2 Quick-Start Guide

Learn the Essentials of Big
Data Computing in the Apache
Hadoop® 2 Ecosystem

Douglas Eadline



Pearson

<https://hemanthrajhemu.github.io>

Using Apache Flume to Acquire Data Streams	148
Flume Example Walk-Through	151
Manage Hadoop Workflows with Apache Oozie	154
Oozie Example Walk-Through	156
Using Apache HBase	163
HBase Data Model Overview	164
HBase Example Walk-Through	164
Summary and Additional Resources	169
8 Hadoop YARN Applications	171
YARN Distributed-Shell	171
Using the YARN Distributed-Shell	172
A Simple Example	174
Using More Containers	175
Distributed-Shell Examples with Shell Arguments	176
Structure of YARN Applications	178
YARN Application Frameworks	179
Distributed-Shell	180
Hadoop MapReduce	181
Apache Tez	181
Apache Giraph	181
Hoya: HBase on YARN	181
Dryad on YARN	182
Apache Spark	182
Apache Storm	182
Apache REEF: Retainable Evaluator Execution Framework	182
Hamster: Hadoop and MPI on the Same Cluster	183
Apache Flink: Scalable Batch and Stream Data Processing	183
Apache Slider: Dynamic Application Management	183
Summary and Additional Resources	184
9 Managing Hadoop with Apache Ambari	185
Quick Tour of Apache Ambari	186
Dashboard View	186

Services View	189
Hosts View	191
Admin View	193
Views View	193
Admin Pull-Down Menu	194
Managing Hadoop Services	194
Changing Hadoop Properties	198
Summary and Additional Resources	204
10 Basic Hadoop Administration Procedures	205
Basic Hadoop YARN Administration	206
Decommissioning YARN Nodes	206
YARN WebProxy	206
Using the JobHistoryServer	207
Managing YARN Jobs	207
Setting Container Memory	207
Setting Container Cores	208
Setting MapReduce Properties	208
Basic HDFS Administration	208
The NameNode User Interface	208
Adding Users to HDFS	211
Perform an FSCK on HDFS	212
Balancing HDFS	213
HDFS Safe Mode	214
Decommissioning HDFS Nodes	214
SecondaryNameNode	214
HDFS Snapshots	215
Configuring an NFSv3 Gateway to HDFS	217
Capacity Scheduler Background	220
Hadoop Version 2 MapReduce Compatibility	222
Enabling ApplicationMaster Restarts	222
Calculating the Capacity of a Node	222
Running Hadoop Version 1 Applications	224
Summary and Additional Resources	225
A Book Webpage and Code Download	227

7. Essential Hadoop Tools

In This Chapter:

- The Pig scripting tool is introduced as a way to quickly examine data both locally and on a Hadoop cluster.
- The Hive SQL-like query tool is explained using two examples.
- The Sqoop RDBMS tool is used to import and export data from MySQL to/from HDFS.
- The Flume streaming data transport utility is configured to capture weblog data into HDFS.
- The Oozie workflow manager is used to run basic and complex Hadoop workflows.
- The distributed HBase database is used to store and access data on a Hadoop cluster.

The Hadoop ecosystem offers many tools to help with data input, high-level processing, workflow management, and creation of huge databases. Each tool is managed as a separate Apache Software foundation project, but is designed to operate with the core Hadoop services including HDFS, YARN, and MapReduce. Background on each tool is provided in this chapter, along with a *start to finish* example.

USING APACHE PIG

Apache Pig is a high-level language that enables programmers to write complex MapReduce transformations using a simple scripting language. Pig Latin (the actual language) defines a set of transformations on a data set such as aggregate, join, and sort. Pig is often used to extract, transform, and load (ETL) data pipelines, quick research on raw data, and iterative data processing.

Apache Pig has several usage modes. The first is a local mode in which all processing is done on the local machine. The non-local (cluster) modes are MapReduce and Tez. These modes execute the job on the cluster using either the MapReduce engine or the optimized Tez engine. (Tez, which is Hindi for “speed,” optimizes multistep Hadoop jobs such as those found in many Pig queries.) There are also interactive and batch modes available; they enable Pig applications to be developed locally in interactive modes, using small amounts of data, and then run at scale on the cluster in a production mode. The modes are summarized in Table 7.1.

	Local Mode	Tez Local Mode	MapReduce Mode	Tez Mode
Interactive Mode	Yes	Experimental	Yes	Yes
Batch Mode	Yes	Experimental	Yes	Yes

Table 7.1 Apache Pig Usage Modes

Pig Example Walk-Through

For this example, the following software environment is assumed. Other environments should work in a similar fashion.

- OS: Linux

- Platform: RHEL 6.6
- Hortonworks HDP 2.2 with Hadoop version: 2.6
- Pig version: 0.14.0

If you are using the pseudo-distributed installation from [Chapter 2, “Installation Recipes,”](#) instructions for installing Pig are provided in that chapter. More information on installing Pig by hand can be found on the Pig website: <http://pig.apache.org/#Getting+Started>. Apache Pig is also installed as part of the Hortonworks HDP Sandbox.

In this simple example, Pig is used to extract user names from the `/etc/passwd` file. A full description of the Pig Latin language is beyond the scope of this introduction, but more information about Pig can be found at <http://pig.apache.org/docs/r0.14.0/start.html>. The following example assumes the user is `hdfs`, but any valid user with access to HDFS can run the example.

To begin the example, copy the `passwd` file to a working directory for local Pig operation:

```
$ cp /etc/passwd .
```

Next, copy the data file into HDFS for Hadoop MapReduce operation:

```
$ hdfs dfs -put passwd passwd
```

You can confirm the file is in HDFS by entering the following command:

```
hdfs dfs -ls passwd
-rw-r--r-- 2 hdfs hdfs 2526 2015-03-17 11:08 passwd
```

In the following example of local Pig operation, all processing is done on the local machine (Hadoop is not used). First, the interactive command line is started:

```
$ pig -x local
```

If Pig starts correctly, you will see a `grunt>` prompt. You may also see a bunch of INFO messages, which you can ignore. Next, enter the following commands to load the `passwd` file and then grab the user name and dump it to the terminal. Note that Pig commands must end with a semicolon (`;`).

```
grunt> A = load 'passwd' using PigStorage(':' );
grunt> B = foreach A generate $0 as id;
grunt> dump B;
```

The processing will start and a list of user names will be printed to the screen. To exit the interactive session, enter the command `quit`.

```
$ grunt> quit
```

To use Hadoop MapReduce, start Pig as follows (or just enter `pig`):

```
$ pig -x mapreduce
```

The same sequence of commands can be entered at the `grunt>` prompt. You may wish to change the `$0` argument to pull out other items in the `passwd` file. In the case of this simple script, you will notice that the MapReduce version takes much longer. Also, because we are running this application under Hadoop, make sure the file is placed in HDFS.

If you are using the Hortonworks HDP distribution with `tez` installed, the `tez` engine can be used as follows:

```
$ pig -x tez
```

Pig can also be run from a script. An example script (`id.pig`) is available from the example code download (see [Appendix A, “Book Webpage and Code Download”](#)). This script, which is repeated here, is designed to do the same things as the interactive version:

[**Click here to view code image**](#)

```
/* id.pig */
A = load 'passwd' using PigStorage(':'); -- load the passwd file
B = foreach A generate $0 as id; -- extract the user IDs
dump B;
store B into 'id.out'; -- write the results to a directory name id.out
```

Comments are delineated by `/* */` and `--` at the end of a line. The script will create a directory called `id.out` for the results. First, ensure that the `id.out` directory is not in your local directory, and then start Pig with the script on the command line:

```
$ /bin/rm -r id.out/
$ pig -x local id.pig
```

If the script worked correctly, you should see at least one data file with the results and a zero-length file with the name `_SUCCESS`. To run the MapReduce version, use the same procedure; the only difference is that now all reading and writing takes place in HDFS.

```
$ hdfs dfs -rm -r id.out
$ pig id.pig
```

If Apache `tez` is installed, you can run the example script using the `-x tez` option. You can learn more about writing Pig scripts at <http://pig.apache.org/docs/r0.14.0/start.html>.

USING APACHE HIVE

Apache Hive is a data warehouse infrastructure built on top of Hadoop for providing data summarization, ad hoc queries, and the analysis of large data sets using a SQL-like language called HiveQL. Hive is considered the de facto standard for interactive SQL queries over petabytes of data using Hadoop and offers the following features:

- Tools to enable easy data extraction, transformation, and loading (ETL)

- A mechanism to impose structure on a variety of data formats
- Access to files stored either directly in HDFS or in other data storage systems such as HBase
- Query execution via MapReduce and Tez (optimized MapReduce)

Hive provides users who are already familiar with SQL the capability to query the data on Hadoop clusters. At the same time, Hive makes it possible for programmers who are familiar with the MapReduce framework to add their custom mappers and reducers to Hive queries. Hive queries can also be dramatically accelerated using the Apache Tez framework under YARN in Hadoop version 2.

Hive Example Walk-Through

For this example, the following software environment is assumed. Other environments should work in a similar fashion.

- OS: Linux
- Platform: RHEL 6.6
- Hortonworks HDP 2.2 with Hadoop version: 2.6
- Hive version: 0.14.0

If you are using the pseudo-distributed installation from [Chapter 2](#), instructions for installing Hive are provided in that chapter. More information on installation can be found on the Hive website: <http://hive.apache.org>. Hive is also installed as part of the Hortonworks HDP Sandbox. Although the following example assumes the user is `hdfs`, any valid user with access to HDFS can run the example.

To start Hive, simply enter the `hive` command. If Hive starts correctly, you should get a `hive>` prompt.

Click here to view code image

```
$ hive  
(some messages may show up here)  
hive>
```

As a simple test, create and drop a table. Note that Hive commands must end with a semicolon (;).

Click here to view code image

```
hive> CREATE TABLE pokes (foo INT, bar STRING);  
OK  
Time taken: 1.705 seconds  
hive> SHOW TABLES;  
OK  
pokes  
Time taken: 0.174 seconds, Fetched: 1 row(s)  
hive> DROP TABLE pokes;  
OK  
Time taken: 4.038 seconds
```

A more detailed example can be developed using a web server log file to summarize message types. First, create a table using the following command:

[**Click here to view code image**](#)

```
hive> CREATE TABLE logs(t1 string, t2 string, t3 string, t4 string, t5 string, t6 string, t7 string)
ROW FORMAT DELIMITED FIELDS TERMINATED BY '';
OK
Time taken: 0.129 seconds
```

Next, load the data—in this case, from the `sample.log` file. This file is available from the example code download (see [Appendix A](#)). Note that the file is found in the local directory and not in HDFS.

[**Click here to view code image**](#)

```
hive> LOAD DATA LOCAL INPATH 'sample.log' OVERWRITE INTO TABLE logs;
Loading data to table default.logs
Table default.logs stats: [numFiles=1, numRows=0, totalSize=99271, rawDataSize=0]
OK
Time taken: 0.953 seconds
```

Finally, apply the select step to the file. Note that this invokes a Hadoop MapReduce operation. The results appear at the end of the output (e.g., totals for the message types `DEBUG`, `ERROR`, and so on).

[**Click here to view code image**](#)

```
hive> SELECT t4 AS sev, COUNT(*) AS cnt FROM logs WHERE t4 LIKE '[%]' GROUP BY t4;
Query ID = hdfs_20150327130000_d1e1a265-a5d7-4ed8-b785-2c6569791368
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1427397392757_0001, Tracking URL = http://norbert:8088/proxy/application_1427397392757_0001/
Kill Command = /opt/hadoop-2.6.0/bin/hadoop job -kill job_1427397392757_0001
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2015-03-27 13:00:17,399 Stage-1 map = 0%, reduce = 0%
2015-03-27 13:00:26,100 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 2.14 sec
2015-03-27 13:00:34,979 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 4.07 sec
MapReduce Total cumulative CPU time: 4 seconds 70 msec
```

```
Ended Job = job_1427397392757_0001
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 4.07 sec HDFS Read: 106384
HDFS Write: 63 SUCCESS
Total MapReduce CPU Time Spent: 4 seconds 70 msec
OK
[DEBUG] 434
[ERROR] 3
[FATAL] 1
[INFO] 96
[TRACE] 816
[WARN] 4
Time taken: 32.624 seconds, Fetched: 6 row(s)
```

To exit Hive, simply type `exit;:`

```
hive> exit;
```

A More Advanced Hive Example

A more advanced usage case from the Hive documentation can be developed using the movie rating data files obtained from the GroupLens Research (<http://grouplens.org/datasets/movielens>) webpage. The data files are collected from Movie-Lens website (<http://movielens.org>). The files contain various numbers of movie reviews, starting at 100,000 and going up to 20 million entries. The data file and queries used in the following example are available from the book website (see [Appendix A](#)).

In this example, 100,000 records will be transformed from `userid, movieid, rating, unixtime` to `userid, movieid, rating, and weekday` using Apache Hive and a Python program (i.e., the UNIX time notation will be transformed to the day of the week). The first step is to download and extract the data:

[Click here to view code image](#)

```
$ wget http://files.grouplens.org/datasets/movielens/ml-100k.zip
$ unzip ml-100k.zip
$ cd ml-100k
```

Before we use Hive, we will create a short Python program called `weekday_mapper.py` with following contents:

[Click here to view code image](#)

```
import sys
import datetime

for line in sys.stdin:
    line = line.strip()
    userid, movieid, rating, unixtime = line.split('\t')
```

```
weekday = datetime.datetime.fromtimestamp(float(unixtime)).isoweekday()
print '\t'.join([userid, movieid, rating, str(weekday)])LOAD DATA LOCAL INPATH
'./u.data' OVERWRITE INTO TABLE u_data;
```

Next, start Hive and create the data table (`u_data`) by entering the following at the `hive>` prompt:

```
CREATE TABLE u_data (
  userid INT,
  movieid INT,
  rating INT,
  unixtime STRING)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '\t'
STORED AS TEXTFILE;
```

Load the movie data into the table with the following command:

[Click here to view code image](#)

```
hive> LOAD DATA LOCAL INPATH './u.data' OVERWRITE INTO TABLE u_data;
```

The number of rows in the table can be reported by entering the following command:

[Click here to view code image](#)

```
hive > SELECT COUNT(*) FROM u_data;
```

This command will start a single MapReduce job and should finish with the following lines:

[Click here to view code image](#)

```
...
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1  Cumulative CPU: 2.26 sec  HDFS Read: 1979380
HDFS Write: 7 SUCCESS
Total MapReduce CPU Time Spent: 2 seconds 260 msec
OK
100000
Time taken: 28.366 seconds, Fetched: 1 row(s)
```

Now that the table data are loaded, use the following command to make the new table (`u_data_new`):

[Click here to view code image](#)

```
hive> CREATE TABLE u_data_new (
  userid INT,
```

```
movieid INT,  
rating INT,  
weekday INT)  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY '\t';
```

The next command adds the `weekday_mapper.py` to Hive resources:

[**Click here to view code image**](#)

```
hive> add FILE weekday_mapper.py;
```

Once `weekday_mapper.py` is successfully loaded, we can enter the transformation query:

[**Click here to view code image**](#)

```
hive> INSERT OVERWRITE TABLE u_data_new  
SELECT  
    TRANSFORM (userid, movieid, rating, unixtime)  
    USING 'python weekday_mapper.py'  
    AS (userid, movieid, rating, weekday)  
FROM u_data;
```

If the transformation was successful, the following final portion of the output should be displayed:

[**Click here to view code image**](#)

```
...  
Table default.u_data_new stats: [numFiles=1, numRows=100000, totalSize=1179173,  
rawDataSize=1079173]  
MapReduce Jobs Launched:  
Stage-Stage-1: Map: 1 Cumulative CPU: 3.44 sec HDFS Read: 1979380 HDFS Write:  
1179256 SUCCESS  
Total MapReduce CPU Time Spent: 3 seconds 440 msec  
OK  
Time taken: 24.06 seconds
```

The final query will sort and group the reviews by weekday:

[**Click here to view code image**](#)

```
hive> SELECT weekday, COUNT(*) FROM u_data_new GROUP BY weekday;
```

Final output for the review counts by weekday should look like the following:

[**Click here to view code image**](#)

...
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 2.39 sec HDFS Read: 1179386
HDFS Write: 56 SUCCESS
Total MapReduce CPU Time Spent: 2 seconds 390 msec
OK
1 13278
2 14816
3 15426
4 13774
5 17964
6 12318
7 12424
Time taken: 22.645 seconds, Fetched: 7 row(s)

As shown previously, you can remove the tables used in this example with the `DROP TABLE` command. In this case, we are also using the `-e` command-line option. Note that queries can be loaded from files using the `-f` option as well.

[**Click here to view code image**](#)

```
$ hive -e 'drop table u_data_new'  
$ hive -e 'drop table u_data'
```

USING APACHE SQOOP TO ACQUIRE RELATIONAL DATA

Sqoop is a tool designed to transfer data between Hadoop and relational databases. You can use Sqoop to import data from a relational database management system (RDBMS) into the Hadoop Distributed File System (HDFS), transform the data in Hadoop, and then export the data back into an RDBMS.

Sqoop can be used with any Java Database Connectivity (JDBC)-compliant database and has been tested on Microsoft SQL Server, PostgreSQL, MySQL, and Oracle. In version 1 of Sqoop, data were accessed using connectors written for specific databases. Version 2 (in beta) does not support connectors or version 1 data transfer from a RDBMS directly to Hive or HBase, or data transfer from Hive or HBase to your RDBMS. Instead, version 2 offers more generalized ways to accomplish these tasks.

The remainder of this section provides a brief overview of how Sqoop works with Hadoop. In addition, a basic Sqoop example walk-through is demonstrated. To fully explore Sqoop, more information can be found by consulting the Sqoop project website: <http://sqoop.apache.org>

Apache Sqoop Import and Export Methods

[Figure 7.1](#) describes the Sqoop data import (to HDFS) process. The data import is done in two steps. In the first step, shown in the figure, Sqoop examines the database to gather the necessary metadata for the data to be imported. The second step is a map-only (no reduce step) Hadoop job that Sqoop submits to the cluster. This job does the actual data transfer using the metadata

captured in the previous step. Note that each node doing the import must have access to the database.

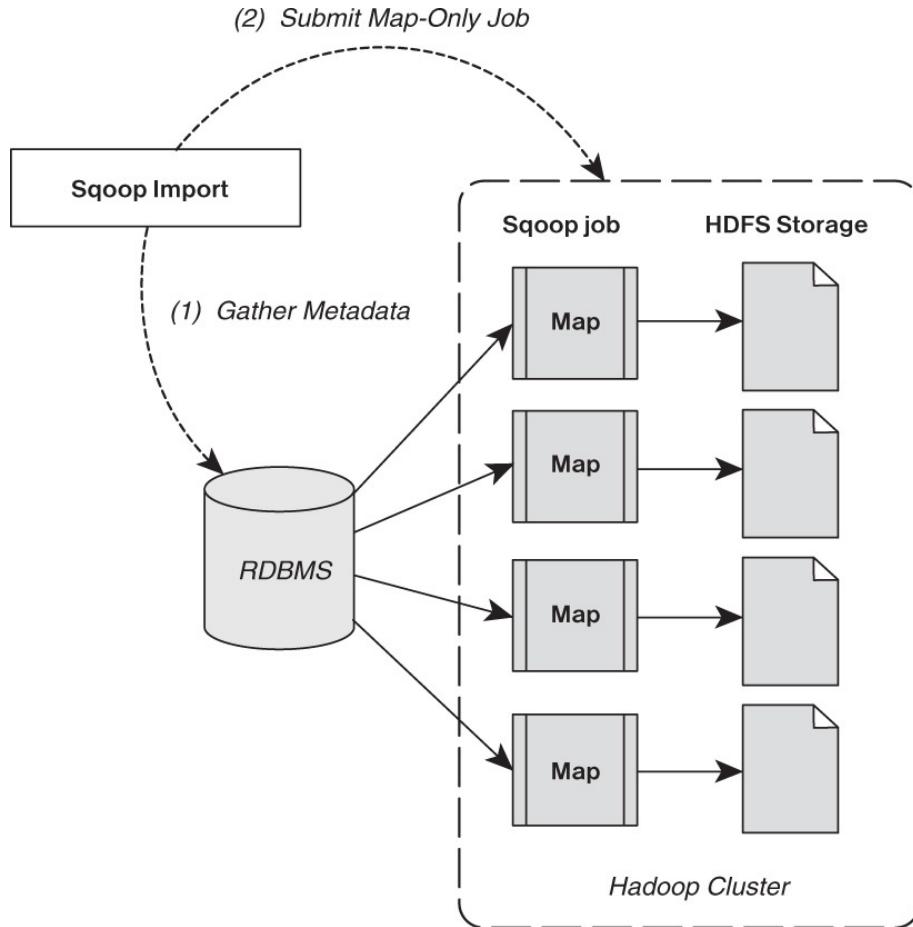


Figure 7.1 Two-step Apache Sqoop data import method (Adapted from Apache Sqoop Documentation)

The imported data are saved in an HDFS directory. Sqoop will use the database name for the directory, or the user can specify any alternative directory where the files should be populated. By default, these files contain comma-delimited fields, with new lines separating different records. You can easily override the format in which data are copied over by explicitly specifying the field separator and record terminator characters. Once placed in HDFS, the data are ready for processing.

Data export from the cluster works in a similar fashion. The export is done in two steps, as shown in [Figure 7.2](#). As in the import process, the first step is to examine the database for metadata. The export step again uses a map-only Hadoop job to write the data to the database. Sqoop divides the input data set into splits, then uses individual map tasks to push the splits to the database. Again, this process assumes the map tasks have access to the database.

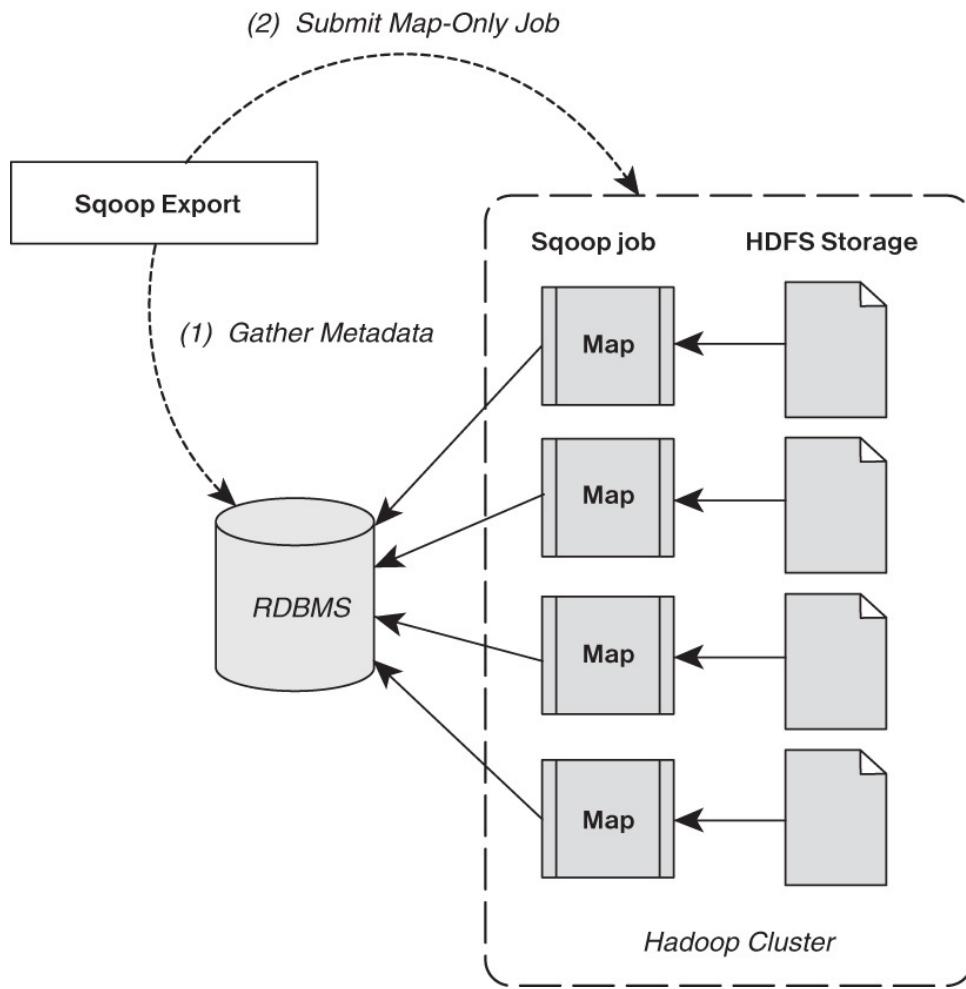


Figure 7.2 Two-step Sqoop data export method (Adapted from Apache Sqoop Documentation)

Apache Sqoop Version Changes

Sqoop Version 1 uses specialized connectors to access external systems. These connectors are often optimized for various RDBMSs or for systems that do not support JDBC. Connectors are plug-in components based on Sqoop's extension framework and can be added to any existing Sqoop installation. Once a connector is installed, Sqoop can use it to efficiently transfer data between Hadoop and the external store supported by the connector. By default, Sqoop version 1 includes connectors for popular databases such as MySQL, PostgreSQL, Oracle, SQL Server, and DB2. It also supports direct transfer to and from the RDBMS to HBase or Hive.

In contrast, to streamline the Sqoop input methods, Sqoop version 2 no longer supports specialized connectors or direct import into HBase or Hive. All imports and exports are done through the JDBC interface. [Table 7.2](#) summarizes the changes from version 1 to version 2. Due to these changes, any new development should be done with Sqoop version 2.

Feature	Sqoop Version 1	Sqoop Version 2
Connectors for all major RDBMSs	Supported.	Not supported. Use the generic JDBC connector.
Kerberos security integration	Supported.	Not supported.
Data transfer from RDBMS to Hive or HBase	Supported.	Not supported. First import data from RDBMS into HDFS, then load data into Hive or HBase manually.
Data transfer from Hive or HBase to RDBMS	Not supported. First export data from Hive or HBase into HDFS, and then use Sqoop for export.	Not supported. First export data from Hive or HBase into HDFS, then use Sqoop for export.

Table 7.2 Apache Sqoop Version Comparison

Sqoop Example Walk-Through

The following simple example illustrates use of Sqoop. It can be used as a foundation from which to explore the other capabilities offered by Apache Sqoop. The following steps will be performed:

1. Download Sqoop.
2. Download and load sample MySQL data.
3. Add Sqoop user permissions for the local machine and cluster.
4. Import data from MySQL to HDFS.
5. Export data from HDFS to MySQL.

For this example, the following software environment is assumed. Other environments should work in a similar fashion.

- OS: Linux
- Platform: RHEL 6.6
- Hortonworks HDP 2.2 with Hadoop version: 2.6
- Sqoop version: 1.4.5
- A working installation of MySQL on the host

If you are using the pseudo-distributed installation from [Chapter 2](#) or want to install Sqoop by hand, see the installation instructions on the Sqoop website: <http://sqoop.apache.org>. Sqoop is also installed as part of the Hortonworks HDP Sandbox.

Step 1: Download Sqoop and Load Sample MySQL Database

If you have not done so already, make sure Sqoop is installed on your cluster. Sqoop is needed on only a single node in your cluster. This Sqoop node will then serve as an entry point for all connecting Sqoop clients. Because the Sqoop node is a Hadoop MapReduce client, it requires both a Hadoop installation and access to HDFS.

To install Sqoop using the HDP distribution RPM files, simply enter:

[**Click here to view code image**](#)

```
# yum install sqoop sqoop-metastore
```

For this example, we will use the world example database from the MySQL site (<http://dev.mysql.com/doc/world-setup/en/index.html>). This database has three tables:

- Country: information about countries of the world
- City: information about some of the cities in those countries
- CountryLanguage: languages spoken in each country

To get the database, use `wget` to download and then extract the file:

[**Click here to view code image**](#)

```
$ wget http://downloads.mysql.com/docs/world_innodb.sql.gz
$ gunzip world_innodb.sql.gz
```

Next, log into MySQL (assumes you have privileges to create a database) and import the desired database by following these steps:

[**Click here to view code image**](#)

```
$ mysql -u root -p
mysql> CREATE DATABASE world;
mysql> USE world;
mysql> SOURCE world_innodb.sql;
mysql> SHOW TABLES;
+-----+
| Tables_in_world |
+-----+
| City      |
| Country   |
| CountryLanguage |
+-----+
3 rows in set (0.01 sec)
```

The following MySQL command will let you see the table details (output omitted for clarity):

[**Click here to view code image**](#)

```
mysql> SHOW CREATE TABLE Country;
mysql> SHOW CREATE TABLE City;
mysql> SHOW CREATE TABLE CountryLanguage;
```

Step 2: Add Sqoop User Permissions for the Local Machine and Cluster

In MySQL, add the following privileges for user `sqoop` to MySQL. Note that you must use both the local host name and the cluster subnet for Sqoop to work properly. Also, for the purposes of this example, the `sqoop` password is `sqoop`.

[Click here to view code image](#)

```
mysql> GRANT ALL PRIVILEGES ON world.* To 'sqoop'@'limulus' IDENTIFIED BY
'sqoop';
mysql> GRANT ALL PRIVILEGES ON world.* To 'sqoop'@'10.0.0.%' IDENTIFIED BY
'sqoop';
mysql> quit
```

Next, log in as `sqoop` to test the permissions:

```
$ mysql -u sqoop -p
mysql> USE world;
mysql> SHOW TABLES;
+-----+
| Tables_in_world |
+-----+
| City      |
| Country   |
| CountryLanguage |
+-----+
3 rows in set (0.01 sec)
```

```
mysql> quit
```

Step 3: Import Data Using Sqoop

As a test, we can use Sqoop to list databases in MySQL. The results appear after the warnings at the end of the output. Note the use of local host name (`limulus`) in the JDBC statement.

[Click here to view code image](#)

```
$ sqoop list-databases --connect jdbc:mysql://limulus/world --username sqoop --password sqoop
Warning: /usr/lib/sqoop/..//accumulo does not exist! Accumulo imports will fail.
Please set $ACCUMULO_HOME to the root of your Accumulo installation.
14/08/18 14:38:55 INFO sqoop.Sqoop: Running Sqoop version: 1.4.4.2.1.2.1-471
14/08/18 14:38:55 WARN tool.BaseSqoopTool: Setting your password on the
command-line is insecure. Consider using -P instead.
14/08/18 14:38:55 INFO manager.MySQLManager: Preparing to use a MySQL streaming
resultset.
```

```
information_schema  
test  
world
```

In a similar fashion, you can use Sqoop to connect to MySQL and list the tables in the world database:

[**Click here to view code image**](#)

```
sqoop list-tables --connect jdbc:mysql://limulus/world --username sqoop --password sqoop  
...  
14/08/18 14:39:43 INFO sqoop.Sqoop: Running Sqoop version: 1.4.4.2.1.2.1-471  
14/08/18 14:39:43 WARN tool.BaseSqoopTool: Setting your password on the  
command-line is insecure. Consider using -P instead.  
14/08/18 14:39:43 INFO manager.MySQLManager: Preparing to use a MySQL streaming  
resultset.  
City  
Country  
CountryLanguage
```

To import data, we need to make a directory in HDFS:

[**Click here to view code image**](#)

```
$ hdfs dfs -mkdir sqoop-mysql-import
```

The following command imports the Country table into HDFS. The option `-table` signifies the table to import, `--target-dir` is the directory created previously, and `-m 1` tells Sqoop to use one map task to import the data.

[**Click here to view code image**](#)

```
$ sqoop import --connect jdbc:mysql://limulus/world --username sqoop --password sqoop --  
table Country -m 1 --target-dir /user/hdfs/sqoop-mysql-import/country  
...  
14/08/18 16:47:15 INFO mapreduce.ImportJobBase: Transferred 30.752 KB in  
12.7348 seconds  
(2.4148 KB/sec)  
14/08/18 16:47:15 INFO mapreduce.ImportJobBase: Retrieved 239 records.
```

The import can be confirmed by examining HDFS:

[**Click here to view code image**](#)

```
$ hdfs dfs -ls sqoop-mysql-import/country  
Found 2 items  
-rw-r--r-- 2 hdfs hdfs 0 2014-08-18 16:47 sqoop-mysql-import/  
world/_SUCCESS
```

```
-rw-r--r-- 2 hdfs hdfs 31490 2014-08-18 16:47 sqoop-mysql-import/world/part-m-00000
```

The file can be viewed using the `hdfs dfs -cat` command:

[**Click here to view code image**](#)

```
$ hdfs dfs -cat sqoop-mysql-import/country/part-m-00000
ABW,Aruba,North America,Caribbean,193.0,null,103000,78.4,828.0,793.0,Aruba,
Nonmetropolitan
Territory of The Netherlands,Beatrix,129,AW
...
ZWE,Zimbabwe,Africa,Eastern Africa,390757.0,1980,11669000,37.8,5951.0,8670.0,
Zimbabwe,
Republic,Robert G. Mugabe,4068,ZW
```

To make the Sqoop command more convenient, you can create an options file and use it on the command line. Such a file enables you to avoid having to rewrite the same options. For example, a file called `world-options.txt` with the following contents will include the `import` command, `--connect`, `--username`, and `--password` options:

```
import
--connect
jdbc:mysql://limulus/world
--username
sqoop
--password
sqoop
```

The same import command can be performed with the following shorter line:

[**Click here to view code image**](#)

```
$ sqoop --options-file world-options.txt --table City -m 1 --target-dir /user/hdfs/sqoop-mysql-import/city
```

It is also possible to include an SQL Query in the import step. For example, suppose we want just cities in Canada:

[**Click here to view code image**](#)

```
SELECT ID,Name from City WHERE CountryCode='CAN'
```

In such a case, we can include the `--query` option in the Sqoop import request. The `--query` option also needs a variable called `$CONDITIONS`, which will be explained next. In the following query example, a single mapper task is designated with the `-m 1` option:

[**Click here to view code image**](#)

```
sqoop --options-file world-options.txt -m 1 --target-dir /user/hdfs/sqoop-mysql-import/canada-city --query "SELECT ID,Name from City WHERE CountryCode='CAN' AND \$CONDITIONS"
```

Inspecting the results confirms that only cities from Canada have been imported:

[**Click here to view code image**](#)

```
$ hdfs dfs -cat sqoop-mysql-import/canada-city/part-m-00000
```

1810,MontrÃ©al

1811,Calgary

1812,Toronto

...

1856,Sudbury

1857,Kelowna

1858,Barrie

Since there was only one mapper process, only one copy of the query needed to be run on the database. The results are also reported in a single file (`part-m-0000`).

Multiple mappers can be used to process the query if the `--split-by` option is used. The `split-by` option is used to parallelize the SQL query. Each parallel task runs a subset of the main query, with the results of each sub-query being partitioned by bounding conditions inferred by Sqoop.

Your query must include the token `$CONDITIONS` that each Sqoop process will replace with a unique condition expression based on the `--split-by` option. Note that `$CONDITIONS` is not an environment variable. Although Sqoop will try to create balanced sub-queries based on the range of your primary key, it may be necessary to split on another column if your primary key is not uniformly distributed.

The following example illustrates the use of the `--split-by` option. First, remove the results of the previous query:

[**Click here to view code image**](#)

```
$ hdfs dfs -rm -r -skipTrash sqoop-mysql-import/canada-city
```

Next, run the query using four mappers (`-m 4`), where we split by the ID number (`--split-by ID`):

[**Click here to view code image**](#)

```
sqoop --options-file world-options.txt -m 4 --target-dir /user/hdfs/sqoop-mysql-import/canada-city --query "SELECT ID,Name from City WHERE CountryCode='CAN' AND \$CONDITIONS" --split-by ID
```

If we look at the number of results files, we find four files corresponding to the four mappers we requested in the command:

[Click here to view code image](#)

```
$ hdfs dfs -ls sqoop-mysql-import/canada-city
Found 5 items
-rw-r--r-- 2 hdfs hdfs 0 2014-08-18 21:31 sqoop-mysql-import/
canada-city/_SUCCESS
-rw-r--r-- 2 hdfs hdfs 175 2014-08-18 21:31 sqoop-mysql-import/canada-city/
part-m-00000
-rw-r--r-- 2 hdfs hdfs 153 2014-08-18 21:31 sqoop-mysql-import/canada-city/
part-m-00001
-rw-r--r-- 2 hdfs hdfs 186 2014-08-18 21:31 sqoop-mysql-import/canada-city/
part-m-00002
-rw-r--r-- 2 hdfs hdfs 182 2014-08-18 21:31 sqoop-mysql-import/canada-city/
part-m-00003
```

Step 4: Export Data from HDFS to MySQL

Sqoop can also be used to export data from HDFS. The first step is to create tables for exported data. There are actually two tables needed for each exported table. The first table holds the exported data (`CityExport`), and the second is used for staging the exported data (`CityExportStaging`). Enter the following MySQL commands to create these tables:

[Click here to view code image](#)

```
mysql> CREATE TABLE 'CityExport' (
    'ID' int(11) NOT NULL AUTO_INCREMENT,
    'Name' char(35) NOT NULL DEFAULT '',
    'CountryCode' char(3) NOT NULL DEFAULT '',
    'District' char(20) NOT NULL DEFAULT '',
    'Population' int(11) NOT NULL DEFAULT '0',
    PRIMARY KEY ('ID'));
mysql> CREATE TABLE 'CityExportStaging' (
    'ID' int(11) NOT NULL AUTO_INCREMENT,
    'Name' char(35) NOT NULL DEFAULT '',
    'CountryCode' char(3) NOT NULL DEFAULT '',
    'District' char(20) NOT NULL DEFAULT '',
    'Population' int(11) NOT NULL DEFAULT '0',
    PRIMARY KEY ('ID'));
```

Next, create a `cities-export-options.txt` file similar to the `world-options.txt` created previously, but use the `export` command instead of the `import` command.

The following command will export the cities data we previously imported back into MySQL:

[Click here to view code image](#)

```
sqoop --options-file cities-export-options.txt --table CityExport --staging-table
CityExportStaging --clear-staging-table -m 4 --export-dir /user/hdfs/sqoop-mysql-import/city
```

Finally, to make sure everything worked correctly, check the table in MySQL to see if the cities are in the table:

[Click here to view code image](#)

```
$ mysql> select * from CityExport limit 10;
+----+-----+-----+-----+
| ID | Name      | CountryCode | District    | Population |
+----+-----+-----+-----+
| 1  | Kabul      | AFG        | Kabul       | 1780000   |
| 2  | Qandahar   | AFG        | Qandahar   | 237500    |
| 3  | Herat      | AFG        | Herat      | 186800    |
| 4  | Mazar-e-Sharif | AFG        | Balkh      | 127800    |
| 5  | Amsterdam   | NLD        | Noord-Holland | 731200   |
| 6  | Rotterdam   | NLD        | Zuid-Holland | 593321    |
| 7  | Haag        | NLD        | Zuid-Holland | 440900    |
| 8  | Utrecht     | NLD        | Utrecht    | 234323    |
| 9  | Eindhoven   | NLD        | Noord-Brabant | 201843   |
| 10 | Tilburg     | NLD        | Noord-Brabant | 193238   |
+----+-----+-----+-----+
10 rows in set (0.00 sec)
```

Some Handy Cleanup Commands

If you are not especially familiar with MySQL, the following commands may be helpful to clean up the examples. To remove the table in MySQL, enter the following command:

[Click here to view code image](#)

```
mysql> drop table 'CityExportStaging';
```

To remove the data in a table, enter this command:

[Click here to view code image](#)

```
mysql> delete from CityExportStaging;
```

To clean up imported files, enter this command:

[Click here to view code image](#)

```
$ hdfs dfs -rm -r -skipTrash sqoop-mysql-import/{country,city, canada-city}
```

USING APACHE FLUME TO ACQUIRE DATA STREAMS

Apache Flume is an independent agent designed to collect, transport, and store data into HDFS. Often data transport involves a number of Flume agents that may traverse a series of machines and locations. Flume is often used for log files, social media-generated data, email messages, and just about any continuous data source.

As shown in Figure 7.3, a Flume agent is composed of three components.

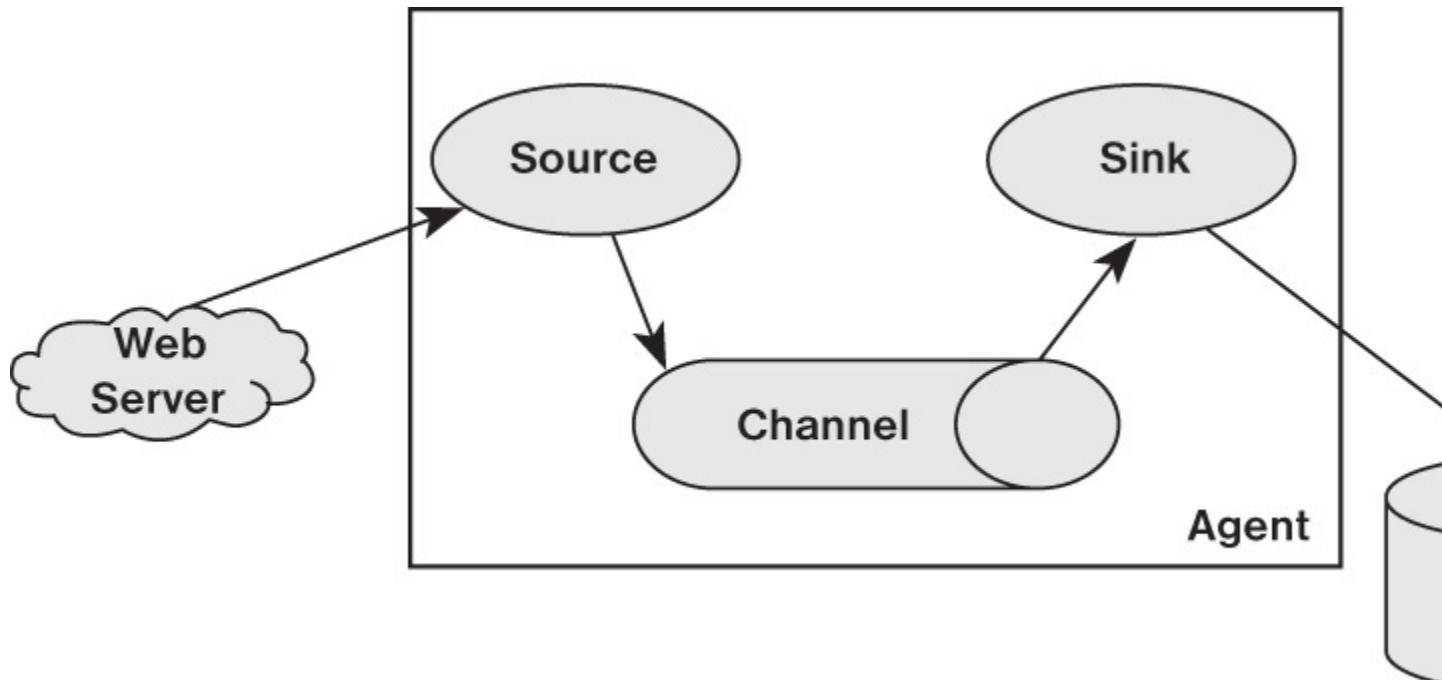


Figure 7.3 Flume agent with source, channel, and sink (Adapted from Apache Flume Documentation)

- **Source.** The source component receives data and sends it to a channel. It can send the data to more than one channel. The input data can be from a real-time source (e.g., weblog) or another Flume agent.
- **Channel.** A channel is a data queue that forwards the source data to the sink destination. It can be thought of as a buffer that manages input (source) and output (sink) flow rates.
- **Sink.** The sink delivers data to destination such as HDFS, a local file, or another Flume agent. A Flume agent must have all three of these components defined. A Flume agent can have several sources, channels, and sinks. Sources can write to multiple channels, but a sink can take data from only a single channel. Data written to a channel remain in the channel until a sink removes the data. By default, the data in a channel are kept in memory but may be optionally stored on disk to prevent data loss in the event of a network failure.

As shown in Figure 7.4, Sqoop agents may be placed in a pipeline, possibly to traverse several machines or domains. This configuration is normally used when data are collected on one machine (e.g., a web server) and sent to another machine that has access to HDFS.

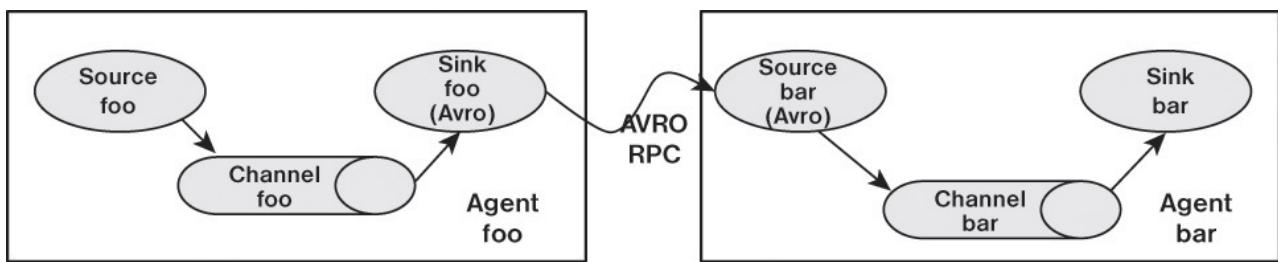


Figure 7.4 Pipeline created by connecting Flume agents (Adapted from Apache Flume Sqoop Documentation)

In a Flume pipeline, the sink from one agent is connected to the source of another. The data transfer format normally used by Flume, which is called Apache Avro, provides several useful features. First, Avro is a data serialization/deserialization system that uses a compact binary format. The schema is sent as part of the data exchange and is defined using JSON (JavaScript Object Notation). Avro also uses remote procedure calls (RPCs) to send data. That is, an Avro sink will contact an Avro source to send data.

Another useful Flume configuration is shown in [Figure 7.5](#). In this configuration, Flume is used to consolidate several data sources before committing them to HDFS.

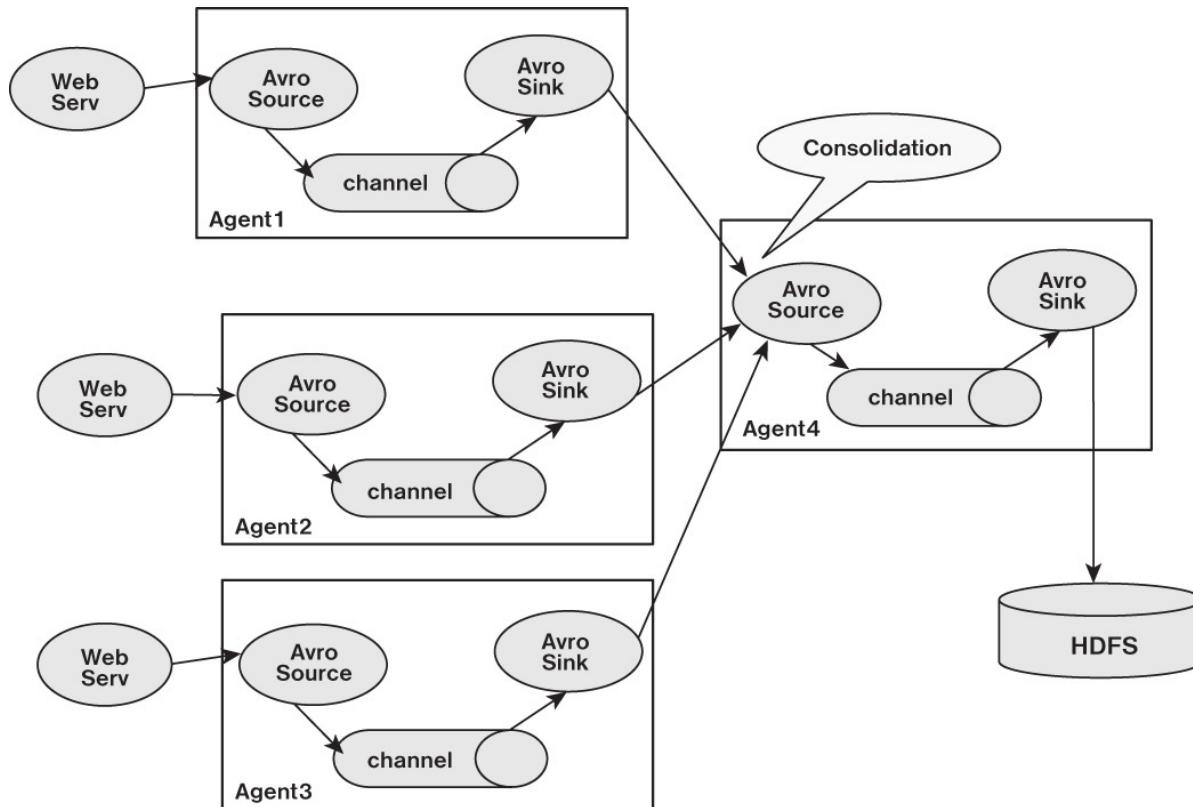


Figure 7.5 A Flume consolidation network (Adapted from Apache Flume Documentation)

There are many possible ways to construct Flume transport networks. In addition, other Flume features not described in depth here include plug-ins and interceptors that can enhance Flume pipelines. For more information and example configurations, see the Flume User Guide at <https://flume.apache.org/FlumeUserGuide.html>.

Flume Example Walk-Through

Follow these steps to walk through a Flume example.

Step 1: Download and Install Apache Flume

For this example, the following software environment is assumed. Other environments should work in a similar fashion.

- OS: Linux

- Platform: RHEL 6.6
- Hortonworks HDP 2.2 with Hadoop version: 2.6
- Flume version: 1.5.2

Flume can be installed by hand if you are using the pseudo-distributed installation from [Chapter 2](#). Consult the installation instructions on the Flume website: <https://flume.apache.org>. Flume is also installed as part of the Hortonworks HDP Sandbox. If Flume is not installed and you are using the Hortonworks HDP repository, you can add Flume with the following command:

[**Click here to view code image**](#)

```
# yum install flume flume-agent
```

In addition, for the simple example, `telnet` will be needed:

```
# yum install telnet
```

The following examples will also require some configuration files. See [Appendix A](#) for download instructions.

Step 2: Simple Test

A simple test of Flume can be done on a single machine. To start the Flume agent, enter the `flume-ng` command shown here. This command uses the `simple-example.conf` file to configure the agent.

[**Click here to view code image**](#)

```
$ flume-ng agent --conf conf --conf-file simple-example.conf --name simple_agent -Dflume.root.logger=INFO,console
```

In another terminal window, use `telnet` to contact the agent:

[**Click here to view code image**](#)

```
$ telnet localhost 4444
Trying ::1...
telnet: connect to address ::1: Connection refused
Trying 127.0.0.1...
Connected to localhost.
Escape character is '^].
testing 1 2 3
OK
```

If Flume is working correctly, the window where the Flume agent was started will show the testing message entered in the telnet window:

[**Click here to view code image**](#)

```
14/08/14 16:20:58 INFO sink.LoggerSink: Event: { headers:{} body: 74 65 73 74 69  
6E 67 20 20 31 20 32 20 33 0D  testing 1 2 3. }
```

Step 3: Weblog Example

In this example, a record from the weblogs from the local machine (Ambari output) will be placed into HDFS using Flume. This example is easily modified to use other weblogs from different machines. Two files are needed to configure Flume. (See the sidebar and [Appendix A](#) for file downloading instructions.)

- `web-server-target-agent.conf`—the target Flume agent that writes the data to HDFS
 - `web-server-source-agent.conf`—the source Flume agent that captures the weblog data
- The weblog is also mirrored on the local file system by the agent that writes to HDFS. To run the example, create the directory as `root`:

[Click here to view code image](#)

```
# mkdir /var/log/flume-hdfs  
# chown hdfs:hadoop /var/log/flume-hdfs/
```

Next, as user `hdfs`, make a Flume data directory in HDFS:

[Click here to view code image](#)

```
$ hdfs dfs -mkdir /user/hdfs/flume-channel/
```

Now that you have created the data directories, you can start the Flume target agent (execute as user `hdfs`):

[Click here to view code image](#)

```
$ flume-ng agent -c conf -f web-server-target-agent.conf -n collector
```

This agent writes the data into HDFS and should be started before the source agent. (The source reads the weblogs.) This configuration enables automatic use of the Flume agent.

The `/etc/flume/conf/{flume.conf, flume-env.sh.template}` files need to be configured for this purpose. For this example, the `/etc/flume/conf/flume.conf` file can be the same as the `web-server-target.conf` file (modified for your environment).

Note

With the HDP distribution, Flume can be started as a service when the system boots (e.g., `service start flume`).

In this example, the source agent is started as `root`, which will start to feed the weblog data to the target agent. Alternatively, the source agent can be on another machine if desired.

[Click here to view code image](#)

```
# flume-ng agent -c conf -f web-server-source-agent.conf -n source_agent
```

To see if Flume is working correctly, check the local log by using the `tail` command. Also confirm that the `flume-ng` agents are not reporting any errors (the file name will vary).

[**Click here to view code image**](#)

```
$ tail -f /var/log/flume-hdfs/1430164482581-1
```

The contents of the local log under `flume-hdfs` should be identical to that written into HDFS. You can inspect this file by using the `hdfs -tail` command (the file name will vary). Note that while running Flume, the most recent file in HDFS may have the extension `.tmp` appended to it. The `.tmp` indicates that the file is still being written by Flume. The target agent can be configured to write the file (and start another `.tmp` file) by setting some or all of the `rollCount`, `rollSize`, `rollInterval`, `idleTimeout`, and `batchSize` options in the configuration file.

[**Click here to view code image**](#)

```
$ hdfs dfs -tail flume-channel/apache_access_combined/150427/FlumeData.1430164801381
```

Both files should contain the same data. For instance, the preceding example had the following data in both files:

[**Click here to view code image**](#)

```
10.0.0.1 -- [27/Apr/2015:16:04:21 -0400] "GET /ambarinagios/nagios/nagios_alerts.php?q1=alerts&alert_type=all HTTP/1.1" 200 30801 "-" "Java/1.7.0_65"  
10.0.0.1 -- [27/Apr/2015:16:04:25 -0400] "POST /cgi-bin/rrd.py HTTP/1.1" 200 784  
"-" "Java/1.7.0_65"  
10.0.0.1 -- [27/Apr/2015:16:04:25 -0400] "POST /cgi-bin/rrd.py HTTP/1.1" 200 508  
"-" "Java/1.7.0_65"
```

You can modify both the target and source files to suit your system.

Flume Configuration Files

A compete explanation of Flume configuration is beyond the scope of this chapter. The Flume website has additional information on Flume configuration: <http://flume.apache.org/FlumeUserGuide.html#configuration>. The configurations used previously also have links to help explain the settings. Some of the important settings used in the preceding example follow.

In `web-server-source-agent.conf`, the following lines set the source. Note that the weblog is acquired by using the `tail` command to record the log file.

[**Click here to view code image**](#)

```
source_agent.sources = apache_server  
source_agent.sources.apache_server.type = exec  
source_agent.sources.apache_server.command = tail -f /etc/httpd/logs/access_log
```

Further down in the file, the sink is defined. The `source_agent.sinks.avro_sink.hostname` is used to assign the Flume node that will write to HDFS. The port number is also set in the target configuration file.

[**Click here to view code image**](#)

```
source_agent.sinks = avro_sink
source_agent.sinks.avro_sink.type = avro
source_agent.sinks.avro_sink.channel = memoryChannel
source_agent.sinks.avro_sink.hostname = 192.168.93.24
source_agent.sinks.avro_sink.port = 4545
```

The HDFS settings are placed in the `web-server-target-agent.conf` file. Note the path that was used in the previous example and the data specification.

[**Click here to view code image**](#)

```
collector.sinks.HadoopOut.type = hdfs
collector.sinks.HadoopOut.channel = mc2
collector.sinks.HadoopOut.hdfs.path = /user/hdfs/flume-channel/%{log_type}/
%y%m%d
collector.sinks.HadoopOut.hdfs.fileType = DataStream
```

The target file also defines the port and two channels (mc1 and mc2). One of these channels writes the data to the local file system, and the other writes to HDFS. The relevant lines are shown here:

[**Click here to view code image**](#)

```
collector.sources.AvroIn.port = 4545
collector.sources.AvroIn.channels = mc1 mc2
collector.sinks.LocalOut.sink.directory = /var/log/flume-hdfs
collector.sinks.LocalOut.channel = mc1
```

The HDFS file rollover counts create a new file when a threshold is exceeded. In this example, that threshold is defined to allow any file size and write a new file after 10,000 events or 600 seconds.

[**Click here to view code image**](#)

```
collector.sinks.HadoopOut.hdfs.rollSize = 0
collector.sinks.HadoopOut.hdfs.rollCount = 10000
collector.sinks.HadoopOut.hdfs.rollInterval = 600
```

A full discussion of Flume can be found on the Flume website.

MANAGE HADOOP WORKFLOWS WITH APACHE OOZIE

Oozie is a workflow director system designed to run and manage multiple related Apache Hadoop jobs. For instance, complete data input and analysis may require several discrete Hadoop jobs to be run as a workflow in which the output of one job serves as the input for a successive job. Oozie is designed to construct and manage these workflows. Oozie is not a substitute for the YARN scheduler. That is, YARN manages resources for individual Hadoop jobs, and Oozie provides a way to connect and control Hadoop jobs on the cluster.

Oozie workflow jobs are represented as directed acyclic graphs (DAGs) of actions. (DAGs are basically graphs that cannot have directed loops.) Three types of Oozie jobs are permitted:

- Workflow—a specified sequence of Hadoop jobs with outcome-based decision points and control dependency. Progress from one action to another cannot happen until the first action is complete.

- Coordinator—a scheduled workflow job that can run at various time intervals or when data become available.

- Bundle—a higher-level Oozie abstraction that will batch a set of coordinator jobs.

Oozie is integrated with the rest of the Hadoop stack, supporting several types of Hadoop jobs out of the box (e.g., Java MapReduce, Streaming MapReduce, Pig, Hive, and Sqoop) as well as system-specific jobs (e.g., Java programs and shell scripts). Oozie also provides a CLI and a web UI for monitoring jobs.

Figure 7.6 depicts a simple Oozie workflow. In this case, Oozie runs a basic MapReduce operation. If the application was successful, the job ends; if an error occurred, the job is killed.

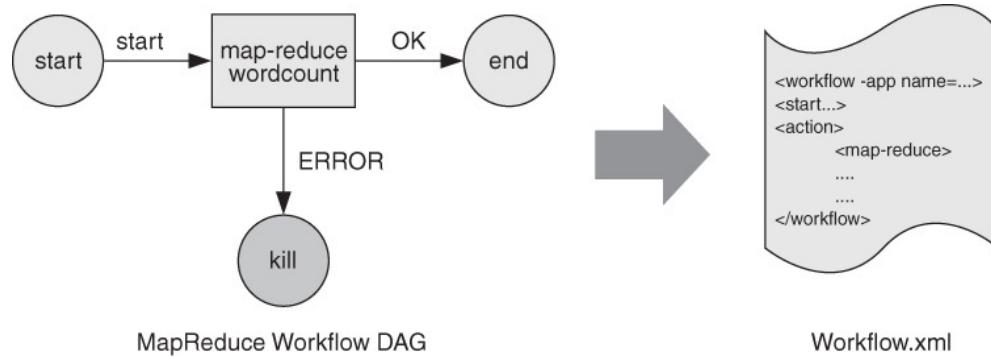


Figure 7.6 A simple Oozie DAG workflow (Adapted from Apache Oozie Documentation)

Oozie workflow definitions are written in hPDL (an XML Process Definition Language). Such workflows contain several types of nodes:

- **Control flow nodes** define the beginning and the end of a workflow. They include start, end, and optional fail nodes.
- **Action nodes** are where the actual processing tasks are defined. When an action node finishes, the remote systems notify Oozie and the next node in the workflow is executed. Action nodes can also include HDFS commands.
- **Fork/join nodes** enable parallel execution of tasks in the workflow. The fork node enables two or more tasks to run at the same time. A join node represents a rendezvous point that must wait until all forked tasks complete.
- **Control flow nodes** enable decisions to be made about the previous task. Control decisions are based on the results of the previous action (e.g., file size or file existence). Decision nodes are essentially switch-case statements that use JSP EL (Java Server Pages—Expression Language) that evaluate to either true or false.

Figure 7.7 depicts a more complex workflow that uses all of these node types. More information on Oozie can be found at <http://oozie.apache.org/docs/4.1.0/index.html>.

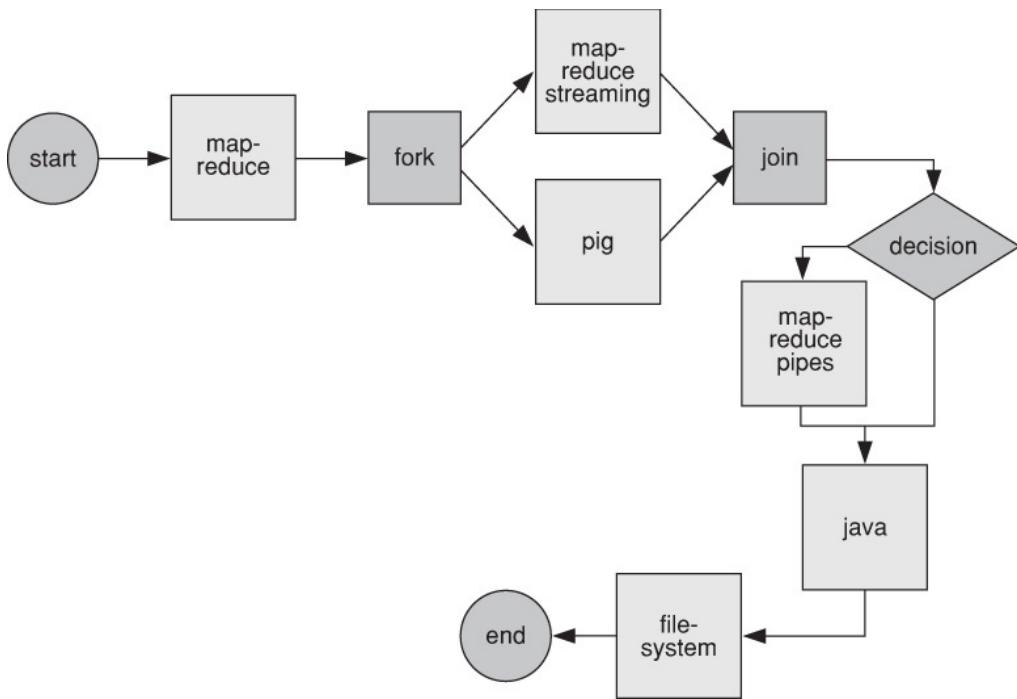


Figure 7.7 A more complex Oozie DAG workflow (Adapted from Apache Oozie Documentation)

Oozie Example Walk-Through

For this example, the following software environment is assumed. Other environments should work in a similar fashion.

- OS: Linux
- Platform: CentOS 6.6
- Hortonworks HDP 2.2 with Hadoop version: 2.6
- Oozie version: 4.1.0

If you are using the pseudo-distributed installation from [Chapter 2](#) or want to install Oozie by hand, see the installation instructions on the Ozzie website: <http://oozie.apache.org>. Oozie is also installed as part of the Hortonworks HDP Sandbox.

Step 1: Download Oozie Examples

The Oozie examples used in this section can be found on the book website (see [Appendix A](#)). They are also available as part of the `oozie-client.noarch` RPM in the Hortonworks HDP 2.x packages. For HDP 2.1, the following command can be used to extract the files into the working directory used for the demo:

[Click here to view code image](#)

```
$ tar xvzf /usr/share/doc/oozie-4.0.0.2.1.2.1/oozie-examples.tar.gz
```

For HDP 2.2, the following command will extract the files:

[Click here to view code image](#)

```
$ tar xvzf /usr/hdp/2.2.4.2-2/oozie/doc/oozie-examples.tar.gz
```

Once extracted, rename the examples directory to `oozie-examples` so that you will not confuse it with the other examples directories.

```
$ mv examples oozie-examples
```

The examples must also be placed in HDFS. Enter the following command to move the example files into HDFS:

[**Click here to view code image**](#)

```
$ hdfs dfs -put oozie-examples/ oozie-examples
```

The Oozie shared library must be installed in HDFS. If you are using the Ambari installation of HDP 2.x, this library is already found in HDFS: `/user/oozie/share/lib`.

Note

In HDP 2.2+, some additional version-tagged directories may appear below this path. If you installed and built Oozie by hand, then make sure `/user/oozie` exists in HDFS and put the `oozie-sharelib` files in this directory as user `oozie` and group `hadoop`.

The example applications are found under the `oozie-examples/app` directory, one directory per example. Each directory contains at least `workflow.xml` and `job.properties` files. Other files needed for each example are also in its directory. The inputs for all examples are in the `oozie-examples/input-data` directory. The examples will create output under the `examples/output-data` directory in HDFS.

Step 2: Run the Simple MapReduce Example

Move to the simple MapReduce example directory:

[**Click here to view code image**](#)

```
$ cd oozie-examples/apps/map-reduce/
```

This directory contains two files and a `lib` directory. The files are:

- The `job.properties` file defines parameters (e.g., path names, ports) for a job. This file may change per job.
- The `workflow.xml` file provides the actual workflow for the job. In this case, it is a simple MapReduce (pass/fail). This file usually stays the same between jobs.

The `job.properties` file included in the examples requires a few edits to work properly. Using a text editor, change the following lines by adding the host name of the NameNode and ResourceManager (indicated by `jobTracker` in the file).

[**Click here to view code image**](#)

```
nameNode=hdfs://localhost:8020  
jobTracker=localhost:8032
```

to the following (note the port change for jobTracker):

Click here to view code image

```
nameNode=dfs://_HOSTNAME_:8020  
jobTracker=_HOSTNAME_:8050
```

For example, for the cluster created with Ambari in [Chapter 2](#), the lines were changed to

```
nameNode=dfs://limulus:8020  
jobTracker=limulus:8050
```

The examplesRoot variable must also be changed to oozie-examples, reflecting the change made previously:

```
examplesRoot=oozie-examples
```

These changes must be done for all the job.properties files in the Oozie examples that you choose to run.

The DAG for the simple MapReduce example is shown in [Figure 7.6](#). The workflow.xml file describes these simple steps and has the following workflow nodes:

```
<start to="mr-node"/>  
<action name="mr-node">  
<kill name="fail">  
<end name="end"/>
```

A complete description of Oozie workflows is beyond the scope of this chapter. However, in the simple case described here, basic aspects of the file can be highlighted. First, under the `<action name="mr-node">` tag, the MapReduce process is set with a `<map-reduce>` tag. As part of this description, the `<prepare>` and `<configuration>` tags set up the job. Note that the `mapred.{mapper, reducer}.class` refer to the local `lib` directory. As shown in [Figure 7.6](#), this simple workflow runs an example MapReduce job and prints an error message if it fails. To run the Oozie MapReduce example job from the `oozie-examples/apps/map-reduce` directory, enter the following line:

Click here to view code image

```
$ oozie job -run -oozie http://limulus:11000/oozie -config job.properties
```

When Oozie accepts the job, a job ID will be printed:

Click here to view code image

```
job: 0000001-150424174853048-oozie-oozi-W
```

You will need to change the “limulus” host name to match the name of the node running your Oozie server. The job ID can be used to track and control job progress.

The “Oozie Is Not Allowed to Impersonate Oozie” Error

When trying to run Oozie, you may get the puzzling error:

[**Click here to view code image**](#)

oozie is not allowed to impersonate oozie

If you receive this message, make sure the following is defined in the core-site.xml file:

[**Click here to view code image**](#)

```
<property>
  <name>hadoop.proxyuser.oozie.hosts</name>
    <value>*</value>
</property>

<property>
  <name>hadoop.proxyuser.oozie.groups</name>
    <value>*</value>
</property>
```

If you are using Ambari, make this change (or add the lines) in the Services/HDFS/Config window and restart Hadoop. Otherwise, make the change by hand and restart all the Hadoop daemons.

This setting is required because Oozie needs to impersonate other users to run jobs. The group property can be set to a specific user group or to a wild card. This setting allows the account that runs the Oozie server to run as part of the user’s group.

To avoid having to provide the `-oozie` option with the Oozie URL every time you run the `oozie` command, set the `OOZIE_URL` environment variable as follows (using your Oozie server host name in place of “limulus”):

[**Click here to view code image**](#)

```
$ export OOZIE_URL="http://limulus:11000/oozie"
```

You can now run all subsequent Oozie commands without specifying the `-oozie` URL option. For instance, using the job ID, you can learn about a particular job’s progress by issuing the following command:

[**Click here to view code image**](#)

```
$ oozie job -info 0000001-150424174853048-oozie-oozi-W
```

The resulting output (line length compressed) is shown in the following listing. Because this job is just a simple test, it may be complete by the time you issue the `-info` command. If it is not complete, its progress will be indicated in the listing.

[**Click here to view code image**](#)

```
Job ID : 0000001-150424174853048-oozie-oozi-W
```

```
-----
```

Workflow Name : map-reduce-wf

App Path : hdfs://limulus:8020/user/hdfs/examples/apps/map-reduce

Status : SUCCEEDED
Run : 0
User : hdfs
Group : -
Created : 2015-04-29 20:52 GMT
Started : 2015-04-29 20:52 GMT
Last Modified : 2015-04-29 20:53 GMT
Ended : 2015-04-29 20:53 GMT
CoordAction ID: -

Actions

ID	Status	Ext ID	Ext Status	Err Code
0000001-150424174853048-oozie -oozi-W@:start:	OK	-	OK	-
0000001-150424174853048-oozie -oozi-W@mr-node	OK	job_1429912013449_0006	SUCCEEDED	-
0000001-150424174853048-oozie -oozi-W@end	OK	-	OK	-

The various steps shown in the output can be related directly to the `workflow.xml` mentioned previously. Note that the MapReduce job number is provided. This job will also be listed in the ResourceManager web user interface. The application output is located in HDFS under the `oozie-examples/output-data/map-reduce` directory.

Step 3: Run the Oozie Demo Application

A more sophisticated example can be found in the demo directory (`oozie-examples/apps/demo`). This workflow includes MapReduce, Pig, and file system tasks as well as fork, join, decision, action, start, stop, kill, and end nodes.

Move to the demo directory and edit the `job.properties` file as described previously. Entering the following command runs the workflow (assuming the `OOZIE_URL` environment variable has been set):

[Click here to view code image](#)

```
$ oozie job -run -config job.properties
```

You can track the job using either the Oozie command-line interface or the Oozie web console. To start the web console from within Ambari, click on the Oozie service, and then click on the Quick Links pull-down menu and select Oozie Web UI. Alternatively, you can start the Oozie web UI by connecting to the Oozie server directly. For example, the following command will bring up the Oozie UI (use your Oozie server host name in place of “limulus”):

[**Click here to view code image**](#)

```
$ firefox http://limulus:11000/oozie/
```

Figure 7.8 shows the main Oozie console window. Note that a link to Oozie documentation is available directly from this window.

The screenshot shows a Mozilla Firefox browser window titled "Oozie Web Console - Mozilla Firefox". The address bar displays "limulus:11000/oozie/". The page content is the "Oozie Web Console" interface. At the top, there is a navigation bar with tabs: Workflow Jobs, Coordinator Jobs, Bundle Jobs, System Info, Instrumentation, and Settings. Below the tabs, there is a sub-navigation bar with buttons: All Jobs (which is selected), Active Jobs, Done Jobs, and Custom Filter. To the right of this bar, it says "Server version [4.1.0.2.2.4.2-2]". The main area is a table listing four workflow jobs:

Job Id	Name	Status	Run	User	Group	Created	Started	Last Modified
1 0000003-150424174853048-oozie-oozi-W	demo-wf	RUNNING	0	hdfs		Wed, 29 Apr 2015 21:05:32 GMT	Wed, 29 Apr 2015 21:05:32 GMT	Wed, 29 Apr 2015 21:05:37 GMT
2 0000002-150424174853048-oozie-oozi-W	map-reduce-wf	SUCCEE...	0	hdfs		Wed, 29 Apr 2015 21:01:23 GMT	Wed, 29 Apr 2015 21:01:23 GMT	Wed, 29 Apr 2015 21:01:49 GMT
3 0000001-150424174853048-oozie-oozi-W	map-reduce-wf	SUCCEE...	0	hdfs		Wed, 29 Apr 2015 20:52:41 GMT	Wed, 29 Apr 2015 20:52:41 GMT	Wed, 29 Apr 2015 20:53:07 GMT
4 0000000-150424174853048-oozie-oozi-W	map-reduce-wf	SUCCEE...	0	ambari-qa		Fri, 24 Apr 2015 21:50:13 GMT	Fri, 24 Apr 2015 21:50:13 GMT	Fri, 24 Apr 2015 21:50:43 GMT

Figure 7.8 Oozie main console window

Workflow jobs are listed in tabular form, with the most recent job appearing first. If you click on a workflow, the Job Info window in Figure 7.9 will be displayed. The job progression results, similar to those printed by the Oozie command line, are shown in the Actions window at the bottom.

The screenshot shows the Oozie Web Console interface in Mozilla Firefox. The title bar reads "Oozie Web Console - Mozilla Firefox". The address bar shows the URL "limulus:11000/oozie/". The main content area displays a job information window for a workflow named "demo-wf". The "Job Info" tab is selected, showing details such as Job Id, Name, App Path, Run, Status, User, Group, Parent Coord, Create Time, Start Time, Last Modified, and End Time. Below this is a table titled "Actions" listing 10 workflow actions with columns for Action Id, Name, Type, Status, Transition, StartTime, and EndTime.

Action Id	Name	Type	Status	Transition	StartTime	EndTime
1	0000003-150424174853048-oozie-oozi-W@:start	:start:	OK	cleanup-node	Wed, 29 Apr 2015 21:05:32 GMT	Wed, 29 Apr 2015 21:05:32 GMT
2	0000003-150424174853048-oozie-oozi-W@cleanup-node	cleanup-node	fs	fork-node	Wed, 29 Apr 2015 21:05:32 GMT	Wed, 29 Apr 2015 21:05:33 GMT
3	0000003-150424174853048-oozie-oozi-W@fork-node	fork-node	:FORK:	*	Wed, 29 Apr 2015 21:05:33 GMT	Wed, 29 Apr 2015 21:05:33 GMT
4	0000003-150424174853048-oozie-oozi-W@pig-node	pig-node	pig	join-node	Wed, 29 Apr 2015 21:05:33 GMT	Wed, 29 Apr 2015 21:06:05 GMT
5	0000003-150424174853048-oozie-oozi-W@streaming-no...	streaming-n...	map-reduce	join-node	Wed, 29 Apr 2015 21:05:36 GMT	Wed, 29 Apr 2015 21:06:01 GMT
6	0000003-150424174853048-oozie-oozi-W@join-node	join-node	:JOIN:	mr-node	Wed, 29 Apr 2015 21:06:06 GMT	Wed, 29 Apr 2015 21:06:06 GMT
7	0000003-150424174853048-oozie-oozi-W@mr-node	mr-node	map-reduce	decision-node	Wed, 29 Apr 2015 21:06:06 GMT	Wed, 29 Apr 2015 21:06:30 GMT
8	0000003-150424174853048-oozie-oozi-W@decision-node	decision-node	switch	hdfs-node	Wed, 29 Apr 2015 21:06:30 GMT	Wed, 29 Apr 2015 21:06:30 GMT
9	0000003-150424174853048-oozie-oozi-W@hdfs-node	hdfs-node	fs	end	Wed, 29 Apr 2015 21:06:31 GMT	Wed, 29 Apr 2015 21:06:31 GMT
10	0000003-150424174853048-oozie-oozi-W@end	end	:END:		Wed, 29 Apr 2015 21:06:31 GMT	Wed, 29 Apr 2015 21:06:31 GMT

Figure 7.9 Oozie workflow information window

Other aspects of the job can be examined by clicking the other tabs in the window. The last tab actually provides a graphical representation of the workflow DAG. If the job is not complete, it will highlight the steps that have been completed thus far. The DAG for the demo example is shown in [Figure 7.10](#). The actual image was split to fit better on the page. As with the previous example, comparing this information to `workflow.xml` file can provide further insights into how Oozie operates.

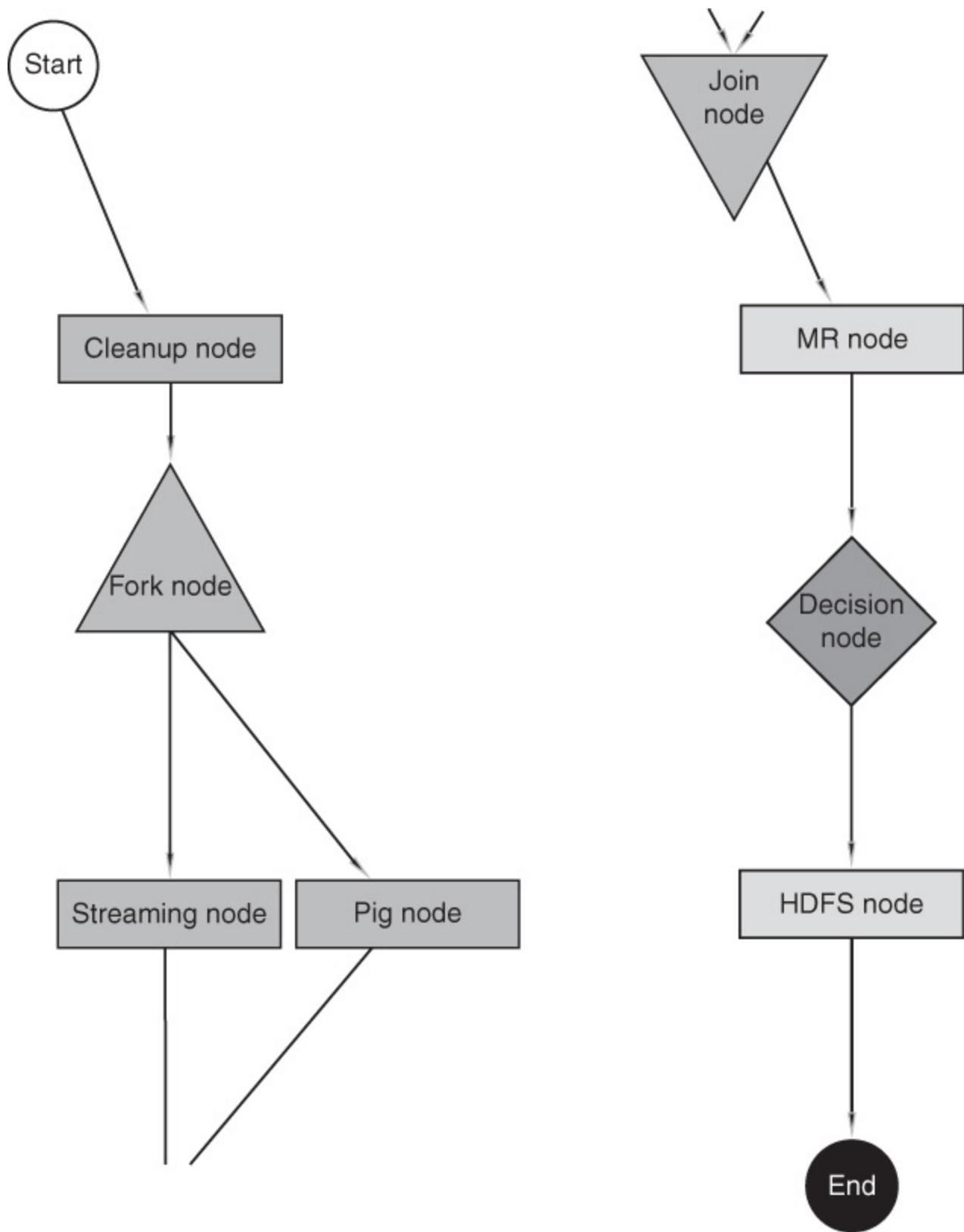


Figure 7.10 Oozie-generated workflow DAG for the demo example, as it appears on the screen

A Short Summary of Oozie Job Commands

The following summary lists some of the more commonly encountered Oozie commands. See the latest documentation at <http://oozie.apache.org> for more information. (Note that the examples here assume `OOZIE_URL` is defined.)

- Run a workflow job (returns `_OOZIE_JOB_ID_`):

[Click here to view code image](#)

```
$ oozie job -run -config JOB_PROPERTIES
```

- Submit a workflow job (returns `_OOZIE_JOB_ID_` but does not start):

[Click here to view code image](#)

```
$ oozie job -submit -config JOB_PROPERTIES
```

- Start a submitted job:

[Click here to view code image](#)

```
$ oozie job -start _OOZIE_JOB_ID_
```

- Check a job's status:

[Click here to view code image](#)

```
$ oozie job -info _OOZIE_JOB_ID_
```

- Suspend a workflow:

[Click here to view code image](#)

```
$ oozie job -suspend _OOZIE_JOB_ID_
```

- Resume a workflow:

[Click here to view code image](#)

```
$ oozie job -resume _OOZIE_JOB_ID_
```

- Rerun a workflow:

[Click here to view code image](#)

```
$ oozie job -rerun _OOZIE_JOB_ID_ -config JOB_PROPERTIES
```

- Kill a job:

[Click here to view code image](#)

```
$ oozie job -kill _OOZIE_JOB_ID_
```

- View server logs:

[Click here to view code image](#)

```
$ oozie job -logs _OOZIE_JOB_ID_
```

Full logs are available at `/var/log/oozie` on the Oozie server.

USING APACHE HBASE

Apache HBase is an open source, distributed, versioned, nonrelational database modeled after Google's Bigtable (<http://research.google.com/archive/bigtable.html>). Like Bigtable, HBase leverages the distributed data storage provided by the underlying distributed file systems spread across commodity servers. Apache HBase provides Bigtable-like capabilities on top of Hadoop and HDFS. Some of the more important features include the following capabilities:

- Linear and modular scalability
- Strictly consistent reads and writes
- Automatic and configurable sharding of tables
- Automatic failover support between RegionServers
- Convenient base classes for backing Hadoop MapReduce jobs with Apache HBase tables
- Easy-to-use Java API for client access

HBase Data Model Overview

A table in HBase is similar to other databases, having rows and columns. Columns in HBase are grouped into column families, all with the same prefix. For example, consider a table of daily stock prices. There may be a column family called "price" that has four members—`price:open`, `price:close`, `price:low`, and `price:high`. A column does not need to be a family. For instance, the stock table may have a column named "volume" indicating how many shares were traded. All column family members are stored together in the physical file system.

Specific HBase cell values are identified by a row key, column (column family and column), and version (timestamp). It is possible to have many versions of data within an HBase cell. A version is specified as a timestamp and is created each time data are written to a cell. Almost anything can serve as a row key, from strings to binary representations of longs to serialized data structures. Rows are lexicographically sorted with the lowest order appearing first in a table. The empty byte array denotes both the start and the end of a table's namespace. All table accesses are via the table row key, which is considered its primary key. More information on HBase can be found on the HBase website: <http://hbase.apache.org>.

HBase Example Walk-Through

For this example, the following software environment is assumed. Other environments should work in a similar fashion.

- OS: Linux
- Platform: CentOS 6.6
- Hortonworks HDP 2.2 with Hadoop version: 2.6
- HBase version: 0.98.4

If you are using the pseudo-distributed installation from [Chapter 2](#) or want to install HBase by hand, see the installation instructions on the HBase website: <http://hbase.apache.org>. HBase is also installed as part of the Hortonworks HDP Sandbox.

The following example illustrates a small subset of HBase commands. Consult the HBase website for more background. HBase provides a shell for interactive use. To enter the shell, type the following as a user:

```
$ hbase shell  
hbase(main):001:0>
```

To exit the shell, type `exit`.

Various commands can be conveniently entered from the shell prompt. For instance, the `status` command provides the system status:

[**Click here to view code image**](#)

```
hbase(main):001:0> status  
4 servers, 0 dead, 1.0000 average load
```

Additional arguments can be added to the `status` command, including '`simple`', '`summary`', or '`detailed`'. The single quotes are needed for proper operation. For example, the following command will provide simple status information for the four HBase servers (actual server statistics have been removed for clarity):

[**Click here to view code image**](#)

```
hbase(main):002:0> status 'simple'  
4 live servers  
n1:60020 1429912048329  
...  
n2:60020 1429912040653  
...  
limulus:60020 1429912041396  
...  
n0:60020 1429912042885  
...  
0 dead servers  
Aggregate load: 0, regions: 4
```

Other basic commands, such as `version` or `whoami`, can be entered directly at the `hbase (main)` prompt. In the example that follows, we will use a small set of daily stock price data for Apple computer. The data have the following form:

Date	Open	High	Low	Close	Volume
6-May-15	126.56	126.75	123.36	125.01	71820387

The data can be downloaded from Google using the following command. Note that other stock prices are available by changing the `NASDAQ:AAPL` argument to any other valid exchange and stock name (e.g., `NYSE: IBM`).

[Click here to view code image](#)

```
$ wget -O Apple-stock.csv  
http://www.google.com/finance/historical?q=NASDAQ:AAPL&authuser=0&output=csv
```

The Apple stock price database is in comma-separated format (csv) and will be used to illustrate some basic operations in the HBase shell.

Create the Database

The next step is to create the database in HBase using the following command:

[Click here to view code image](#)

```
hbase(main):006:0> create 'apple', 'price' , 'volume'  
0 row(s) in 0.8150 seconds
```

In this case, the table name is `apple`, and two columns are defined. The date will be used as the row key. The `price` column is a family of four values (`open`, `close`, `low`, `high`).

The `put` command is used to add data to the database from within the shell. For instance, the preceding data can be entered by using the following commands:

[Click here to view code image](#)

```
put 'apple','6-May-15','price:open','126.56'  
put 'apple','6-May-15','price:high','126.75'  
put 'apple','6-May-15','price:low','123.36'  
put 'apple','6-May-15','price:close','125.01'  
put 'apple','6-May-15','volume','71820387'
```

Note that these commands can be copied and pasted into HBase shell and are available from the book download files (see [Appendix A](#)). The shell also keeps a history for the session, and previous commands can be retrieved and edited for resubmission.

Inspect the Database

The entire database can be listed using the `scan` command. Be careful when using this command with large databases. This example is for one row.

[Click here to view code image](#)

```
scan 'apple'  
hbase(main):006:0> scan 'apple'  
ROW          COLUMN+CELL  
6-May-15    column=price:close, timestamp=1430955128359, value=125.01  
6-May-15    column=price:high, timestamp=1430955126024, value=126.75
```

```
6-May-15    column=price:low, timestamp=1430955126053, value=123.36
6-May-15    column=price:open, timestamp=1430955125977, value=126.56
6-May-15    column=volume:, timestamp=1430955141440, value=71820387
```

Get a Row

You can use the row key to access an individual row. In the stock price database, the date is the row key.

[Click here to view code image](#)

```
hbase(main):008:0> get 'apple', '6-May-15'
COLUMN          CELL
price:close     timestamp=1430955128359, value=125.01
price:high      timestamp=1430955126024, value=126.75
price:low       timestamp=1430955126053, value=123.36
price:open      timestamp=1430955125977, value=126.56
volume:         timestamp=1430955141440, value=71820387
5 row(s) in 0.0130 seconds
```

Get Table Cells

A single cell can be accessed using the `get` command and the `COLUMN` option:

[Click here to view code image](#)

```
hbase(main):013:0> get 'apple', '5-May-15', {COLUMN => 'price:low'}
COLUMN          CELL
price:low       timestamp=1431020767444, value=125.78
1 row(s) in 0.0080 seconds
```

In a similar fashion, multiple columns can be accessed as follows:

[Click here to view code image](#)

```
hbase(main):012:0> get 'apple', '5-May-15', {COLUMN => ['price:low', 'price:high']}
COLUMN          CELL
price:high      timestamp=1431020767444, value=128.45
price:low       timestamp=1431020767444, value=125.78
2 row(s) in 0.0070 seconds
```

Delete a Cell

A specific cell can be deleted using the following command:

[Click here to view code image](#)

```
hbase(main):009:0> delete 'apple', '6-May-15', 'price:low'
```

If the row is inspected using `get`, the `price:low` cell is not listed.

[Click here to view code image](#)

```
hbase(main):010:0> get 'apple', '6-May-15'
COLUMN          CELL
price:close    timestamp=1430955128359, value=125.01
price:high     timestamp=1430955126024, value=126.75
price:open      timestamp=1430955125977, value=126.46
volume:        timestamp=1430955141440, value=71820387
4 row(s) in 0.0130 seconds
```

Delete a Row

You can delete an entire row by giving the `deleteall` command as follows:

[Click here to view code image](#)

```
hbase(main):009:0> deleteall 'apple', '6-May-15'
```

Remove a Table

To remove (drop) a table, you must first disable it. The following two commands remove the `apple` table from Hbase:

[Click here to view code image](#)

```
hbase(main):009:0> disable 'apple'
hbase(main):010:0> drop 'apple'
```

Scripting Input

Commands to the HBase shell can be placed in bash scripts for automated processing. For instance, the following can be placed in a bash script:

[Click here to view code image](#)

```
echo "put 'apple','6-May-15','price:open','126.56'" | hbase shell
```

The book software page includes a script (`input_to_hbase.sh`) that imports the `Apple-stock.csv` file into HBase using this method. It also removes the column titles in the first line. The script will load the entire file into HBase when you issue the following command:

[Click here to view code image](#)

```
$ input_to_hbase.sh Apple-stock.csv
```

While the script can be easily modified to accommodate other types of data, it is not recommended for production use because the upload is very inefficient and slow. Instead, this script is best used to experiment with small data files and different types of data.

Adding Data in Bulk

There are several ways to efficiently load bulk data into HBase. Covering all of these methods is beyond the scope of this chapter. Instead, we will focus on the `ImportTsv` utility, which loads data in tab-separated values (tsv) format into HBase. It has two distinct usage modes:

- Loading data from a tsv-format file in HDFS into HBase via the `put` command
- Preparing StoreFiles to be loaded via the `completebulkload` utility

The following example shows how to use `ImportTsv` for the first option, loading the tsv-format file using the `put` command. The second option works in a two-step fashion and can be explored by consulting <http://hbase.apache.org/book.html#importtsv>.

The first step is to convert the `Apple-stock.csv` file to tsv format. The following script, which is included in the book software, will remove the first line and do the conversion. In doing so, it creates a file named `Apple-stock.tsv`.

Click here to view code image

```
$ convert-to-tsv.sh Apple-stock.csv
```

Next, the new file is copied to HDFS as follows:

Click here to view code image

```
$ hdfs dfs -put Apple-stock.tsv /tmp
```

Finally, `ImportTsv` is run using the following command line. Note the column designation in the `-Dimporttsv.columns` option. In the example, the `HBASE_ROW_KEY` is set as the first column—that is, the date for the data.

Click here to view code image

```
$ hbase org.apache.hadoop.hbase.mapreduce.ImportTsv -  
Dimporttsv.columns=HBASE_ROW_KEY,price:open,price:high,price:low,price:close,volume  
apple /tmp/Apple-stock.tsv
```

The `ImportTsv` command will use MapReduce to load the data into HBase. To verify that the command works, drop and re-create the `apple` database, as described previously, before running the import command.

Apache HBase Web Interface

Like many of the Hadoop ecosystem tools, HBase has a web interface. To start the HBase console, shown in [Figure 7.11](#), from within Ambari, click on the HBase service, and then click on the Quick Links pull-down menu and select HBase Master UI. Alternatively, you can connect to the HBase master directly to start the HBase web UI. For example, the following command will bring up the HBase UI (use your HBase master server host name in place of “limulus”):

Click here to view code image

```
$ firefox http://limulus:60010/master-status
```

The screenshot shows the HBase master status interface in Mozilla Firefox. The title bar reads "Master: limulus - Mozilla Firefox". The address bar shows "limulus:60010/master-status". The top navigation bar includes links for Home, Table Details, Local Logs, Log Level, Debug Dump, Metrics Dump, and HBase Configuration. Below this, the main content area is titled "Master limulus".

Region Servers

Base Stats Memory Requests Storefiles Compactions

ServerName	Start time	Requests Per Second	Num. Regions
limulus,60020,1429912041396	Fri Apr 24 17:47:21 EDT 2015	0	1
n0,60020,1429912042885	Fri Apr 24 17:47:22 EDT 2015	0	2
n1,60020,1429912048329	Fri Apr 24 17:47:28 EDT 2015	0	0
n2,60020,1429912040653	Fri Apr 24 17:47:20 EDT 2015	0	1
Total:4		0	4

Backup Masters

ServerName	Port	Start Time

Total:0

Tables

User Tables Catalog Tables Snapshots

Figure 7.11 HBase web GUI

The HBase master status (including links to the region servers) is reported by the GUI in addition to statistics for the data tables in the database.

SUMMARY AND ADDITIONAL RESOURCES

This chapter introduced several essential Hadoop tools and provided at least one complete end-to-end example application of each tool. For each tool, the tool version and Hadoop and operating system distribution were provided, along with references and links to any example data.

The Apache Pig scripting tool enables you to quickly examine data and preprocess raw data for later Hadoop analysis. Apache Hive is an SQL-like interface for Hadoop. The Apache Sqoop application can be used to import and export RDBMS data into HDFS. Apache Flume is used to capture and transport weblog data.

To explore the concept of Hadoop workflows, this chapter also introduced the Oozie workflow management tool. Two examples of Oozie workflows were described, along with the Oozie web GUI.

Finally, the Apache HBase distributed database was introduced. Basic commands were covered, and an example demonstrated the bulk import of external stock market data into HBase.

Additional information and background on each of the tools can be obtained from the following resources.

- **Apache Pig scripting language**

- <http://pig.apache.org/>
- <http://pig.apache.org/docs/r0.14.0/start.html>

- **Apache Hive SQL-like query language**

- <https://hive.apache.org/>
- <https://cwiki.apache.org/confluence/display/Hive/GettingStarted>
- <http://grouplens.org/datasets/movielens> (data for example)

- **Apache Sqoop RDBMS import/export**

- <http://sqoop.apache.org>
- <http://dev.mysql.com/doc/world-setup/en/index.html> (data for example)

- **Apache Flume steaming data and transport utility**

- <https://flume.apache.org>
- <https://flume.apache.org/FlumeUserGuide.html>

- **Apache Oozie workflow manager**

- <http://oozie.apache.org>
- <http://oozie.apache.org/docs/4.0.0/index.html>

- **Apache HBase distributed database**

- <http://hbase.apache.org/book.html>
- <http://hbase.apache.org>
- <http://research.google.com/archive/bigtable.html> (Google Big Table paper)
- <http://www.google.com/finance/historical?q=NASDAQ:AAPL&authuser=0&output=csv> (data for example)

8. Hadoop YARN Applications

In This Chapter:

- The YARN Distributed-Shell is introduced as a non-MapReduce application.
- The Hadoop YARN application and operation structure is explained.
- A summary of YARN application frameworks is provided.

The introduction of Hadoop version 2 has drastically increased the number and scope of new applications. By splitting the version 1 monolithic MapReduce engine into two parts, a scheduler and the MapReduce framework, Hadoop has become a general-purpose large-scale data analytics platform. A simple example of a non-MapReduce Hadoop application is the YARN Distributed-Shell described in this chapter. As the number of non-MapReduce application frameworks continues to grow, the user's ability to navigate the data lake increases.

YARN DISTRIBUTED-SHELL

The Hadoop YARN project includes the Distributed-Shell application, which is an example of a Hadoop non-MapReduce application built on top of YARN. Distributed-Shell is a simple mechanism for running shell commands and scripts in containers on multiple nodes in a Hadoop cluster. This application is not meant to be a production administration tool, but rather a demonstration of the non-MapReduce capability that can be implemented on top of YARN. There are multiple mature implementations of a distributed shell that administrators typically use to manage a cluster of machines.

In addition, Distributed-Shell can be used as a starting point for exploring and building Hadoop YARN applications. This chapter offers guidance on how the Distributed-Shell can be used to understand the operation of YARN applications.

USING THE YARN DISTRIBUTED-SHELL

For the purpose of the examples presented in the remainder of this chapter, we assume and assign the following installation path, based on Hortonworks HDP 2.2, the Distributed-Shell application:

[Click here to view code image](#)

```
$ export YARN_DS=/usr/hdp/current/hadoop-yarn-client/hadoop-yarn-applications-distributedshell.jar
```

For the pseudo-distributed install using Apache Hadoop version 2.6.0, the following path will run the Distributed-Shell application (assuming \$HADOOP_HOME is defined to reflect the location Hadoop):

[Click here to view code image](#)

```
$ export YARN_DS=$HADOOP_HOME/share/hadoop/yarn/hadoop-yarn-applications-distributedshell-2.6.0.jar
```

If another distribution is used, search for the file `hadoop-yarn-applications-distributedshell*.jar` and set `$YARN_DS` based on its location. Distributed-Shell exposes various options that can be found by running the following command:

[Click here to view code image](#)

```
$ yarn org.apache.hadoop.yarn.applications.distributedshell.Client -jar $YARN_DS -help
```

The output of this command follows; we will explore some of these options in the examples illustrated in this chapter.

[Click here to view code image](#)

usage: Client	
-appname <arg>	Application Name. Default value - DistributedShell
-attempt_failures_validity_interval <arg>	when attempt_failures_validity_ interval in milliseconds is set to > 0, the failure number will not take failures which happen out of the validityInterval into failure count. If failure count reaches to maxAppAttempts, the application will be failed.
-container_memory <arg>	Amount of memory in MB to be requested to run the shell command
-container_vcores <arg>	Amount of virtual cores to be requested to run the shell command
-create	Flag to indicate whether to create the domain specified with -domain.
-debug	Dump out debug information
-domain <arg>	ID of the timeline domain where the timeline entities will be put
-help	Print usage
-jar <arg>	Jar file containing the application master
-keep_containers_across_application_attempts	Flag to indicate whether to keep containers across application attempts. If

-log_properties <arg>	the flag is true, running containers will not be killed when application attempt fails and these containers will be retrieved by the new application attempt
-master_memory <arg>	log4j.properties file Amount of memory in MB to be requested to run the application master
-master_vcores <arg>	Amount of virtual cores to be requested to run the application master
-modify_acls <arg>	Users and groups that allowed to modify the timeline entities in the given domain
-node_label_expression <arg>	Node label expression to determine the nodes where all the containers of this application will be allocated, "" means containers can be allocated anywhere, if you don't specify the option, default node_label_expression of queue will be used.
-num_containers <arg>	No. of containers on which the shell command needs to be executed
-priority <arg>	Application Priority. Default 0
-queue <arg>	RM Queue in which this application is to be submitted
-shell_args <arg>	Command line args for the shell script. Multiple args can be separated by empty space.
-shell_cmd_priority <arg>	Priority for the shell command containers
-shell_command <arg>	Shell command to be executed by the Application Master. Can

	only specify either --shell_command or --shell_script
-shell_env <arg>	Environment for shell script. Specified as env_key=env_val pairs
-shell_script <arg>	Location of the shell script to be executed. Can only specify either --shell_command or --shell_script
-timeout <arg>	Application timeout in milliseconds
-view_acls <arg>	Users and groups that allowed to view the timeline entities in the given domain

A Simple Example

The simplest use-case for the Distributed-Shell application is to run an arbitrary shell command in a container. We will demonstrate the use of the `uptime` command as an example. This command is run on the cluster using Distributed-Shell as follows:

[Click here to view code image](#)

```
$ yarn org.apache.hadoop.yarn.applications.distributedshell.Client -jar $YARN_DS -shell_command uptime
```

By default, Distributed-Shell spawns only one instance of a given shell command. When this command is run, you can see progress messages on the screen but nothing about the actual shell command. If the shell command succeeds, the following should appear at the end of the output:

[Click here to view code image](#)

```
15/05/27 14:48:53 INFO distributedshell.Client: Application completed successfully
```

If the shell command did not work for whatever reason, the following message will be displayed:

[Click here to view code image](#)

```
15/05/27 14:58:42 ERROR distributedshell.Client: Application failed to complete successfully
```

The next step is to examine the output for the application. Distributed-Shell redirects the output of the individual shell commands run on the cluster nodes into the log files, which are found either on the individual nodes or aggregated onto HDFS, depending on whether log aggregation is enabled.

Assuming log aggregation is enabled, the results for each instance of the command can be found by using the `yarn logs` command. For the previous `uptime` example, the following command can be used to inspect the logs:

[**Click here to view code image**](#)

```
$ yarn logs -applicationId application_1432831236474_0001
```

Note

The applicationId can be found from the program output or by using the `yarn application` command (see the “[Managing YARN Jobs](#)” section in [Chapter 10](#), “[Basic Hadoop Administration Procedures](#)”).

The abbreviated output follows:

[**Click here to view code image**](#)

```
Container: container_1432831236474_0001_01_000001 on n0_45454
```

```
=====
```

```
LogType:AppMaster.stderr
```

```
Log Upload Time:Thu May 28 12:41:58 -0400 2015
```

```
LogLength:3595
```

```
Log Contents:
```

```
15/05/28 12:41:52 INFO distributedshell.ApplicationMaster: Initializing  
ApplicationMaster
```

```
[...]
```

```
Container: container_1432831236474_0001_01_000002 on n1_45454
```

```
=====
```

```
LogType:stderr
```

```
Log Upload Time:Thu May 28 12:41:59 -0400 2015
```

```
LogLength:0
```

```
Log Contents:
```

```
LogType:stdout
```

```
Log Upload Time:Thu May 28 12:41:59 -0400 2015
```

```
LogLength:71
```

```
Log Contents:
```

```
12:41:56 up 33 days, 19:28, 0 users, load average: 0.08, 0.06, 0.01
```

Notice that there are two containers. The first container (`con..._000001`) is the ApplicationMaster for the job. The second container (`con..._000002`) is the actual shell script. The output for the `uptime` command is located in the second containers `stdout` after the `Log Contents:` label.

Using More Containers

Distributed-Shell can run commands to be executed on any number of containers by way of the `-num_containers` argument. For example, to see on which nodes the Distributed-Shell command was run, the following command can be used:

[Click here to view code image](#)

```
$ yarn org.apache.hadoop.yarn.applications.distributedshell.Client -jar $YARN_DS -  
shell_command hostname -num_containers 4
```

If we now examine the results for this job, there will be five containers in the log. The four command containers (2 through 5) will print the name of the node on which the container was run.

Distributed-Shell Examples with Shell Arguments

Arguments can be added to the shell command using the `-shell_args` option. For example, to do a `ls -l` in the directory from where the shell command was run, we can use the following commands:

[Click here to view code image](#)

```
$ yarn org.apache.hadoop.yarn.applications.distributedshell.Client -jar $YARN_DS -  
shell_command ls -shell_args -l
```

The resulting output from the log file is as follows:

[Click here to view code image](#)

```
total 20  
-rw-r--r-- 1 yarn hadoop 74 May 28 10:37 container_tokens  
-rwx----- 1 yarn hadoop 643 May 28 10:37 default_container_executor_session.sh  
-rwx----- 1 yarn hadoop 697 May 28 10:37 default_container_executor.sh  
-rwx----- 1 yarn hadoop 1700 May 28 10:37 launch_container.sh  
drwx--x-- 2 yarn hadoop 4096 May 28 10:37 tmp
```

As can be seen, the resulting files are new and not located anywhere in HDFS or the local file system. When we explore further by giving a `pwd` command for Distributed-Shell, the following directory is listed and created on the node that ran the shell command:

[Click here to view code image](#)

```
/hdfs2/hadoop/yarn/local/usercache/hdfs/appcache/application_1432831236474_0003/  
container_1432831236474_0003_01_000002/
```

Searching for this directory will prove to be problematic because these transient files are used by YARN to run the Distributed-Shell application and are removed once the application finishes. You can preserve these files for a specific interval by adding the following lines to the `yarn-site.xml` configuration file and restarting YARN:

[Click here to view code image](#)

```
<property>  
  <name>yarn.nodemanager.delete.debug-delay-sec</name>
```

```
<value>100000</value>
</property>
```

Choose a delay, in seconds, to preserve these files, and remember that all applications will create these files. If you are using Ambari, look on the YARN Configs tab under the Advanced yarn-site options, make the change and restart YARN. (See [Chapter 9, “Managing Hadoop with Apache Ambari,”](#) for more information on Ambari administration.) These files will be retained on the individual nodes only for the duration of the specified delay.

When debugging or investigating YARN applications, these files—in particular, `launch_container.sh`—offer important information about YARN processes. Distributed-Shell can be used to see what this file contains. Using DistributedShell, the contents of the `launch_container.sh` file can be printed with the following command:

[**Click here to view code image**](#)

```
$ yarn org.apache.hadoop.yarn.applications.distributedshell.Client -jar $YARN_DS -
shell_command cat -shell_args launch_container.sh
```

This command prints the `launch_container.sh` file that is created and run by YARN. The contents of this file are shown in [Listing 8.1](#). The file basically exports some important YARN variables and then, at the end, “execs” the command (`cat launch_container.sh`) directly and sends any output to logs.

Listing 8.1 Distributed-Shell `launch_container.sh` File

[**Click here to view code image**](#)

```
#!/bin/bash

export NM_HTTP_PORT="8042"
export
LOCAL_DIRS="/opt/hadoop/yarn/local/usercache/hdfs/appcache/
application_1432816241597_0004,/hdfs1/hadoop/yarn/local/usercach
e/hdfs/appcache/
application_1432816241597_0004,/hdfs2/hadoop/yarn/local/usercach
e/hdfs/appcache/
application_1432816241597_0004"
export JAVA_HOME="/usr/lib/jvm/java-1.7.0-openjdk.x86_64"
export
NM_AUX_SERVICE_mapreduce_shuffle="AAA0+gAAAAAAAAAAAAAA
AAAAAA
AAA=
"
export HADOOP_YARN_HOME="/usr/hdp/current/hadoop-yarn-client"
export
HADOOP_TOKEN_FILE_LOCATION="/hdfs2/hadoop/yarn/local/usercache/h
dfs/
```

```

appcache/application_1432816241597_0004/container_1432816241597_
0004_01_000002/
container_tokens"
export NM_HOST="limulus"
export JVM_PID="$@"
export USER="hdfs"
export PWD="/hdfs2/hadoop/yarn/local/usercache/hdfs/appcache/
application_1432816241597_0004/container_1432816241597_0004_01_0
00002"
export CONTAINER_ID="container_1432816241597_0004_01_000002"
export NM_PORT="45454"
export HOME="/home/"
export LOGNAME="hdfs"
export HADOOP_CONF_DIR="/etc/hadoop/conf"
export MALLOC_ARENA_MAX="4"
export
LOG_DIRS="/opt/hadoop/yarn/log/application_1432816241597_0004/
container_1432816241597_0004_01_000002,/hdfs1/hadoop/yarn/log/
application_1432816241597_0004/container_1432816241597_0004_01_0
00002,/hdfs2/
hadoop/yarn/log/application_1432816241597_0004/
container_1432816241597_0004_01_000002"
exec /bin/bash -c "cat launch_container.sh
1>/hdfs2/hadoop/yarn/log/application_1432816241597_0004/
container_1432816241597_0004_01_000002/stdout
2>/hdfs2/hadoop/yarn/log/
application_1432816241597_0004/container_1432816241597_0004_01_0
00002/stderr "
hadoop_shell_errorcode=$?
if [ $hadoop_shell_errorcode -ne 0 ]
then
    exit $hadoop_shell_errorcode
fi

```

There are more options for the Distributed-Shell that you can test. The real value of the Distributed-Shell application is its ability to demonstrate how applications are launched within the Hadoop YARN infrastructure. It is also a good starting point when you are creating YARN applications.

STRUCTURE OF YARN APPLICATIONS

A full explanation of writing YARN programs is beyond the scope of this book. The structure and operation of a YARN application are covered briefly in this section. For further information on writing YARN applications, consult *Apache Hadoop YARN: Moving beyond MapReduce and Batch Processing with Apache Hadoop 2* (see the references listed at the end of this chapter).

As mentioned in [Chapter 1](#), “[Background and Concepts](#),” the central YARN ResourceManager runs as a scheduling daemon on a dedicated machine and acts as the central authority for

allocating resources to the various competing applications in the cluster. The ResourceManager has a central and global view of all cluster resources and, therefore, can ensure fairness, capacity, and locality are shared across all users. Depending on the application demand, scheduling priorities, and resource availability, the ResourceManager dynamically allocates resource containers to applications to run on particular nodes. A container is a logical bundle of resources (e.g., memory, cores) bound to a particular cluster node. To enforce and track such assignments, the ResourceManager interacts with a special system daemon running on each node called the NodeManager. Communications between the ResourceManager and NodeManagers are heartbeat based for scalability. NodeManagers are responsible for local monitoring of resource availability, fault reporting, and container life-cycle management (e.g., starting and killing jobs). The ResourceManager depends on the NodeManagers for its “global view” of the cluster.

User applications are submitted to the ResourceManager via a public protocol and go through an admission control phase during which security credentials are validated and various operational and administrative checks are performed. Those applications that are accepted pass to the scheduler and are allowed to run. Once the scheduler has enough resources to satisfy the request, the application is moved from an accepted state to a running state. Aside from internal bookkeeping, this process involves allocating a container for the single ApplicationMaster and spawning it on a node in the cluster. Often called container 0, the ApplicationMaster does not have any additional resources at this point, but rather must request additional resources from the ResourceManager.

The ApplicationMaster is the “master” user job that manages all application life-cycle aspects, including dynamically increasing and decreasing resource consumption (i.e., containers), managing the flow of execution (e.g., in case of MapReduce jobs, running reducers against the output of maps), handling faults and computation skew, and performing other local optimizations. The ApplicationMaster is designed to run arbitrary user code that can be written in any programming language, as all communication with the ResourceManager and NodeManager is encoded using extensible network protocols (i.e., Google Protocol Buffers, <http://code.google.com/p/protobuf/>).

YARN makes few assumptions about the ApplicationMaster, although in practice it expects most jobs will use a higher-level programming framework. By delegating all these functions to ApplicationMasters, YARN’s architecture gains a great deal of scalability, programming model flexibility, and improved user agility. For example, upgrading and testing a new MapReduce framework can be done independently of other running MapReduce frameworks.

Typically, an ApplicationMaster will need to harness the processing power of multiple servers to complete a job. To achieve this, the ApplicationMaster issues resource requests to the ResourceManager. The form of these requests includes specification of locality preferences (e.g., to accommodate HDFS use) and properties of the containers. The ResourceManager will attempt to satisfy the resource requests coming from each application according to availability and scheduling policies. When a resource is scheduled on behalf of an ApplicationMaster, the ResourceManager generates a lease for the resource, which is acquired by a subsequent ApplicationMaster heartbeat. The ApplicationMaster then works with the NodeManagers to start the resource. A token-based security mechanism guarantees its authenticity when the ApplicationMaster presents the container lease to the NodeManager. In a typical situation, running containers will communicate with the ApplicationMaster through an application-specific protocol to report status and health information and to receive framework-specific commands. In

this way, YARN provides a basic infrastructure for monitoring and life-cycle management of containers, while each framework manages application-specific semantics independently. This design stands in sharp contrast to the original Hadoop version 1 design, in which scheduling was designed and integrated around managing only MapReduce tasks.

Figure 8.1 illustrates the relationship between the application and YARN components. The YARN components appear as the large outer boxes (ResourceManager and NodeManagers), and the two applications appear as smaller boxes (containers), one dark and one light. Each application uses a different ApplicationMaster; the darker client is running a Message Passing Interface (MPI) application and the lighter client is running a traditional MapReduce application.

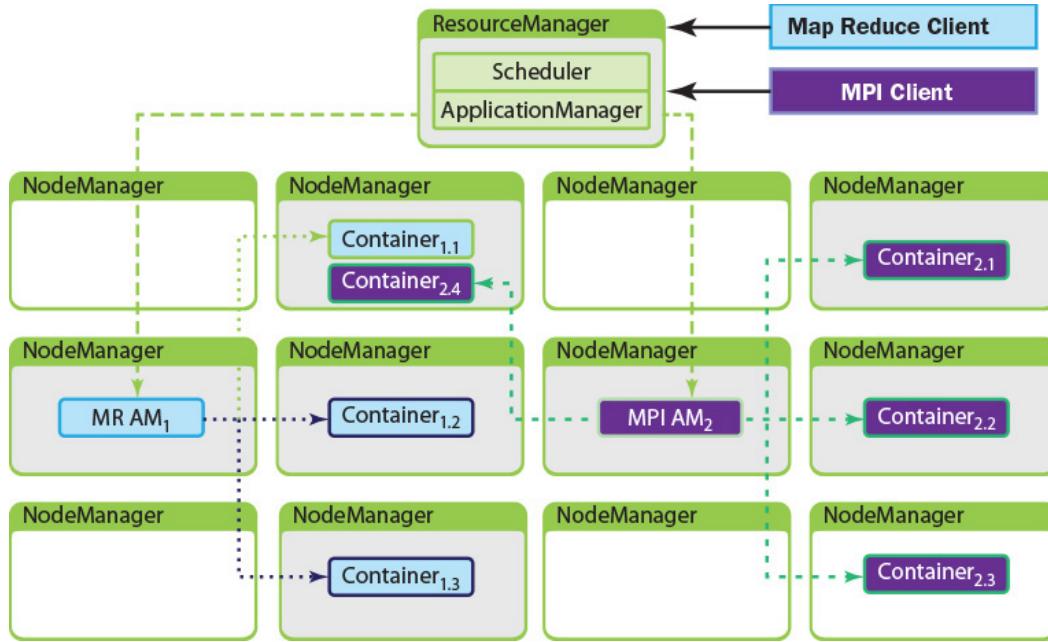


Figure 8.1 YARN architecture with two clients (MapReduce and MPI). The darker client (MPI AM₂) is running an MPI application, and the lighter client (MR AM₁) is running a MapReduce application. (From Arun C. Murthy, et al., *Apache Hadoop™ YARN*, copyright © 2014, p. 45. Reprinted and electronically reproduced by permission of Pearson Education, Inc., New York, NY.)

YARN APPLICATION FRAMEWORKS

One of the most exciting aspects of Hadoop version 2 is the capability to run all types of applications on a Hadoop cluster. In Hadoop version 1, the only processing model available to users is MapReduce. In Hadoop version 2, MapReduce is separated from the resource management layer of Hadoop and placed into its own application framework. Indeed, the growing number of YARN applications offers a high level and multifaceted interface to the Hadoop data lake discussed in Chapter 1.

YARN presents a resource management platform, which provides services such as scheduling, fault monitoring, data locality, and more to MapReduce and other frameworks. Figure 8.2 illustrates some of the various frameworks that will run under YARN. Note that the Hadoop version 1 applications (e.g., Pig and Hive) run under the MapReduce framework.

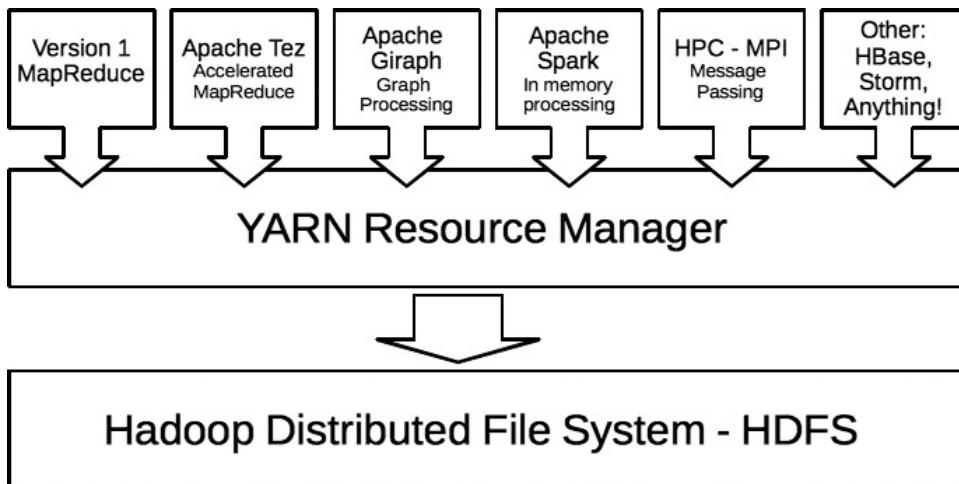


Figure 8.2 Example of the Hadoop version 2 ecosystem. Hadoop version 1 supports batch MapReduce applications only.

This section presents a brief survey of emerging open source YARN application frameworks that are being developed to run under YARN. As of this writing, many YARN frameworks are under active development and the framework landscape is expected to change rapidly. Commercial vendors are also taking advantage of the YARN platform. Consult the webpage for each individual framework for full details of its current stage of development and deployment.

Distributed-Shell

As described earlier in this chapter, Distributed-Shell is an example application included with the Hadoop core components that demonstrates how to write applications on top of YARN. It provides a simple method for running shell commands and scripts in containers in parallel on a Hadoop YARN cluster.

Hadoop MapReduce

MapReduce was the first YARN framework and drove many of YARN's requirements. It is integrated tightly with the rest of the Hadoop ecosystem projects, such as Apache Pig, Apache Hive, and Apache Oozie.

Apache Tez

One great example of a new YARN framework is Apache Tez. Many Hadoop jobs involve the execution of a complex directed acyclic graph (DAG) of tasks using separate MapReduce stages. Apache Tez generalizes this process and enables these tasks to be spread across stages so that they can be run as a single, all-encompassing job.

Tez can be used as a MapReduce replacement for projects such as Apache Hive and Apache Pig. No changes are needed to the Hive or Pig applications. For more information, see <https://tez.apache.org>.

Apache Giraph

Apache Giraph is an iterative graph processing system built for high scalability. Facebook, Twitter, and LinkedIn use it to create social graphs of users. Giraph was originally written to run on standard Hadoop V1 using the MapReduce framework, but that approach proved inefficient and totally unnatural for various reasons. The native Giraph implementation under YARN provides the user with an iterative processing model that is not directly available with MapReduce. Support for YARN has been present in Giraph since its own version 1.0 release. In addition, using the flexibility of YARN, the Giraph developers plan on implementing their own web interface to monitor job progress. For more information, see <http://giraph.apache.org>.

Hoya: HBase on YARN

The Hoya project creates dynamic and elastic Apache HBase clusters on top of YARN. A client application creates the persistent configuration files, sets up the HBase cluster XML files, and then asks YARN to create an ApplicationMaster. YARN copies all files listed in the client's application-launch request from HDFS into the local file system of the chosen server, and then executes the command to start the Hoya ApplicationMaster. Hoya also asks YARN for the number of containers matching the number of HBase region servers it needs. For more information, see <http://hortonworks.com/blog/introducing-hoya-hbase-on-yarn>.

Dryad on YARN

Similar to Apache Tez, Microsoft's Dryad provides a DAG as the abstraction of execution flow. This framework is ported to run natively on YARN and is fully compatible with its non-YARN version. The code is written completely in native C++ and C# for worker nodes and uses a thin layer of Java within the application. For more information, see <http://research.microsoft.com/en-us/projects/dryad>.

Apache Spark

Spark was initially developed for applications in which keeping data in memory improves performance, such as iterative algorithms, which are common in machine learning, and interactive data mining. Spark differs from classic MapReduce in two important ways. First, Spark holds intermediate results in memory, rather than writing them to disk. Second, Spark supports more than just MapReduce functions; that is, it greatly expands the set of possible analyses that can be executed over HDFS data stores. It also provides APIs in Scala, Java, and Python.

Since 2013, Spark has been running on production YARN clusters at Yahoo!. The advantage of porting and running Spark on top of YARN is the common resource management and a single underlying file system. For more information, see <https://spark.apache.org>.

Apache Storm

Traditional MapReduce jobs are expected to eventually finish, but Apache Storm continuously processes messages until it is stopped. This framework is designed to process unbounded streams of data in real time. It can be used in any programming language. The basic Storm use-cases include real-time analytics, online machine learning, continuous computation, distributed RPC (remote procedure calls), ETL (extract, transform, and load), and more. Storm provides fast performance, is scalable, is fault tolerant, and provides processing guarantees. It works directly

under YARN and takes advantage of the common data and resource management substrate. For more information, see <http://storm.apache.org>.

Apache REEF: Retainable Evaluator Execution Framework

YARN's flexibility sometimes requires significant effort on the part of application implementers. The steps involved in writing a custom application on YARN include building your own ApplicationMaster, performing client and container management, and handling aspects of fault tolerance, execution flow, coordination, and other concerns. The REEF project by Microsoft recognizes this challenge and factors out several components that are common to many applications, such as storage management, data caching, fault detection, and checkpoints. Framework designers can build their applications on top of REEF more easily than they can build those same applications directly on YARN, and can reuse these common services/libraries. REEF's design makes it suitable for both MapReduce and DAG-like executions as well as iterative and interactive computations. For more information, see <http://www.reef-project.org/welcome>.

Hamster: Hadoop and MPI on the Same Cluster

The Message Passing Interface (MPI) is widely used in high-performance computing (HPC). MPI is primarily a set of optimized message-passing library calls for C, C++, and Fortran that operate over popular server interconnects such as Ethernet and InfiniBand. Because users have full control over their YARN containers, there is no reason why MPI applications cannot run within a Hadoop cluster. The Hamster effort is a work-in-progress that provides a good discussion of the issues involved in mapping MPI to a YARN cluster (see <https://issues.apache.org/jira/browse/MAPREDUCE-2911>). Currently, an alpha version of MPICH2 is available for YARN that can be used to run MPI applications. For more information, see <https://github.com/clarkyzl/mpich2-yarn>.

Apache Flink: Scalable Batch and Stream Data Processing

Apache Flink is a platform for efficient, distributed, general-purpose data processing. It features powerful programming abstractions in Java and Scala, a high-performance run time, and automatic program optimization. It also offers native support for iterations, incremental iterations, and programs consisting of large DAGs of operations.

Flink is primarily a stream-processing framework that can look like a batch-processing environment. The immediate benefit from this approach is the ability to use the same algorithms for both streaming and batch modes (exactly as is done in Apache Spark). However, Flink can provide low-latency similar to that found in Apache Storm, but which is not available in Apache Spark.

In addition, Flink has its own memory management system, separate from Java's garbage collector. By managing memory explicitly, Flink almost eliminates the memory spikes often seen on Spark clusters. For more information, see <https://flink.apache.org>.

Apache Slider: Dynamic Application Management

Apache Slider (incubating) is a YARN application to deploy existing distributed applications on YARN, monitor them, and make them larger or smaller as desired in real time.

Applications can be stopped and then started; the distribution of the deployed application across the YARN cluster is persistent and allows for best-effort placement close to the previous locations. Applications that remember the previous placement of data (such as HBase) can exhibit fast startup times by capitalizing on this feature.

YARN monitors the health of “YARN containers” that are hosting parts of the deployed applications. If a container fails, the Slider manager is notified. Slider then requests a new replacement container from the YARN ResourceManager. Some of Slider’s other features include user creation of on-demand applications, the ability to stop and restart applications as needed (preemption), and the ability to expand or reduce the number of application containers as needed. The Slider tool is a Java command-line application. For more information, see <http://slider.incubator.apache.org>.

SUMMARY AND ADDITIONAL RESOURCES

The Hadoop YARN Distributed-Shell is a simple demonstration of a non-MapReduce program. It can be used to learn about how YARN operates and launches jobs across the cluster. The structure and operation of YARN programs are designed to provide a highly scalable and flexible method to create Hadoop applications. Other resources listed here describe application development with this framework.

YARN application frameworks offer many new capabilities for the Hadoop data lake. Many new algorithms and programming models are available, including an optimized MapReduce engine using Apache Tez. Each framework description in this chapter provides a webpage reference from which to find more information.

■ Apache Hadoop YARN Development

■ Book: Murthy, A., et al. 2014. *Apache Hadoop YARN: Moving beyond MapReduce and Batch Processing with Apache Hadoop 2*, Boston, MA: Addison-Wesley. <http://www.informit.com/store/apache-hadoop-yarn-moving-beyond-mapreduce-and-batch-9780321934505>

■ <http://hadoop.apache.org/docs/r2.7.0/hadoop-yarn/hadoop-yarn-site/WritingYarnApplications.html>

■ MemcacheD on Yarn; <http://hortonworks.com/blog/how-to-deploy-memcached-on-yarn/>

■ Hortonworks YARN resources; <http://hortonworks.com/get-started/yarn>

■ Apache Hadoop YARN Frameworks

■ See the webpage reference at the end of each individual description.

9. Managing Hadoop with Apache Ambari

In This Chapter:

- A tour of the Apache Ambari graphical management tool is provided.
- The procedure for restarting a stopped Hadoop service is explained.
- The procedure for changing Hadoop properties and tracking configurations is presented.

Managing a Hadoop installation by hand can be tedious and time consuming. In addition to keeping configuration files synchronized across a cluster, starting, stopping, and restarting Hadoop services and dependent services in the right order is not a simple task. The Apache Ambari graphical management tool is designed to help you easily manage these and other Hadoop administrative issues. This chapter provides some basic navigation and usage scenarios for Apache Ambari.

Apache Ambari is an open source graphical installation and management tool for Apache Hadoop version 2. Ambari was used in [Chapter 2](#), “[Installation Recipes](#),” to install Hadoop and related packages across a four-node cluster. In particular, the following packages were installed: HDFS, YARN, MapReduce2, Tez, Nagios, Ganglia, Hive, HBase, Pig, Sqoop, Oozie, Zookeeper, and Flume. These packages have been described in other chapters and provide basic Hadoop functionality. As noted in [Chapter 2](#), other packages are available for installation (refer to [Figure 2.18](#)). Finally, to use Ambari as a management tool, the entire installation process must be done using Ambari. It is not possible to use Ambari for Hadoop clusters that have been installed by other means.

Along with being an installation tool, Ambari can be used as a centralized point of administration for a Hadoop cluster. Using Ambari, the user can configure cluster services, monitor the status of cluster hosts (nodes) or services, visualize hotspots by service metric, start or stop services, and add new hosts to the cluster. All of these features infuse a high level of agility into the processes of managing and monitoring a distributed computing environment. Ambari also attempts to provide real-time reporting of important metrics.

Apache Ambari continues to undergo rapid change. The description in this chapter is based on version 1.7. The major aspects of Ambari, which will not change, are explained in the following sections. Further detailed information can be found at <https://ambari.apache.org>.

QUICK TOUR OF APACHE AMBARI

After completing the initial installation and logging into Ambari (as explained in [Chapter 2](#)), a dashboard similar to that shown in [Figure 9.1](#) is presented. The same four-node cluster as created in [Chapter 2](#) will be used to explore Ambari in this chapter. If you need to reopen the Ambari dashboard interface, simply enter the following command (which assumes you are using the Firefox browser, although other browsers may also be used):

```
$ firefox localhost:8080
```

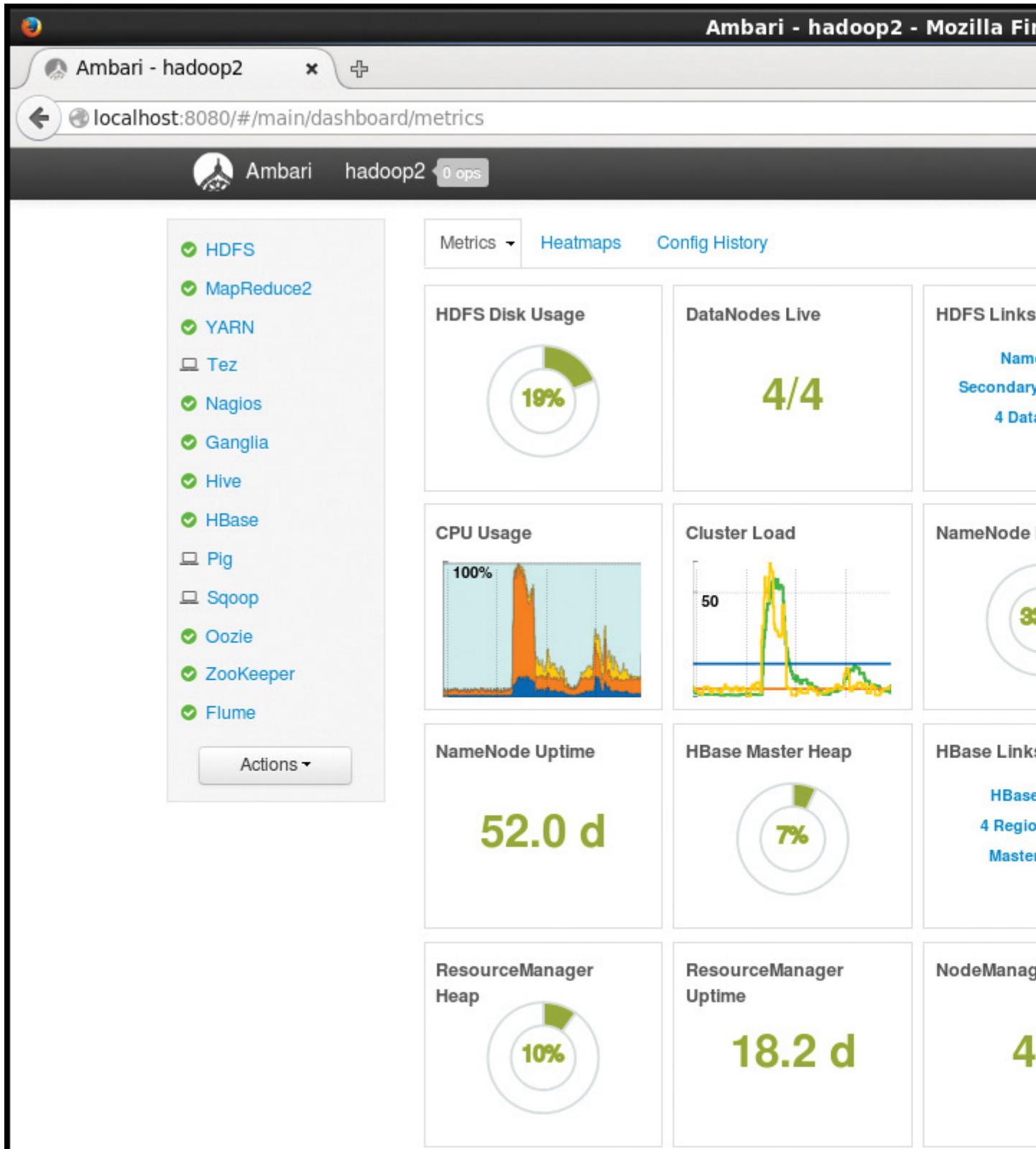


Figure 9.1 Apache Ambari dashboard view of a Hadoop cluster

The default login and password are `admin` and `admin`, respectively. Before continuing any further, you should change the default password. To change the password, select Manage Ambari from the Admin pull-down menu in the upper-right corner. In the management window, click Users under User + Group Management, and then click the `admin` user name. Select Change Password and enter a new password. When you are finished, click the Go To Dashboard link on the left side of the window to return to the dashboard view.

To leave the Ambari interface, select the Admin pull-down menu at the left side of the main menu bar and click Sign out.

The dashboard view provides a number of high-level metrics for many of the installed services. A glance at the dashboard should allow you to get a sense of how the cluster is performing.

The top navigation menu bar, shown in [Figure 9.1](#), provides access to the Dashboard, Services, Hosts, Admin, and Views features (the 3×3 cube is the Views menu). The status (up/down) of various Hadoop services is displayed on the left using green/orange dots. Note that two of the services managed by Ambari are Nagios and Ganglia; the standard cluster management services installed by Ambari, they are used to provide cluster monitoring (Nagios) and metrics (Ganglia).

Dashboard View

The Dashboard view provides small status widgets for many of the services running on the cluster. The actual services are listed on the left-side vertical menu. These services correspond to what was installed in [Chapter 2](#). You can move, edit, remove, or add these widgets as follows:

- **Moving:** Click and hold a widget while it is moved about the grid.
- **Edit:** Place the mouse on the widget and click the gray edit symbol in the upper-right corner of the widget. You can change several different aspects (including thresholds) of the widget.
- **Remove:** Place the mouse on the widget and click the X in the upper-left corner.
- **Add:** Click the small triangle next to the Metrics tab and select Add. The available widgets will be displayed. Select the widgets you want to add and click Apply.

Some widgets provide additional information when you move the mouse over them. For instance, the DataNodes widget displays the number of live, dead, and decommissioning hosts. Clicking directly on a graph widget provides an enlarged view. For instance, [Figure 9.2](#) provides a detailed view of the CPU Usage widget from [Figure 9.1](#).

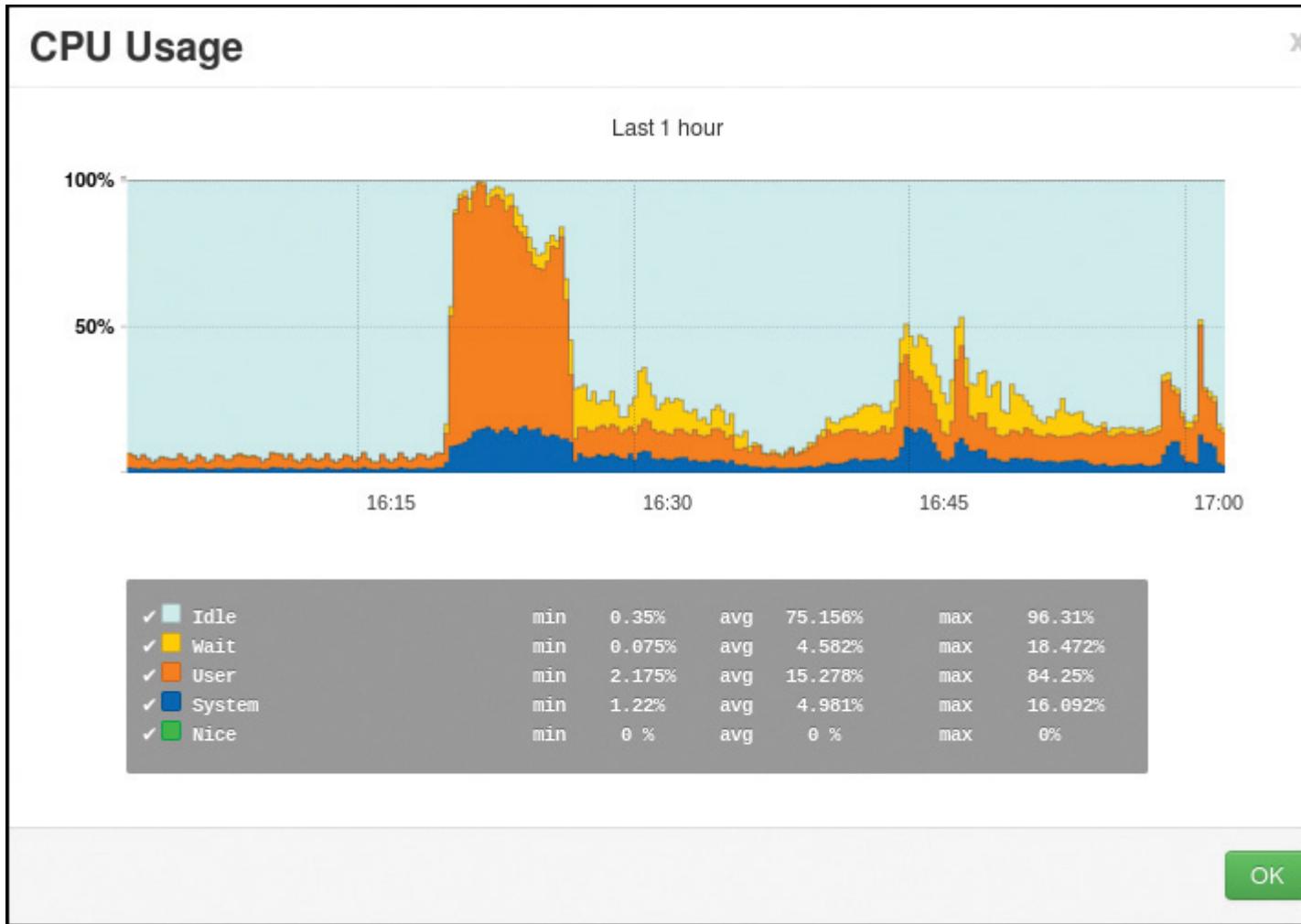


Figure 9.2 Enlarged view of Ambari CPU Usage widget

The Dashboard view also includes a heatmap view of the cluster. Cluster heatmaps physically map selected metrics across the cluster. When you click the Heatmaps tab, a heatmap for the cluster will be displayed. To select the metric used for the heatmap, choose the desired option from the Select Metric pull-down menu. Note that the scale and color ranges are different for each metric. The heatmap for percentage host memory used is displayed in [Figure 9.3](#).

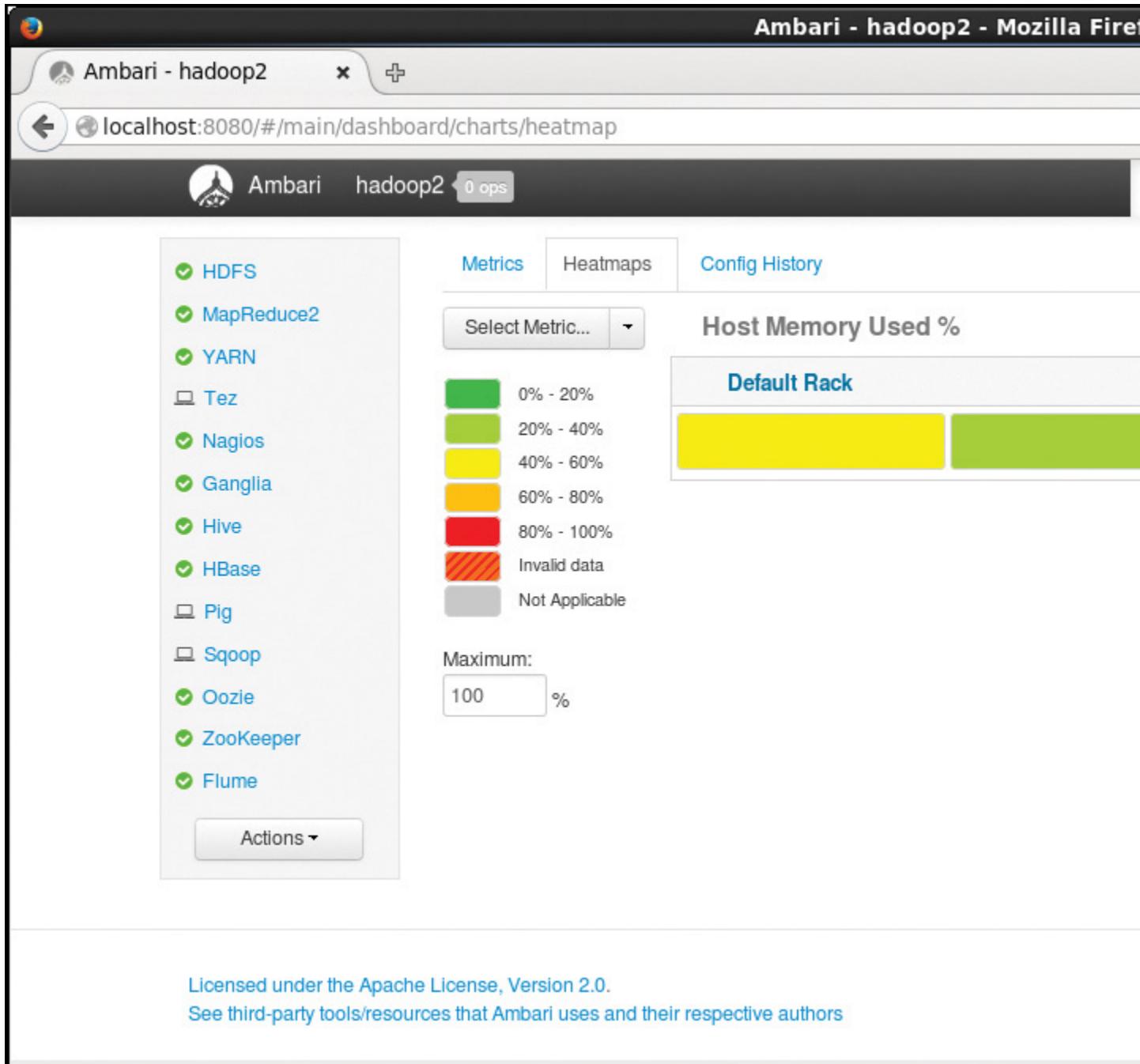


Figure 9.3 Ambari heatmap for Host memory usage

Configuration history is the final tab in the dashboard window. This view provides a list of configuration changes made to the cluster. As shown in [Figure 9.4](#), Ambari enables configurations to be sorted by Service, Configuration Group, Data, and Author. To find the specific configuration settings, click the service name. More information on configuration settings is provided later in the chapter.

The screenshot shows the Ambari master configuration changes list. The left sidebar lists services: HDFS, MapReduce2, YARN, Tez, Nagios, Ganglia, Hive, HBase, Pig, Sqoop, Oozie, ZooKeeper, and Flume. The main area has tabs for Metrics, Heatmaps, and Config History, with Config History selected. It shows a table of configuration versions for the YARN service. The table includes columns for Service (YARN), Config Group (YARN Default), and Created Date. The table lists 11 versions from V2 to V11, all of which are YARN Default and were created on May 26 or 27, 2015.

Service	Config Group	Created
V11 YARN	YARN Default Current	Thu, May 28, 2015
V10 YARN	YARN Default	Thu, May 28, 2015
V9 YARN	YARN Default	Thu, May 28, 2015
V8 YARN	YARN Default	Thu, May 28, 2015
V7 YARN	YARN Default	Thu, May 28, 2015
V6 YARN	YARN Default	Thu, May 28, 2015
V5 YARN	YARN Default	Wed, May 27, 2015
V4 YARN	YARN Default	Tue, May 26, 2015
V3 YARN	YARN Default	Tue, May 26, 2015
V2 YARN	YARN Default	Tue, May 26, 2015

Figure 9.4 Ambari master configuration changes list

Services View

The Services menu provides a detailed look at each service running on the cluster. It also provides a graphical method for configuring each service (i.e., instead of hand-editing the /etc/hadoop/confXML files). The summary tab provides a current Summary view of important service metrics and an Alerts and Health Checks sub-window.

Similar to the Dashboard view, the currently installed services are listed on the left-side menu. To select a service, click the service name in the menu. When applicable, each service will have its own Summary, Alerts and Health Monitoring, and Service Metrics windows. For example, [Figure 9.5](#) shows the Service view for HDFS. Important information such as the status of NameNode, SecondaryNameNode, DataNodes, uptime, and available disk space is displayed in the Summary window. The Alerts and Health Checks window provides the latest status of the service and its component systems. Finally, several important real-time service metrics are displayed as widgets at the bottom of the screen. As on the dashboard, these widgets can be expanded to display a more detailed view.

Ambari - hadoop2 - Mozilla Firefox

Ambari - hadoop2

localhost:8080/#/main/services/HDFS/summary

Ambari hadoop2 0 ops

HDFS

MapReduce2

YARN

Tez

Nagios

Ganglia

Hive

HBase

Pig

Sqoop

Oozie

ZooKeeper

Flume

Actions ▾

Summary Configs Quick Links

Summary

NameNode Started

SNameNode Started

DataNodes 4/4 DataNodes Live

NameNode Uptime 51.97 days

NameNode Heap 317.4 MB / 1004.0 MB (31.6% used)

DataNodes Status 4 live / 0 dead / 0 decommissioning

Disk Usage (DFS Used) 78.9 GB / 1.4 TB (5.64%)

Disk Usage (Non DFS Used) 211.7 GB / 1.4 TB (15.12%)

Disk Usage (Remaining) 1.1 TB / 1.4 TB (79.25%)

Blocks (total) 1350

Block Errors 0 corrupt / 0 missing / 2 under replicated

Total Files + Directories 300622

Upgrade Status No pending upgrade

Safe Mode Status Not in safe mode

HDFS Service Metrics

Total Space Utilization

931.3 GB

File Operations

1 ms

0.5 ms

2 ms

This screenshot shows the Ambari HDFS Summary page. On the left, a sidebar lists various services: HDFS (selected), MapReduce2, YARN, Tez, Nagios, Ganglia, Hive, HBase, Pig, Sqoop, Oozie, ZooKeeper, and Flume. Below the sidebar is an 'Actions' dropdown. At the top right are tabs for 'Summary' (selected) and 'Configs'. A 'Quick Links' section is also present. The main content area has two sections: 'Summary' and 'HDFS Service Metrics'. The 'Summary' section displays various system statistics. The 'HDFS Service Metrics' section contains three charts: 'Total Space Utilization' (a line chart showing utilization over time with a value of 931.3 GB highlighted), 'File Operations' (a bar chart showing operation rates in ops/s), and two smaller line charts at the bottom showing response times of 0.5 ms, 1 ms, and 2 ms.

Figure 9.5 HDFS service summary window

Clicking the Configs tab will open an options form, shown in [Figure 9.6](#), for the service. The options (properties) are the same ones that are set in the Hadoop XML files. When using Ambari, the user has complete control over the XML files and should manage them only through the Ambari interface—that is, the user should *not* edit the files by hand.

Ambari - hadoop2 - Mozilla Firefox

Ambari - hadoop2

localhost:8080/#/main/services/HDFS/configs

Ambari hadoop2 0 ops

HDFS

MapReduce2

YARN

Tez

Nagios

Ganglia

Hive

HBase

Pig

Sqoop

Oozie

ZooKeeper

Flume

Actions ▾

Summary Configs Quick Links

Group HDFS Default (4) Manage Config Groups

V1 admin 2 months ago Current

NameNode

NameNode hosts limulus

NameNode directories /hdfs1/hadoop/hdfs/namenode/

NameNode Java heap size 1024 MB

NameNode new generation size 200 MB

NameNode maximum new generation size 200 MB

NameNode permanent generation size 128 MB

NameNode maximum permanent generation size 256 MB

Secondary NameNode

Secondary NameNode

<https://hemanthrajhemu.github.io>

Figure 9.6 Ambari service options for HDFS

The current settings available for each service are shown in the form. The administrator can set each of these properties by changing the values in the form. Placing the mouse in the input box of the property displays a short description of each property. Where possible, properties are grouped by functionality. The form also has provisions for adding properties that are not listed. An example of changing service properties and restarting the service components is provided in the “[Managing Hadoop Services](#)” section.

If a service provides its own graphical interface (e.g., HDFS, YARN, Oozie), then that interface can be opened in a separate browser tab by using the Quick Links pull-down menu located in top middle of the window.

Finally, the Service Action pull-down menu in the upper-left corner provides a method for starting and stopping each service and/or its component daemons across the cluster. Some services may have a set of unique actions (such as rebalancing HDFS) that apply to only certain situations. Finally, every service has a Service Check option to make sure the service is working properly. The service check is initially run as part of the installation process and can be valuable when diagnosing problems. (See [Appendix B, “Getting Started Flowchart and Troubleshooting Guide.”](#))

Hosts View

Selecting the Hosts menu item provides the information shown in [Figure 9.7](#). The host name, IP address, number of cores, memory, disk usage, current load average, and Hadoop components are listed in this window in tabular form.

The screenshot shows the Ambari interface for managing hosts. At the top, there's a navigation bar with tabs for 'Ambari - hadoop2' and 'hadoop2'. Below the navigation is a toolbar with 'Actions' and 'Filter: All (4)'. The main area displays a table of hosts:

Name	IP Address	Cores (CPU)	RAM
<input type="checkbox"/> Any	Any	Any	Any
<input checked="" type="checkbox"/> limulus	10.0.0.1	4 (4)	23.51GB
<input checked="" type="checkbox"/> n0	10.0.0.10	4 (4)	15.60GB
<input checked="" type="checkbox"/> n1	10.0.0.11	4 (4)	15.60GB
<input checked="" type="checkbox"/> n2	10.0.0.12	4 (4)	15.60GB

At the bottom of the host list, it says '4 of 4 hosts showing - [clear filters](#)'.

Licensed under the Apache License, Version 2.0.
See third-party tools/resources that Ambari uses and their respective authors

Figure 9.7 Ambari main Hosts screen

To display the Hadoop components installed on each host, click the links in the rightmost columns. You can also add new hosts by using the Actions pull-down menu. The new host must be running the Ambari agent (or the root SSH key must be entered) and have the base software described in [Chapter 2](#) installed. The remaining options in the Actions pull-down menu provide control over the various service components running on the hosts.

Further details for a particular host can be found by clicking the host name in the left column. As shown in [Figure 9.8](#), the individual host view provides three sub-windows: Components, Host Metrics, and Summary information. The Components window lists the services that are currently running on the host. Each service can be stopped, restarted, decommissioned, or placed in maintenance mode. The Metrics window displays widgets that provide important metrics (e.g., CPU, memory, disk, and network usage). Clicking the widget displays a larger version of the

graphic. The Summary window provides basic information about the host, including the last time a heartbeat was received.

Ambari - hadoop2 - Mozilla Firefox

Ambari - hadoop2

localhost:8080/#/main/hosts/n0/summary

Ambari hadoop2 0 ops

n0 No alerts

Back

Summary Configs

Components

+ Add

✓ DataNode / HDFS	Started
✓ Flume / Flume	Started
✓ Ganglia Monitor / Ganglia	Started
✓ RegionServer / HBase	Started
✓ NodeManager / YARN	Started
Clients / HBase Client , HCat Client , HDFS Client , Hive Client , MapReduce2 Client , Oozie Client , Pig , Sqoop , Tez Client , YARN Client , ZooKeeper Client	Installed

Host Metrics

Summary

Hostname: n0
IP Address: 10.0.0.10
OS: centos6 (x86_64)
Cores (CPU): 4 (4)
Disk: 138.42GB/378.07GB (36.61% used)
Memory: 15.60GB
Load Avg: 0.07
Heartbeat: a moment ago

Licensed under the Apache License, Version 2.0.
See third-party tools/resources that Ambari uses and their respective authors

Figure 9.8 Ambari cluster host detail view

Admin View

The Administration (Admin) view provides three options. The first, as shown in [Figure 9.9](#), displays a list of installed software. This Repositories listing generally reflects the version of Hortonworks Data Platform (HDP) used during the installation process. The Service Accounts option lists the service accounts added when the system was installed. These accounts are used to run various services and tests for Ambari. The third option, Security, sets the security on the cluster. A fully secured Hadoop cluster is important in many instances and should be explored if a secure environment is needed. This aspect of Ambari is beyond the scope of this book.

Ambari - hadoop2 - Mozilla Firefox

Ambari - hadoop2

localhost:8080/#/main/admin/repositories

Ambari hadoop2 0 ops

Repositories

Service Accounts

Security

Cluster Stack Version: HDP-2.2

Service	Version	Description
Falcon	0.6.0.2.2.0.0	Data management and processing platform
Flume	1.5.2.2.2.0.0	A distributed service for collecting, aggregating, and moving large volumes of log data
Ganglia	3.5.0	Ganglia Metrics Collection system (RRD)
HBase	0.98.4.2.2.0.0	Non-relational distributed database and search engine
HDFS	2.6.0.2.2.0.0	Apache Hadoop Distributed File System
Hive	0.14.0.2.2.0.0	Data warehouse system for ad-hoc querying over large datasets
Kafka	0.8.1.2.2.0.0	A high-throughput distributed messaging system
Knox	0.5.0.2.2.0.0	Provides a single point of authentication and access control
Nagios	3.5.0	Nagios Monitoring and Alerting system
Oozie	4.1.0.2.2.0.0	System for workflow coordination and execution of distributed work flows. It includes an optional Oozie Web Console which relies on Apache Geronimo
Pig	0.14.0.2.2.0.0	Scripting platform for analyzing large data sets
Slider	0.60.0.2.2.0.0	A framework for deploying, managing and scaling distributed services
Sqoop	1.4.5.2.2.0.0	Tool for transferring bulk data between relational databases and HDFS
Storm	0.9.3.2.2.0.0	Apache Hadoop Stream processing framework
Tez	0.5.2.2.2.0.0	Tez is the next generation Hadoop Queue
YARN + MapReduce2	2.6.0.2.2.0.0	Apache Hadoop NextGen MapReduce
ZooKeeper	3.4.6.2.2.0.0	Centralized service which provides highly reliable distributed synchronization

Repositories

2.2

<https://hemanthrajhemu.github.io>

Figure 9.9 Ambari installed packages with versions, numbers, and descriptions

Views View

Ambari Views is a framework offering a systematic way to plug in user interface capabilities that provide for custom visualization, management, and monitoring features in Ambari. Views allows you to extend and customize Ambari to meet your specific needs. You can find more information about Ambari Views from the following

source: <https://cwiki.apache.org/confluence/display/AMBARI/Views>.

Admin Pull-Down Menu

The Administrative (Admin) pull-down menu provides the following options:

- **About**—Provides the current version of Ambari.
- **Manage Ambari**—Open the management screen where Users, Groups, Permissions, and Ambari Views can be created and configured.
- **Settings**—Provides the option to turn off the progress window. (See [Figure 9.15](#).)
- **Sign Out**—Exits the interface.

MANAGING HADOOP SERVICES

During the course of normal Hadoop cluster operation, services may fail for any number of reasons. Ambari monitors all of the Hadoop services and reports any service interruption to the dashboard. In addition, when the system was installed, an administrative email for the Nagios monitoring system was required. All service interruption notifications are sent to this email address.

[Figure 9.10](#) shows the Ambari dashboard reporting a down DataNode. The service error indicator numbers next to the HDFS service and Hosts menu item indicate this condition. The DataNode widget also has turned red and indicates that 3/4 DataNodes are operating.

Ambari - hadoop2 - Mozilla Firefox

Ambari - hadoop2 x Namenode information x +

localhost:8080/#/main/dashboard/metrics

Ambari hadoop2 0 ops

HDFS 2
MapReduce2
YARN
Tez
Nagios
Ganglia
Hive
HBase
Pig
Sqoop
Oozie
ZooKeeper
Flume

Metrics Heatmaps Config History

HDFS Disk Usage
19%

DataNodes Live
3/4

CPU Usage
100%

Cluster Load

NameNode Uptime
52.9 d

HBase Master Heap
4%

ResourceManager Heap
12%

ResourceManager Uptime
19.1 d

NameNode
Secondary
4 Data

NameNode
38

HBase Links
HBase
4 Region
Master

NodeManag
4

Licensed under the Apache License, Version 2.0.
See third-party-tools/resources that Ambari uses and their respective authors.

<https://nemanthrajhemu.github.io>

Figure 9.10 Ambari main dashboard indicating a DataNode issue

Clicking the HDFS service link in the left vertical menu will bring up the service summary screen shown in [Figure 9.11](#). The Alerts and Health Checks window confirms that a DataNode is down.

Ambari - hadoop2 - Mozilla Firefox

Ambari - hadoop2 x Namenode information x +

localhost:8080/#/main/services/HDFS/summary

Ambari hadoop2 0 ops

HDFS 2

MapReduce2

YARN

Tez

Nagios

Ganglia

Hive

HBase

Pig

Sqoop

Oozie

ZooKeeper

Flume

Actions ▾

Summary Configs Quick Links

Summary

NameNode Started

SNameNode Started

DataNodes 3/4 DataNodes Live

NameNode Uptime 52.86 days

NameNode Heap 355.6 MB / 1004.0 MB (35.4% used)

DataNodes Status 3 live / 1 dead / 0 decommissioning

Disk Usage (DFS Used) 73.8 GB / 1.0 TB (7.03%)

Disk Usage (Non DFS Used) 126.8 GB / 1.0 TB (12.08%)

Disk Usage (Remaining) 849.5 GB / 1.0 TB (80.90%)

Blocks (total) 1352

Block Errors 0 corrupt / 0 missing / 385 under replicated

Total Files + Directories 303574

Upgrade Status No pending upgrade

Safe Mode Status Not in safe mode

This screenshot shows the Ambari interface for the HDFS service. On the left, a sidebar lists various services with their status (e.g., HDFS 2, MapReduce2, YARN). The main area is titled 'Summary' and displays key metrics for the NameNode and DataNodes. It includes details like NameNode uptime, heap usage, disk usage (DFS, Non DFS, Remaining), total blocks, and upgrade status.

HDFS Service Metrics

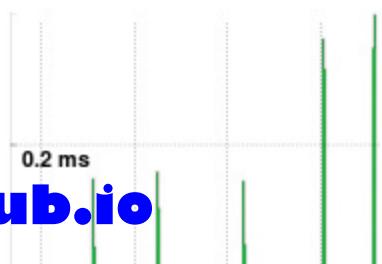
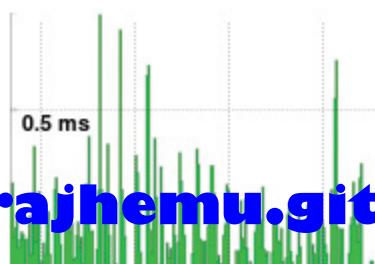
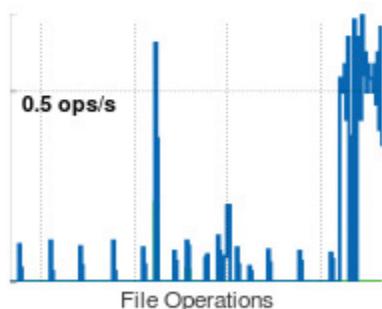
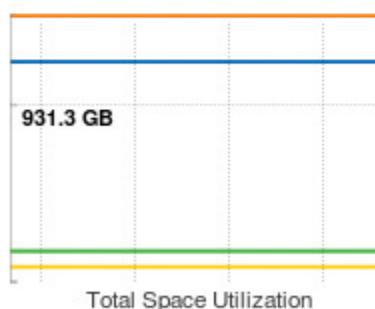


Figure 9.11 Ambari HDFS service summary window indicating a down DataNode

The specific host (or hosts) with an issue can be found by examining the Hosts window. As shown in [Figure 9.12](#), the status of host n1 has changed from a green dot with a check mark inside to a yellow dot with a dash inside. An orange dot with a question mark inside indicates the host is not responding and is probably down. Other service interruption indicators may also be set as a result of the unresponsive node.

The screenshot shows the Ambari Hosts screen. At the top, there are tabs for 'Ambari - hadoop2' and 'Namenode information'. Below the tabs, the URL is 'localhost:8080/#/main/hosts'. The main area has a header with 'Ambari', 'hadoop2', and '0 ops'. There are two buttons: 'Actions' and 'Filter: All (4)'. A table lists four hosts:

Name	IP Address	Cores (CPU)	RAM
limulus	10.0.0.1	4 (4)	23.51
n0	10.0.0.10	4 (4)	15.60
n1	10.0.0.11	4 (4)	15.60
n2	10.0.0.12	4 (4)	15.60

Below the table, it says '4 of 4 hosts showing - [clear filters](#)'. At the bottom, there is a note: 'Licensed under the Apache License, Version 2.0. See third-party tools/resources that Ambari uses and their respective authors'.

Figure 9.12 Ambari Hosts screen indicating an issue with host n1

Clicking on the n1 host link opens the view in [Figure 9.13](#). Inspecting the Components sub-window reveals that the DataNode daemon has stopped on the host. At this point, checking the DataNode logs on host n1 will help identify the actual cause of the failure. Assuming the failure is resolved, the DataNode daemon can be started using the Start option in the pull-down menu next to the service name.

Ambari - hadoop2 - Mozilla Firefox

Ambari - hadoop2 Namenode information

localhost:8080/#/main/hosts/n1/summary

Ambari hadoop2 0 ops

n1 2
Back

Summary Configs

Components

⚠ DataNode / HDFS

Stopped ▾

✓ Flume / Flume

Started ▾

✓ Ganglia Monitor / Ganglia

Started ▾

✓ RegionServer / HBase

Started ▾

✓ NodeManager / YARN

Started ▾

Clients / HBase Client , HCat Client ,
HDFS Client , Hive Client ,
MapReduce2 Client , Oozie
Client , Pig , Sqoop , Tez
Client , YARN Client ,
ZooKeeper Client

Installed ▾

Host Metrics

100%
50%

0.1

39.0 KB
19.5 KB

Summary

Hostname: n1
IP Address: 10.0.0.11
OS: centos6 (x86_64)
Cores (CPU): 4 (4)
Disk: 46.37GB/378.07GB (12.26% used)
Memory: 15.60GB
Load Avg: 0.06
Heartbeat: a moment ago

Licensed under the Apache License, Version 2.0.
See third-party tools/resources that Ambari uses and their respective authors

<https://hemanthrajhemu.github.io>

Figure 9.13 Ambari window for host n1 indicating the DataNode/HDFS service has stopped

When the DataNode daemon is restarted, a confirmation similar to [Figure 9.14](#) is required from the user.



Figure 9.14 Ambari restart confirmation

When a service daemon is started or stopped, a progress window similar to [Figure 9.15](#) is opened. The progress bar indicates the status of each action. Note that previous actions are part of this window. If something goes wrong during the action, the progress bar will turn red. If the system generates a warning about the action, the process bar will turn orange.

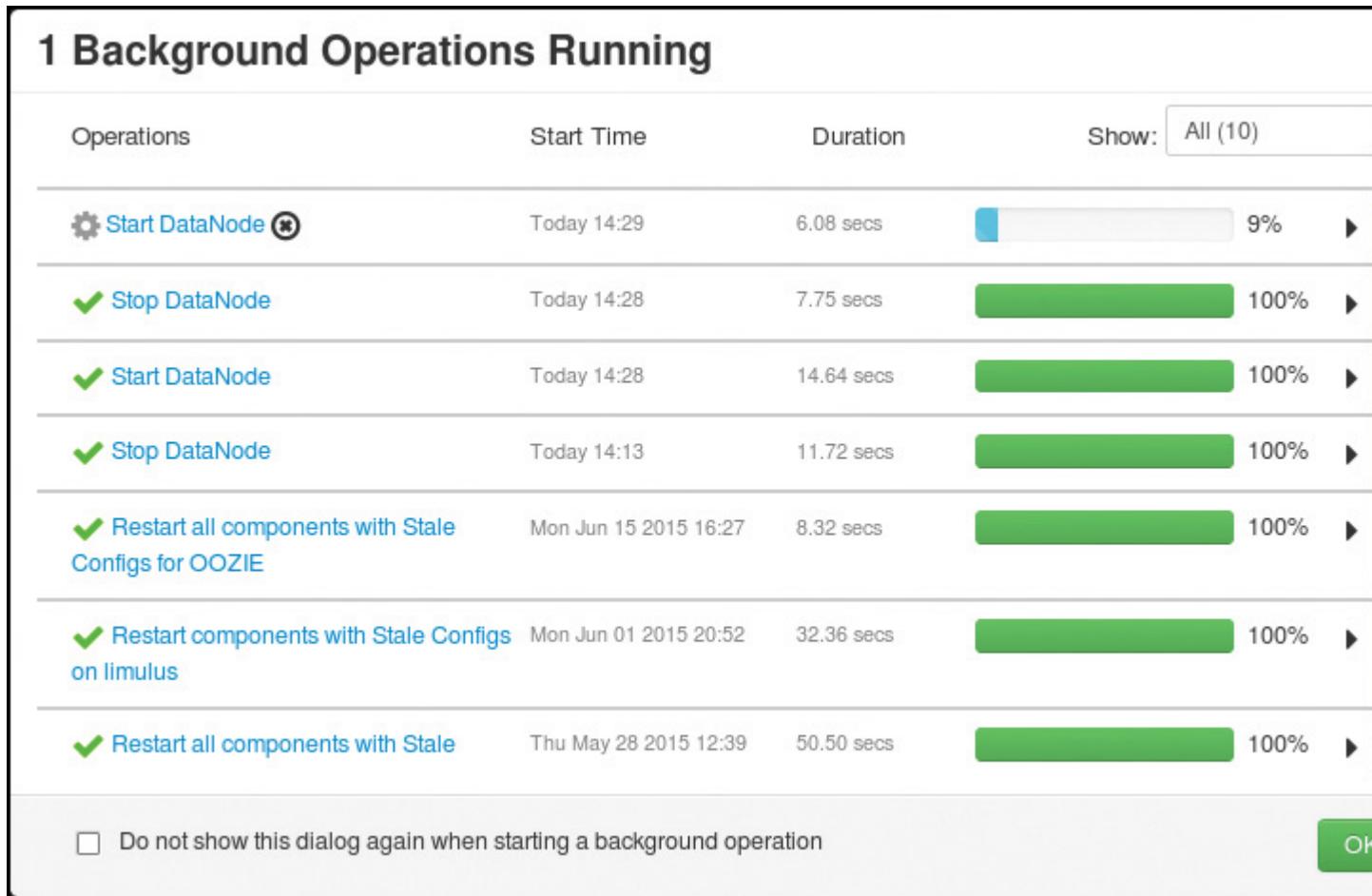


Figure 9.15 Ambari progress window for DataNode restart

When these background operations are running, the small ops (operations) bubble on the top menu bar will indicate how many operations are running. (If different service daemons are started or stopped, each process will be run to completion before the next one starts.)

Once the DataNode has been restarted successfully, the dashboard will reflect the new status (e.g., 4/4 DataNodes are Live). As shown in [Figure 9.16](#), all four DataNodes are now working and the service error indicators are beginning to slowly disappear. The service error indicators may lag behind the real-time widget updates for several minutes.

Ambari - hadoop2 - Mozilla Firefox

Ambari - hadoop2 x Namenode information x +

localhost:8080/#/main/dashboard/metrics

Ambari hadoop2 0 ops

HDFS 1
MapReduce2
YARN
Tez
Nagios
Ganglia
Hive
HBase
Pig
Sqoop
Oozie
ZooKeeper
Flume

Metrics Heatmaps Config History

HDFS Disk Usage: 18%
DataNodes Live: 4/4

CPU Usage: 100%
Cluster Load

NameNode Uptime: 52.9 d
HBase Master Heap: 8%

ResourceManager Heap: 11%
ResourceManager Uptime: 19.1 d
NodeManag

HDFS Links
Name
Secondary
4 Data

NameNode H
46

HBase Links
HBase
4 Region
Master

Licensed under the Apache License, Version 2.0.
See third-party-licenses/resources that Ambari uses and their respective authors.

<https://nemanthrajhemu.github.io>

Figure 9.16 Ambari dashboard indicating all DataNodes are running (The service error indicators will slowly drop off the screen.)

CHANGING HADOOP PROPERTIES

One of the challenges of managing a Hadoop cluster is managing changes to cluster-wide configuration properties. In addition to modifying a large number of properties, making changes to a property often requires restarting daemons (and dependent daemons) across the entire cluster. This process is tedious and time consuming. Fortunately, Ambari provides an easy way to manage this process.

As described previously, each service provides a Configs tab that opens a form displaying all the possible service properties. Any service property can be changed (or added) using this interface. As an example, the configuration properties for the YARN scheduler are shown in [Figure 9.17](#).

Ambari - hadoop2 - Mozilla Firefox

Ambari - hadoop2 x Namenode information x +

localhost:8080/#/main/services/YARN/configs

Ambari hadoop2 0 ops

Summary Configs Quick Links

Group: YARN Default (4) Manage Config Groups

V11 admin 19 days ago Current

V10 admin 19 days ago

V9 admin 19 days ago

Actions ▾

Resource Manager

ResourceManager	limulus
ResourceManager Java heap size	1024 MB
yarn.acl.enable	<input type="checkbox"/>
yarn.admin.acl	
yarn.log-aggregation-enable	<input checked="" type="checkbox"/>

Node Manager

yarn.nodemanager.resource.memory-mb	12288
yarn.nodemanager.vmem-pmem-ratio	2.1
yarn.nodemanager.log-dirs	/opt/hadoop/yarn/log,/hdfs1/hadoop/yarn/log,/local-dir

<https://hemanthrajhemu.github.io>

Figure 9.17 Ambari YARN properties view

The number of options offered depends on the service; the full range of YARN properties can be viewed by scrolling down the form. Both [Chapter 4, “Running Example Programs and Benchmarks,”](#) and [Chapter 6, “MapReduce Programming,”](#) discussed the YARN property `yarn.log-aggregation-enable`. To easily view the application logs, this property must be set to true. This property is normally on by default. As an example for our purposes here, we will use the Ambari interface to disable this feature. As shown in [Figure 9.18](#), when a property is changed, the green Save button becomes activated.

The screenshot shows the Ambari interface for managing YARN properties. At the top, there are tabs for 'V11' and 'Current' (admin authored on Thu, May 28, 2015 12:39). Below this, under the 'Resource Manager' section, there is a table with two rows:

ResourceManager	Value
ResourceManager Java heap size	1024 MB

Below the table are three property settings:

- yarn.acl.enable**: Value is off (unchecked), with a lock icon and a green plus sign icon.
- yarn.admin.acl**: Value is empty, with a lock icon and a green plus sign icon.
- yarn.log-aggregation-enable**: Value is off (unchecked), with a lock icon and a green plus sign icon.

The 'Save' button at the bottom right is green, indicating changes have been made.

Figure 9.18 YARN properties with log aggregation turned off

Changes do not become permanent until the user clicks the Save button. A save/notes window will then be displayed. It is highly recommended that historical notes concerning the change be added to this window.

Once the user adds any notes and clicks the Save button, another window, shown in [Figure 9.20](#), is presented. This window confirms that the properties have been saved.

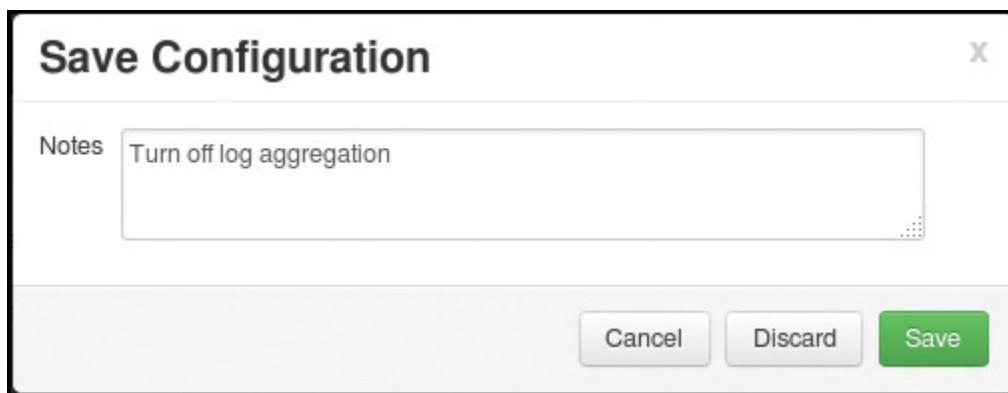


Figure 9.19 Ambari configuration save/notes window

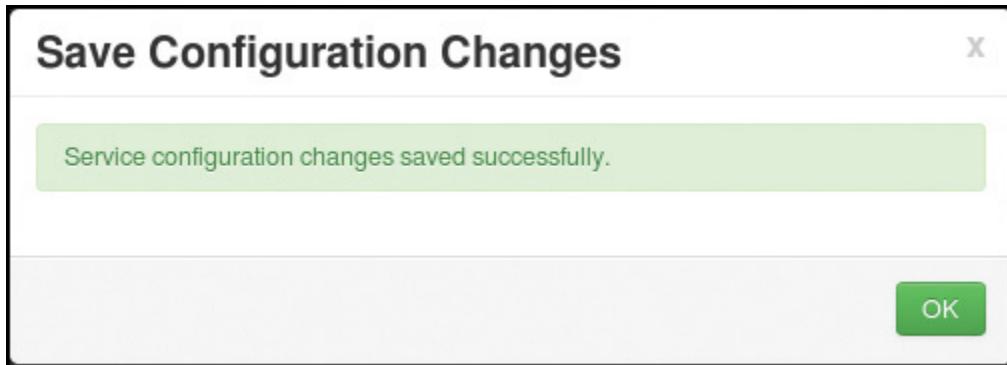


Figure 9.20 Ambari configuration change notification

Once the new property is changed, an orange Restart button will appear at the top left of the window. The new property will not take effect until the required services are restarted. As shown in [Figure 9.21](#), the Restart button provides two options: Restart All and Restart NodeManagers. To be safe, the Restart All should be used. Note that Restart All does not mean all the Hadoop services will be restarted; rather, only those that use the new property will be restarted.

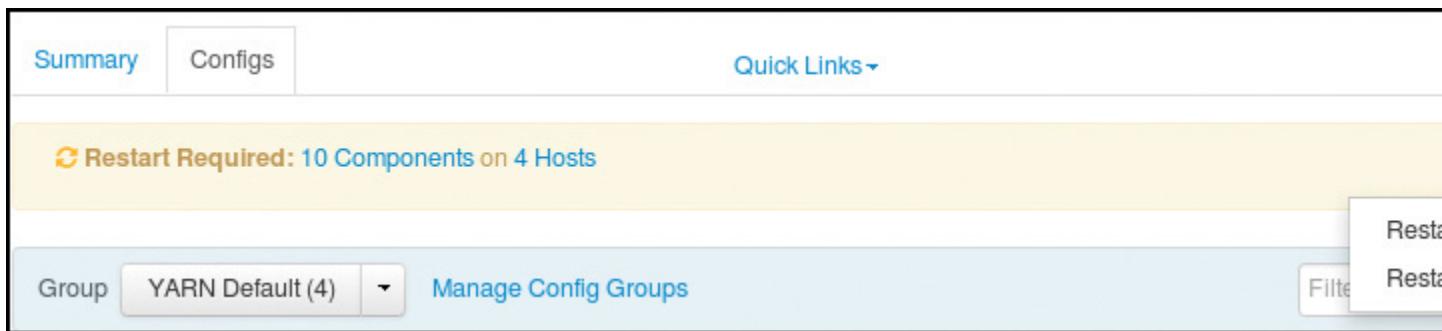


Figure 9.21 Ambari Restart function appears after changes in service properties

After the user clicks Restart All, a confirmation window, shown in [Figure 9.22](#) will be displayed. Click Confirm Restart All to begin the cluster-wide restart.

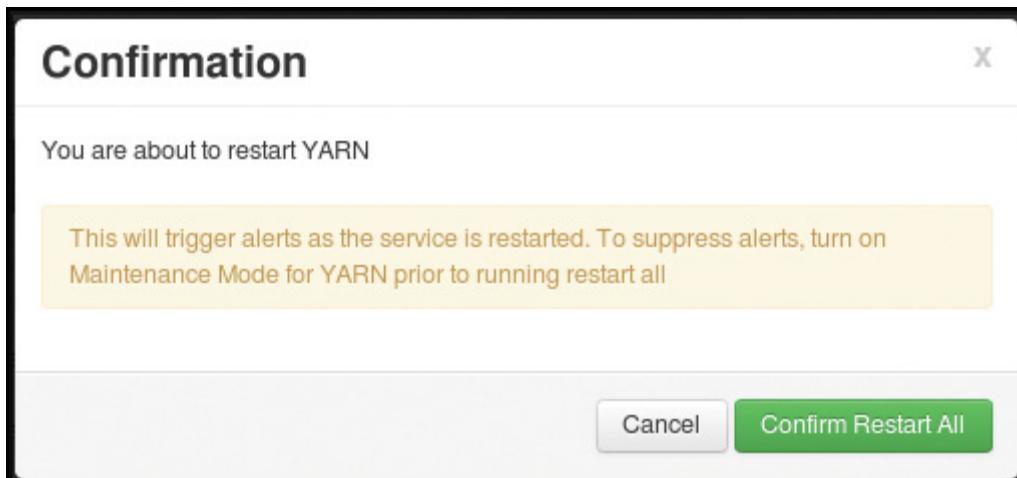


Figure 9.22 Ambari confirmation box for service restart

Similar to the DataNode restart example, a progress window will be displayed. Again, the progress bar is for the entire YARN restart. Details from the logs can be found by clicking the arrow to the right of the bar (see [Figure 9.23](#)).

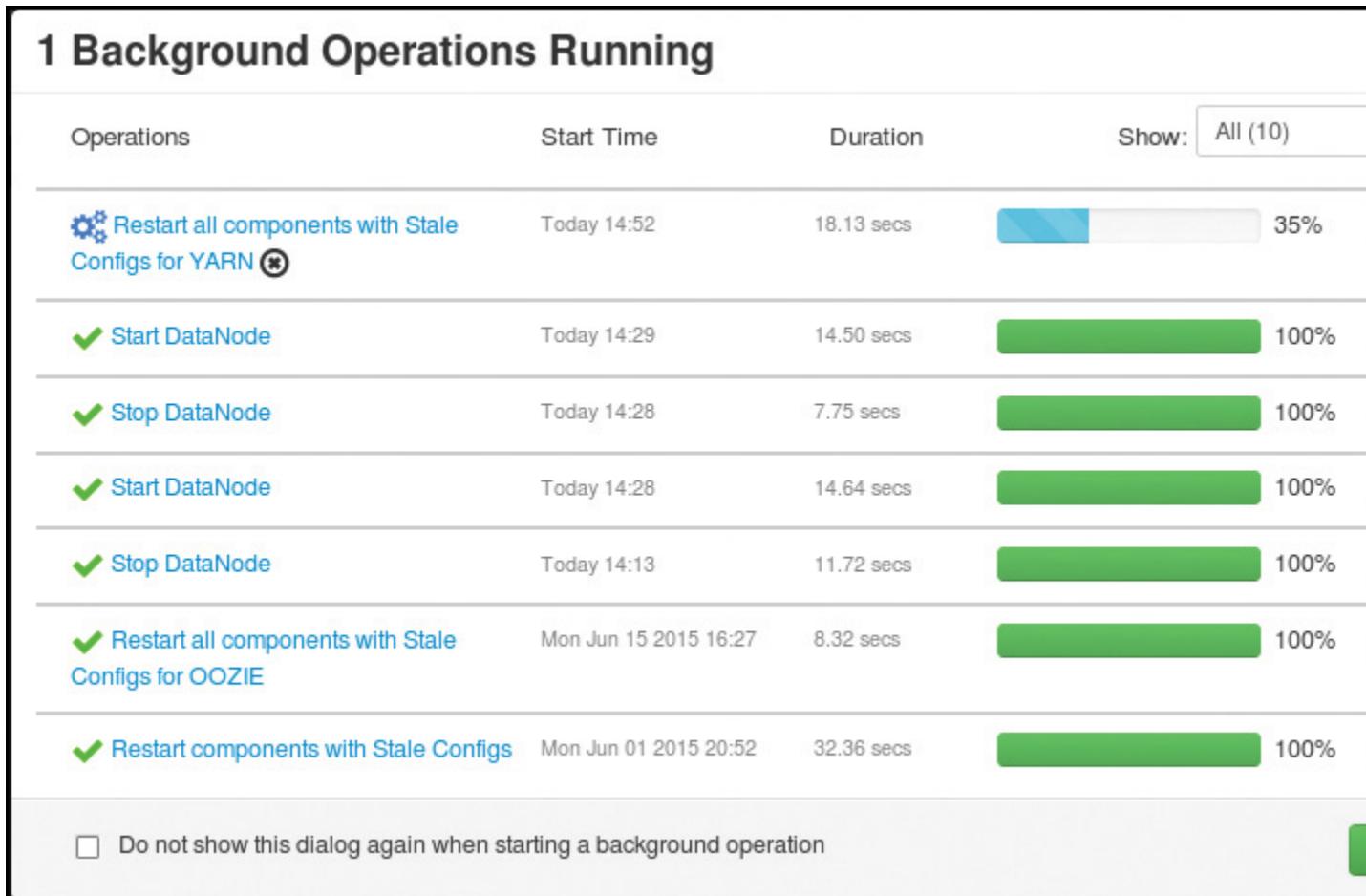


Figure 9.23 Ambari progress window for cluster-wide YARN restart

Once the restart is complete, run a simple example (see [Chapter 4](#)) and attempt to view the logs using the YARN ResourceManager Applications UI. (You can access the UI from the Quick Links pull-down menu in the middle of the YARN series window.) A message similar to that in [Figure 9.24](#) will be displayed (compare this message to the log data in [Figure 6.1](#)).

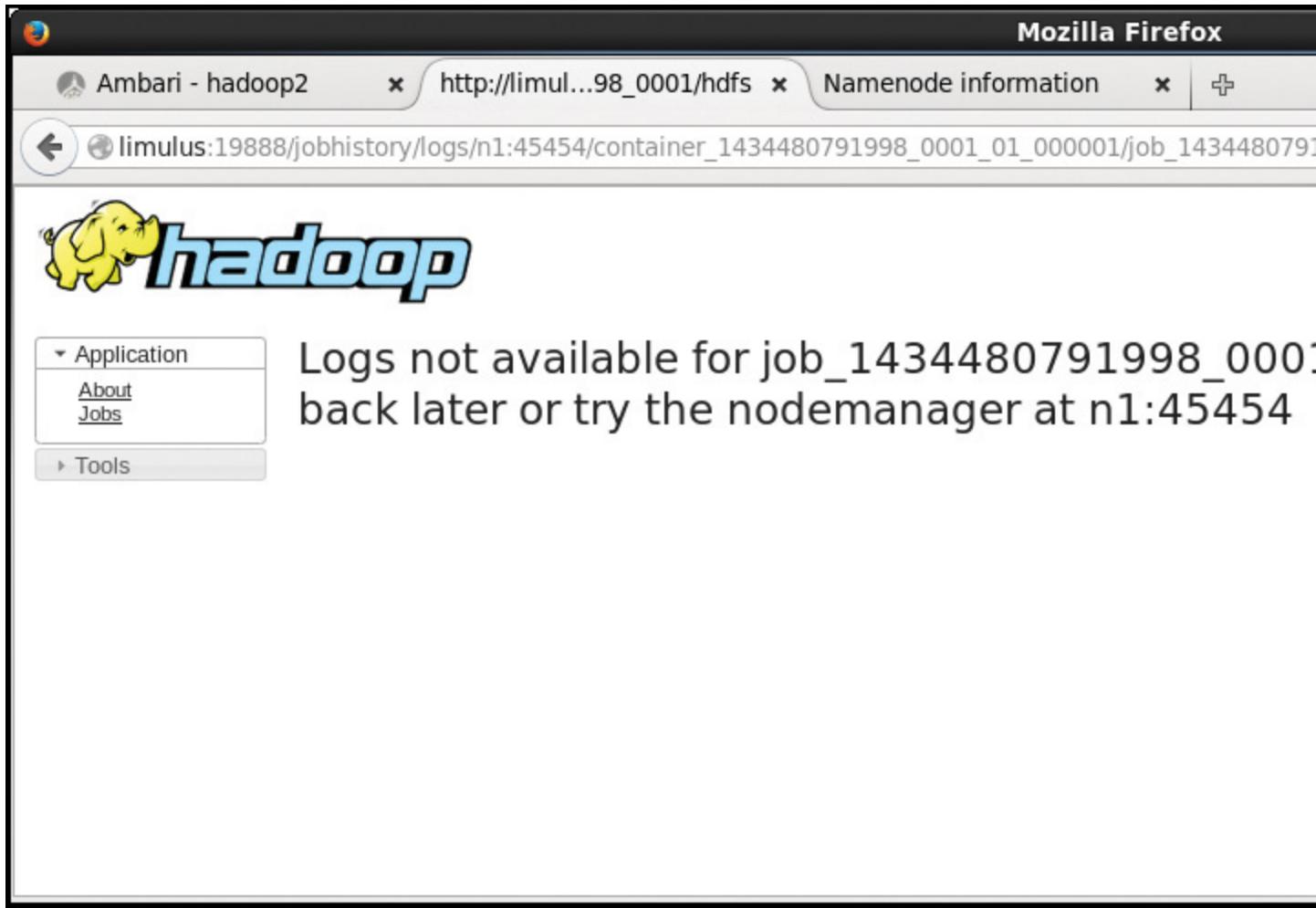


Figure 9.24 YARN ResourceManager interface with log aggregation turned off (compare to [Figure 6.1](#))

Ambari tracks all changes made to system properties. As can be seen in [Figure 9.17](#) and in more detail in [Figure 9.25](#), each time a configuration is changed, a new version is created. Reverting back to a previous version results in a new version. You can reduce the potential for version confusion by providing meaningful comments for each change (e.g., [Figure 9.19](#) and [Figure 9.27](#)). In the preceding example, we created version 12 (V12). The current version is indicated by a green Current label in the horizontal version boxes or in the dark horizontal bar. Scrolling through the version boxes or pulling down the menu on the left-hand side of the dark horizontal bar will display the previous configuration versions.

The screenshot shows the Ambari configuration management interface for the YARN service. At the top, there are tabs for 'Summary' and 'Configs', with 'Configs' being the active tab. A 'Quick Links' dropdown is also present. Below the tabs, a header bar includes a 'Group' dropdown set to 'YARN Default (4)', a 'Manage Config Groups' button, and a 'Filter...' button. A navigation bar with left and right arrows is shown above a list of configuration versions. The list includes:

- V12 admin a moment ago (Current)
- V11 admin 19 days ago
- V10 admin 19 days ago
- V9 admin 19 days ago
- V8 admin 19 days ago

Below the list, a status bar indicates 'V12 Current admin authored on Tue, Jun 16, 2015 14:49'. The main content area is titled 'Resource Manager' and contains the following configuration details:

Configuration	Value	Unit
ResourceManager	limulus	
ResourceManager Java heap size	1024	MB
yarn.acl.enable	<input type="checkbox"/>	
yarn.admin.acl	<input type="text"/>	
yarn.log-aggregation-enable	<input type="checkbox"/>	

Figure 9.25 Ambari configuration change management for YARN service (Version V12 is current)

To revert to a previous version, simply select the version from the version boxes or the pull-down menu. In [Figure 9.26](#), the user has selected the previous version by clicking the Make Current button in the information box. This configuration will return to the previous state where log aggregation is enabled.

The screenshot shows the Ambari UI for managing YARN configurations. At the top, there are tabs for 'Summary' and 'Configs'. Below them, a 'Group' dropdown is set to 'YARN Default (4)'. A list of configurations is shown, with 'V11' highlighted as the 'Current' version. Other versions listed are V12 (7 minutes ago), V11 (19 days ago), V10 (19 days ago), V9 (19 days ago), and V8 (19 days ago). Each configuration entry includes the author ('admin'), the date it was authored, and a 'View', 'Compare', or 'Make Current' button. A 'Filter...' button is also present.

A 'Resource Manager' section is expanded, showing the following configuration details:

- ResourceManager: limulus
- ResourceManager Java heap size: 1024 MB
- yarn.acl.enable:
- yarn.admin.acl: (empty input field)
- yarn.log-aggregation-enable:

Figure 9.26 Reverting to previous YARN configuration (V11) with Ambari

As shown in Figure 9.27, a confirmation/notes window will open before the new configuration is saved. Again, it is suggested that you provide notes about the change in the Notes text box. When the save step is complete, the Make Current button will restore the previous configuration. The orange Restart button will appear and indicate that a service restart is needed before the changes take effect.

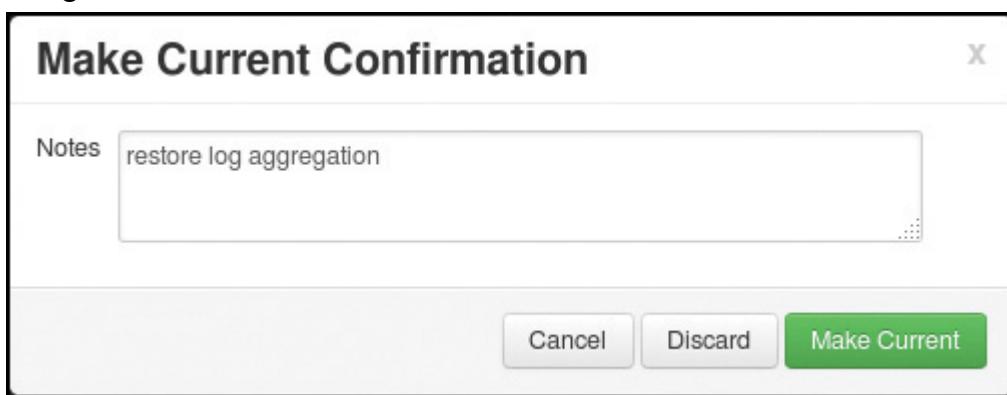


Figure 9.27 Ambari confirmation window for a new configuration

There are several important points to remember about the Ambari versioning tool:

- Every time you change the configuration, a new version is created. Reverting to a previous version creates a new version.
- You can view or compare a version to other versions without having to change or restart services. (See the buttons in the V11 box in [Figure 9.26](#).)
- Each service has its own version record.
- Every time you change the properties, you must restart the service by using the Restart button. When in doubt, restart all services.

SUMMARY AND ADDITIONAL RESOURCES

Apache Ambari provides a single control panel for the entire Hadoop cluster. The Ambari dashboard provides a quick overview of the cluster. Each service has its own summary/status and configuration window. A host's window provides detailed metrics for each host, along with the ability to manage the services running on a particular host.

Other services, including adding users and groups, are part of the Ambari administration features. Each Hadoop service is monitored and all issues are reported via email and displayed on the interface for easy identification. All services can be stopped and started directly from the Ambari interface. In addition, changes to service properties are accomplished by using a simplified form with contextual help. Ambari also provides a versioning system that allows previous cluster configurations to be easily restored.

Apache Ambari is an open source tool that is gaining new features with each release. Currently, it is an efficient and useful tool for all Hadoop 2 clusters. To learn more about Ambari, consult the project website: <https://ambari.apache.org>.

10. Basic Hadoop Administration Procedures

In This Chapter:

- Several basic Hadoop YARN administration topics are presented, including decommissioning YARN nodes, managing YARN applications, and important YARN properties.
- Basic HDFS administration procedures are described, including using the NameNode UI, adding users, performing file system checks, balancing DataNodes, taking HDFS snapshots, and using the HDFS NFSv3 gateway.
- The Capacity scheduler is discussed.
- Hadoop version 2 MapReduce compatibility and node capacity are discussed.

In Chapter 9, “[Managing Hadoop with Apache Ambari](#),” the Apache Ambari web management tool was described in detail. Much of the day-to-day management of a Hadoop cluster can be accomplished using the Ambari interface. Indeed, whenever possible, Ambari should be used to manage the cluster because it keeps track of the cluster state.

Hadoop has two main areas of administration: the YARN resource manager and the HDFS file system. Other application frameworks (e.g., the MapReduce framework) and tools have their own management files. As mentioned in Chapter 2, “[Installation Recipes](#),” Hadoop configuration is accomplished through the use of XML configuration files. The basic files and their function are as follows:

- `core-default.xml`: System-wide properties
- `hdfs-default.xml`: Hadoop Distributed File System properties
- `mapred-default.xml`: Properties for the YARN MapReduce framework
- `yarn-default.xml`: YARN properties

You can find a complete list of properties for all these files

at <http://hadoop.apache.org/docs/current/> (look at the lower-left side of the page under “Configuration”). A full discussion of all options is beyond the scope of this book. The Apache Hadoop documentation does provide helpful comments and defaults for each property.

If you are using Ambari, you should manage the configuration files through the interface, rather than editing them by hand. If some other management tool is employed, it should be used as described. Installation of Hadoop by hand (as in a pseudo-distributed mode in Chapter 2) requires that you edit the configuration files by hand and then copy them to all nodes in the cluster if applicable.

The following sections cover some useful administration tasks that may fall outside of Ambari, require special configuration, or require more explanation. By no means does this discussion cover all possible topics related to Hadoop administration; rather, it is designed to help jump-start your administration of Hadoop version 2.

BASIC HADOOP YARN ADMINISTRATION

YARN has several built-in administrative features and commands. To find out more about them, examine the YARN commands documentation

at <https://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn->

[site/YarnCommands.html#Administration_Commands](#). The main administration command is `yarn rmadmin` ([resource manager administration](#)). Enter `yarn rmadmin -help` to learn more about the various options.

Decommissioning YARN Nodes

If a NodeManager host/node needs to be removed from the cluster, it should be decommissioned first. Assuming the node is responding, you can easily decommission it from the Ambari web UI. Simply go to the Hosts view, click on the host, and select Decommission from the pull-down menu next to the NodeManager component. Note that the host may also be acting as a HDFS DataNode. Use the Ambari Hosts view to decommission the HDFS host in a similar fashion.

YARN WebProxy

The Web Application Proxy is a separate proxy server in YARN that addresses security issues with the cluster web interface on ApplicationMasters. By default, the proxy runs as part of the Resource Manager itself, but it can be configured to run in a stand-alone mode by adding the configuration property `yarn.web-proxy.address` to `yarn-site.xml`. (Using Ambari, go to the YARN Configs view, scroll to the bottom, and select Custom `yarn-site.xml`/Add property.) In stand-alone mode, `yarn.web-proxy.principal` and `yarn.web-proxy.keytab` control the Kerberos principal name and the corresponding keytab, respectively, for use in secure mode. These elements can be added to the `yarn-site.xml` if required.

Using the JobHistoryServer

The removal of the JobTracker and migration of MapReduce from a system to an application-level framework necessitated creation of a place to store MapReduce job history. The JobHistoryServer provides all YARN MapReduce applications with a central location in which to aggregate completed jobs for historical reference and debugging. The settings for the JobHistoryServer can be found in the `mapred-site.xml` file.

Managing YARN Jobs

YARN jobs can be managed using the `yarn application` command. The following options, including `-kill`, `-list`, and `-status`, are available to the administrator with this command. MapReduce jobs can also be controlled with the `mapred job` command.

[Click here to view code image](#)

usage: application

`-appTypes <Comma-separated list of application types>` Works with

--list to filter applications based on
their type.

`-help` Displays help for all commands.

`-kill <Application ID>` Kills the application.

`-list` Lists applications from the RM. Supports optional
use of `-appTypes` to filter applications based
on application type.

`-status <Application ID>` Prints the status of the application.

Neither the YARN ResourceManager UI nor the Ambari UI can be used to kill YARN applications. If a job needs to be killed, give the `yarn application` command to find the Application ID and then use the `-kill` argument.

Setting Container Memory

YARN manages application resource containers over the entire cluster. Controlling the amount of container memory takes place through three important values in the `yarn-site.xml` file:

- `yarn.nodemanager.resource.memory-mb` is the amount of memory the NodeManager can use for containers.
- `scheduler.minimum-allocation-mb` is the smallest container allowed by the ResourceManager. A requested container smaller than this value will result in an allocated container of this size (default 1024MB).
- `yarn.scheduler.maximum-allocation-mb` is the largest container allowed by the ResourceManager (default 8192MB).

Setting Container Cores

You can set the number of cores for containers using the following properties in the `yarn-site.xml`:

- `yarn.scheduler.minimum-allocation-vcores`: The minimum allocation for every container request at the ResourceManager, in terms of virtual CPU cores. Requests smaller than this allocation will not take effect, and the specified value will be allocated the minimum number of cores. The default is 1 core.
- `yarn.scheduler.maximum-allocation-vcores`: The maximum allocation for every container request at the ResourceManager, in terms of virtual CPU cores. Requests larger than this allocation will not take effect, and the number of cores will be capped at this value. The default is 32.
- `yarn.nodemanager.resource.cpu-vcores`: The number of CPU cores that can be allocated for containers. The default is 8.

Setting MapReduce Properties

As noted throughout this book, MapReduce now runs as a YARN application. Consequently, it may be necessary to adjust some of the `mapred-site.xml` properties as they relate to the map and reduce containers. The following properties are used to set some Java arguments and memory size for both the map and reduce containers:

- `mapred.child.java.opts` provides a larger or smaller heap size for child JVMs of maps (e.g., `--Xmx2048m`).
- `mapreduce.map.memory.mb` provides a larger or smaller resource limit for maps (default = 1536MB).
- `mapreduce.reduce.memory.mb` provides a larger heap size for child JVMs of maps (default = 3072MB).
- `mapreduce.reduce.java.opts` provides a larger or smaller heap size for child reducers.

BASIC HDFS ADMINISTRATION

The following section covers some basic administration aspects of HDFS. Advanced topics such as HDFS Federation or high availability are beyond the scope of this chapter. More information on these and other HDFS topics can be found at <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsUserGuide.html>. (Note the HDFS topic menu available on the left-hand side of this webpage.)

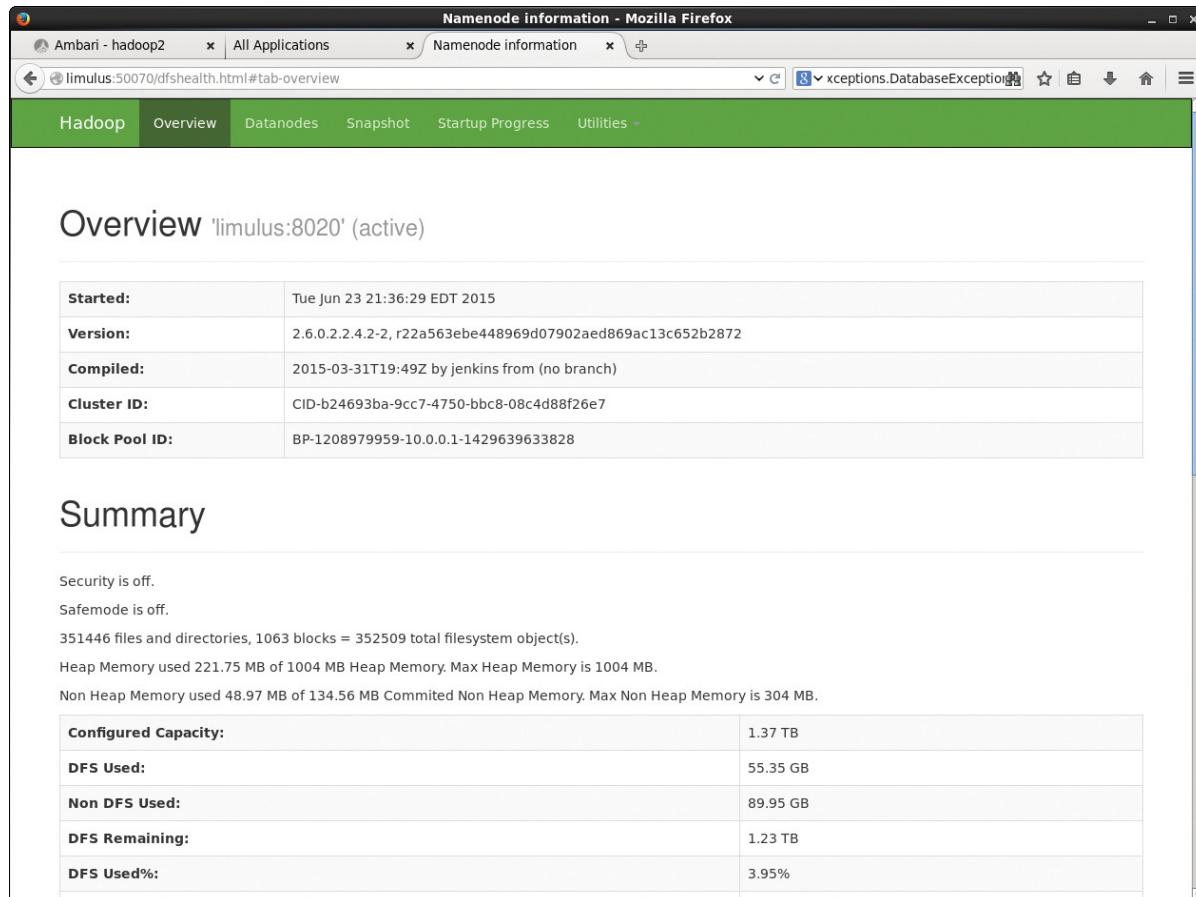
The NameNode User Interface

Monitoring HDFS can be done in several ways. One of the more convenient ways to get a quick view of HDFS status is through the NameNode user interface. This web-based tool provides essential information about HDFS and offers the capability to browse the HDFS namespace and logs.

The web-based UI can be started from within Ambari or from a web browser connected to the NameNode. In Ambari, simply select the HDFS service window and click on the Quick Links pull-down menu in the top middle of the page. Select NameNode UI. A new browser tab will open with the UI shown in [Figure 10.1](#). You can also start the UI directly by entering the following command (the command given here assumes the Firefox browser is used, but other browsers should work as well):

[**Click here to view code image**](#)

```
$ firefox http://localhost:50070
```



The screenshot shows a Mozilla Firefox browser window titled "Namenode information - Mozilla Firefox". The address bar displays "Ambari - hadoop2" and "limulus:50070/dfshealth.html#tab-overview". The main content area is a green header bar with tabs: "Hadoop" (selected), "Overview", "Datanodes", "Snapshot", "Startup Progress", and "Utilities". Below the header is a section titled "Overview 'limulus:8020' (active)". It contains a table with the following data:

Started:	Tue Jun 23 21:36:29 EDT 2015
Version:	2.6.0.2.4.2-2, r22a563ebe448969d07902aed869ac13c652b2872
Compiled:	2015-03-31T19:49Z by jenkins from (no branch)
Cluster ID:	CID-b24693ba-9cc7-4750-bbc8-08c4d88f26e7
Block Pool ID:	BP-1208979959-10.0.0.1-1429639633828

Below this is a "Summary" section with the following text:

Security is off.
Safemode is off.
351446 files and directories, 1063 blocks = 352509 total filesystem object(s).
Heap Memory used 221.75 MB of 1004 MB Heap Memory. Max Heap Memory is 1004 MB.
Non Heap Memory used 48.97 MB of 134.56 MB Committed Non Heap Memory. Max Non Heap Memory is 304 MB.

Configured Capacity:	1.37 TB
DFS Used:	55.35 GB
Non DFS Used:	89.95 GB
DFS Remaining:	1.23 TB
DFS Used%:	3.95%

Figure 10.1 Overview page for NameNode user interface

There are five tabs on the UI: Overview, Datanodes, Snapshot, Startup Progress, and Utilities. The Overview page provides much of the essential information that the command-line tools also offer, but in a much easier-to-read format. The Datanodes tab displays node information like that shown in [Figure 10.2](#).

The screenshot shows a Mozilla Firefox browser window titled "Namenode information - Mozilla Firefox". The address bar shows "Ambari - hadoop2" and "Namennode information". The main content area has a green header bar with tabs: Hadoop, Overview, Datanodes, Snapshot, Startup Progress, and Utilities. The "Datanodes" tab is selected. Below the header, the title "Datanode Information" is displayed. A section titled "In operation" contains a table with the following data:

Node	Last contact	Admin State	Capacity	Used	Non DFS Used	Remaining	Blocks	Block pool used	Failed Volumes	Version
n1 (10.0.0.11:50010)	1	In Service	350.1 GB	7.87 GB	24.48 GB	317.75 GB	352	7.87 GB (2.25%)	0	2.6.0.2.2.4.2-2
n0 (10.0.0.10:50010)	1	In Service	350.1 GB	11.55 GB	24.52 GB	314.03 GB	174	11.55 GB (3.3%)	0	2.6.0.2.2.4.2-2
limulus (10.0.0.1:50010)	1	In Service	349.85 GB	19.06 GB	16.47 GB	314.31 GB	956	19.06 GB (5.45%)	0	2.6.0.2.2.4.2-2
n2 (10.0.0.12:50010)	1	In Service	350.1 GB	16.86 GB	24.48 GB	308.77 GB	634	16.86 GB (4.81%)	0	2.6.0.2.2.4.2-2

A section titled "Decommissioning" is present but empty. At the bottom left, it says "Hadoop, 2014." and at the bottom right, it says "Legacy UI".

Figure 10.2 NameNode web interface showing status of DataNodes

The Snapshot window (shown later in this chapter in [Figure 10.5](#)) lists the “snapshottable” directories and the snapshots. Further information on snapshots can be found in the “[HDFS Snapshots](#)” section.

[Figure 10.3](#) provides a NameNode startup progress view. As noted in [Chapter 3](#), “[Hadoop Distributed File System Basics](#),” when the NameNode starts, it reads the previous file system image file (`fsimage`); applies any new edits to the file system image, thereby creating a new file system image; and drops into safe mode until enough DataNodes come online. This progress is shown in real time in the UI as the NameNode starts. Completed phases are displayed in bold text. The currently running phase is displayed in italicics. Phases that have not yet begun are displayed in gray text. In [Figure 10.3](#), all the phases have been completed, and as indicated in the overview window in [Figure 10.1](#), the system is out of safe mode.

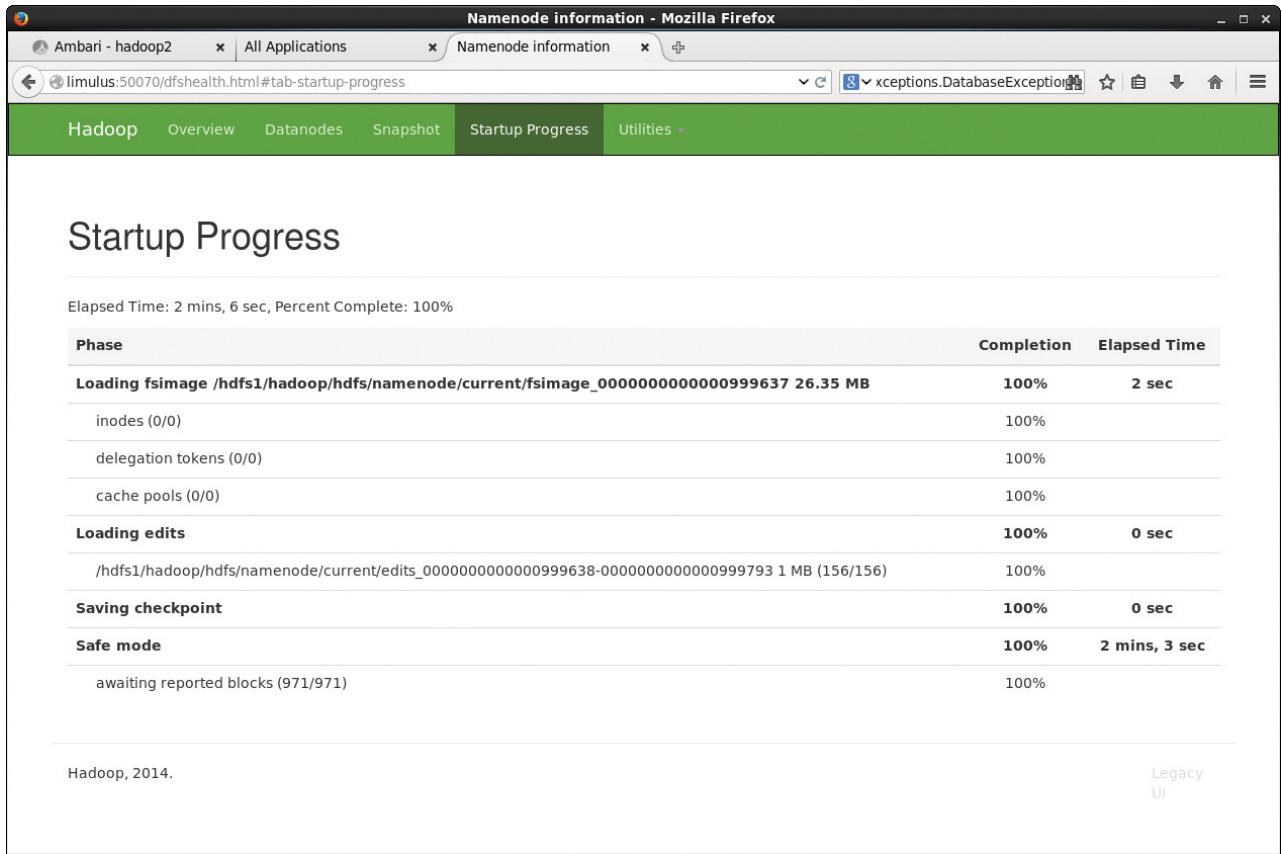


Figure 10.3 NameNode web interface showing startup progress

The Utilities menu offers two options. The first, as shown in [Figure 10.4](#), is a file system browser. From this window, you can easily explore the HDFS namespace. The second option, which is not shown, links to the various NameNode logs.

Permission	Owner	Group	Size	Replication	Block Size	Name
drwxrwxrwx	yarn	hadoop	0 B	0	0 B	app-logs
drwxr-xr-x	hdfs	hdfs	0 B	0	0 B	apps
drwxr-xr-x	hdfs	hdfs	0 B	0	0 B	benchmarks
drwxr-xr-x	hdfs	hdfs	0 B	0	0 B	hdp
drwxr-xr-x	mapred	hdfs	0 B	0	0 B	mapred
drwxr-xr-x	hdfs	hdfs	0 B	0	0 B	mr-history
drwxr-xr-x	hdfs	hdfs	0 B	0	0 B	system
drwxrwxrwx	hdfs	hdfs	0 B	0	0 B	tmp
drwxr-xr-x	hdfs	hdfs	0 B	0	0 B	user
drwx-wx-wx	hdfs	hdfs	0 B	0	0 B	var

Hadoop, 2014.

Figure 10.4 NameNode web interface directory browser

Adding Users to HDFS

For a full explanation of HDFS permissions, see the following document: <http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsPermissionsGuide.html>. Keep in mind that errors that crop up while Hadoop applications are running are often due to file permissions.

To quickly create user accounts manually on a Linux-based system, perform the following steps:

1. Add the user to the group for your operating system on the HDFS client system. In most cases, the `groupname` should be that of the HDFS superuser, which is often `hadoop` or `hdfs`.

[Click here to view code image](#)

```
useradd -G <groupname> <username>
```

2. Create the `username` directory in HDFS.

[Click here to view code image](#)

```
hdfs dfs -mkdir /user/<username>
```

3. Give that account ownership over its directory in HDFS.

[Click here to view code image](#)

```
hdfs dfs -chown <username>:<groupname> /user/<username>
```

Perform an FSCK on HDFS

To check the health of HDFS, you can issue the `hdfs fsck <path>` (file system check) command. The entire HDFS namespace can be checked, or a subdirectory can be entered as an argument to the command. The following example checks the entire HDFS namespace.

[Click here to view code image](#)

```
$ hdfs fsck /
```

```
Connecting to namenode via http://limulus:50070
FSCK started by hdfs (auth:SIMPLE) from /10.0.0.1 for path / at Fri May 29
14:48:01 EDT 2015
.....
.....
Status: HEALTHY
Total size: 100433565781 B (Total open files size: 498 B)
Total dirs: 201331
Total files: 1003
Total symlinks: 0 (Files currently being written: 6)
Total blocks (validated): 1735 (avg. block size 57886781 B) (Total open file
blocks (not validated): 6)
Minimally replicated blocks: 1735 (100.0 %)
Over-replicated blocks: 0 (0.0 %)
Under-replicated blocks: 0 (0.0 %)
Mis-replicated blocks: 0 (0.0 %)
Default replication factor: 2
Average block replication: 1.7850144
Corrupt blocks: 0
Missing replicas: 0 (0.0 %)
Number of data-nodes: 4
Number of racks: 1
FSCK ended at Fri May 29 14:48:03 EDT 2015 in 1853 milliseconds
```

The filesystem under path '/' is HEALTHY

Other options provide more detail, include snapshots and open files, and management of corrupted files.

- `-move` moves corrupted files to `/lost+found`.
- `-delete` deletes corrupted files.
- `-files` prints out files being checked.
- `-openforwrite` prints out files opened for writes during check.
- `-includesnapshots` includes snapshot data. The path indicates the existence of a snapshottable directory or the presence of snapshottable directories under it.
- `-list-corruptfileblocks` prints out a list of missing blocks and the files to which they belong.

- `-blocks` prints out a block report.
- `-locations` prints out locations for every block.
- `-racks` prints out network topology for data-node locations.

Balancing HDFS

Based on usage patterns and DataNode availability, the number of data blocks across the DataNodes may become unbalanced. To avoid over-utilized DataNodes, the HDFS balancer tool rebalances data blocks across the available DataNodes. Data blocks are moved from over-utilized to under-utilized nodes to within a certain percent threshold. Rebalancing can be done when new DataNodes are added or when a DataNode is removed from service. This step does not create more space in HDFS, but rather improves efficiency.

The HDFS superuser must run the balancer. The simplest way to run the balancer is to enter the following command:

```
$ hdfs balancer
```

By default, the balancer will continue to rebalance the nodes until the number of data blocks on all DataNodes are within 10% of each other. The balancer can be stopped, without harming HDFS, at any time by entering a Ctrl-C. Lower or higher thresholds can be set using the `-threshold` argument. For example, giving the following command sets a 5% threshold:

```
$ hdfs balancer -threshold 5
```

The lower the threshold, the longer the balancer will run. To ensure the balancer does not swamp the cluster networks, you can set a bandwidth limit before running the balancer, as follows:

[Click here to view code image](#)

```
$ dfsadmin -setBalancerBandwidth newbandwidth
```

The `newbandwidth` option is the maximum amount of network bandwidth, in bytes per second, that each DataNode can use during the balancing operation.

Balancing data blocks can also break HBase locality. When HBase regions are moved, some data locality is lost, and the RegionServers will then request the data over the network from remote DataNode(s). This condition will persist until a major HBase compaction event takes place (which may either occur at regular intervals or be initiated by the administrator).

See <http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HDFSCCommands.html> for more information on balancer options.

HDFS Safe Mode

As mentioned in [Chapter 3, “Hadoop Distributed File System Basics,”](#) when the NameNode starts, it loads the file system state from the `fsimage` and then applies the edits log file. It then waits for DataNodes to report their blocks. During this time, the NameNode stays in a read-only Safe Mode. The NameNode leaves Safe Mode automatically after the DataNodes have reported that most file system blocks are available.

The administrator can place HDFS in Safe Mode by giving the following command:

[Click here to view code image](#)

```
$ hdfs dfsadmin -safemode enter
```

Entering the following command turns off Safe Mode:

[Click here to view code image](#)

```
$ hdfs dfsadmin -safemode leave
```

HDFS may drop into Safe Mode if a major issue arises within the file system (e.g., a full DataNode). The file system will not leave Safe Mode until the situation is resolved. To check whether HDFS is in Safe Mode, enter the following command:

[Click here to view code image](#)

```
$ hdfs dfsadmin -safemode get
```

Decommissioning HDFS Nodes

If you need to remove a DataNode host/node from the cluster, you should decommission it first. Assuming the node is responding, it can be easily decommissioned from the Ambari web UI. Simply go to the Hosts view, click on the host, and select Decommission from the pull-down menu next to the DataNode component. Note that the host may also be acting as a Yarn NodeManager. Use the Ambari Hosts view to decommission the YARN host in a similar fashion.

SecondaryNameNode

To avoid long NameNode restarts and other issues, the performance of the SecondaryNameNode should be verified. Recall that the SecondaryNameNode takes the previous file system image file (`fsimage*`) and adds the NameNode file system edits to create a new file system image file for the NameNode to use when it restarts. The `hdfs-site.xml` defines a property called `fs.checkpoint.period` (called HDFS Maximum Checkpoint Delay in Ambari). This property provides the time in seconds between the SecondaryNameNode checkpoints.

When a checkpoint occurs, a new `fsimage*` file is created in the directory corresponding to the value of `dfs.namenode.checkpoint.dir` in the `hdfs-site.xml` file. This file is also placed in the NameNode directory corresponding to the `dfs.namenode.name.dir` path designated in the `hdfs-site.xml` file. To test the checkpoint process, a short time period (e.g., 300 seconds) can be used for `fs.checkpoint.period` and HDFS restarted. After five minutes, two identical `fsimage*` files should be present in each of the two previously mentioned directories. If these files are not recent or are missing, consult the NameNode and SecondaryNameNode logs. Once the SecondaryNameNode process is confirmed to be working correctly, reset the `fs.checkpoint.period` to the previous value and restart HDFS. (Ambari versioning is helpful with this type of procedure.) If the SecondaryNameNode is not running, a checkpoint can be forced by running the following command:

[Click here to view code image](#)

```
$ hdfs secondarynamenode -checkpoint force
```

HDFS Snapshots

HDFS snapshots are read-only, point-in-time copies of HDFS. Snapshots can be taken on a subtree of the file system or the entire file system. Some common use-cases for snapshots are data backup, protection against user errors, and disaster recovery.

Snapshots can be taken on any directory once the directory has been set as **snapshotable**. A snapshotable directory is able to accommodate 65,536 simultaneous snapshots. There is no limit on the number of snapshotable directories. Administrators may set any directory to be snapshotable, but nested snapshotable directories are not allowed. For example, a directory cannot be set to snapshotable if one of its ancestors/descendants is a snapshotable directory. More details can be found at <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsSnapshots.html>.

The following example walks through the procedure for creating a snapshot. The first step is to declare a directory as “snapshotable” using the following command:

[Click here to view code image](#)

```
$ hdfs dfsadmin -allowSnapshot /user/hdfs/war-and-peace-input  
Allowing snapshot on /user/hdfs/war-and-peace-input succeeded
```

Once the directory has been made snapshotable, the snapshot can be taken with the following command. The command requires the directory path and a name for the snapshot—in this case, wapi-snap-1.

[Click here to view code image](#)

```
$ hdfs dfs -createSnapshot /user/hdfs/war-and-peace-input wapi-snap-1  
Created snapshot /user/hdfs/war-and-peace-input/.snapshot/wapi-snap-1
```

The path of the snapshot is /user/hdfs/war-and-peace-input/.snapshot/wapi-snap-1. The /user/hdfs/war-and-peace-input directory has one file, as shown by issuing the following command:

[Click here to view code image](#)

```
$ hdfs dfs -ls /user/hdfs/war-and-peace-input/  
Found 1 items  
-rw-r--r-- 2 hdfs hdfs 3288746 2015-06-24 19:56 /user/hdfs/war-and-peace-  
input/war-and-peace.txt
```

If the file is deleted, it can be restored from the snapshot:

[Click here to view code image](#)

```
$ hdfs dfs -rm -skipTrash /user/hdfs/war-and-peace-input/war-and-peace.txt  
Deleted /user/hdfs/war-and-peace-input/war-and-peace.txt
```

```
$ hdfs dfs -ls /user/hdfs/war-and-peace-input/
```

The restoration process is basically a simple copy from the snapshot to the previous directory (or anywhere else). Note the use of the `~/.snapshot/wapi-snap-1` path to restore the file:

[**Click here to view code image**](#)

```
$ hdfs dfs -cp /user/hdfs/war-and-peace-input/.snapshot/wapi-snap-1/war-and-peace.txt /user/hdfs/war-and-peace-input
```

Confirmation that the file has been restored can be obtained by issuing the following command:

[**Click here to view code image**](#)

```
$ hdfs dfs -ls /user/hdfs/war-and-peace-input/  
Found 1 items  
-rw-r--r-- 2 hdfs hdfs 3288746 2015-06-24 21:12 /user/hdfs/war-and-peace-input/war-and-peace.txt
```

The NameNode UI provides a listing of snapshottable directories and the snapshots that have been taken. [Figure 10.5](#) shows the results of creating the previous snapshot.

The screenshot shows a Mozilla Firefox browser window with the title "Namenode information - Mozilla Firefox". The address bar displays "Ambari - hadoop2" and "limulus:50070/dfshealth.html#tab-snapshot". The page content is titled "Snapshot Summary". It shows "Snapshotable directories: 1" with a table:

Path	Snapshot Number	Snapshot Quota	Max Size
/user/hdfs/war-and-peace-input	1	65536	6,553,600

It also shows "Snapshotted directories: 1" with a table:

Snapshot ID	Snapshot Directory
wapi-snap-1	/user/hdfs/war-and-peace-input/.snapshot/wapi-snap-1

A footer note says "Hadoop, 2014."

Figure 10.5 Apache NameNode web interface showing snapshot information

To delete a snapshot, give the following command:

[Click here to view code image](#)

```
$ hdfs dfs -deleteSnapshot /user/hdfs/war-and-peace-input wapi-snap-1
```

To make a directory “*un-snapshottable*” (or go back to the default state), use the following command:

[Click here to view code image](#)

```
$ hdfs dfsadmin -disallowSnapshot /user/hdfs/war-and-peace-input  
Disallowing snapshot on /user/hdfs/war-and-peace-input succeeded
```

Configuring an NFSv3 Gateway to HDFS

HDFS supports an NFS version 3 (NFSv3) gateway. This feature enables files to be easily moved between HDFS and client systems. The NFS gateway supports NFSv3 and allows HDFS to be mounted as part of the client’s local file system. Currently the NFSv3 gateway supports the following capabilities:

- Users can browse the HDFS file system through their local file system using an NFSv3 client-compatible operating system.
- Users can download files from the HDFS file system to their local file system.
- Users can upload files from their local file system directly to the HDFS file system.
- Users can stream data directly to HDFS through the mount point. File append is supported, but random write is *not* supported.

The gateway must be run on the same host as a DataNode, NameNode, or any HDFS client.

More information about the NFSv3 gateway can be found

at <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsNfsGateway.html>.

In the following example, a simple four-node cluster is used to demonstrate the steps for enabling the NFSv3 gateway. Other potential options, including those related to security, are not addressed in this example. A DataNode is used as the gateway node in this example, and HDFS is mounted on the main (login) cluster node.

Step 1: Set Configuration Files

Several Hadoop configuration files need to be changed. In this example, the Ambari GUI will be used to alter the HDFS configuration files. See [Chapter 9](#) for information on using Ambari. Do not save the changes or restart HDFS until all the following changes are made. If you are not using Ambari, you must change these files by hand and then restart the appropriate services across the cluster. The following environment is assumed:

- OS: Linux
- Platform: RHEL 6.6
- Hortonworks HDP 2.2 with Hadoop version: 2.6

Several properties need to be added to the `/etc/hadoop/config/core-site.xml` file. Using Ambari, go to the HDFS service window and select the Configs tab. Toward the bottom of the

screen, select the Add Property link in the Custom core-site.xml section. Add the following two properties (the item used for the key field in Ambari is the `name` field included in this code):

[**Click here to view code image**](#)

```
<property>
  <name>hadoop.proxyuser.root.groups</name>
  <value>*</value>
</property>

<property>
  <name>hadoop.proxyuser.root.hosts</name>
  <value>*</value>
</property>
```

The name of the user who will start the Hadoop NFSv3 gateway is placed in the `name` field. In the previous example, `root` is used for this purpose. This setting can be any user who starts the gateway. If, for instance, user `nfsadmin` starts the gateway, then the two names would be `hadoop.proxyuser.nfsadmin.groups` and `hadoop.proxyuser.nfsadmin.hosts`. The `*` value, entered in the preceding lines, opens the gateway to all groups and allows it to run on any host. Access is restricted by entering groups (comma separated) in the group's property. Entering a host name for the host's property can restrict the host running the gateway.

Next, move to the Advanced hdfs-site.xml section and set the following property:

[**Click here to view code image**](#)

```
<property>
  <name>dfs.namenode.accesstime.precision</name>
  <value>3600000</value>
</property>
```

This property ensures client mounts with access time updates work properly. (See the `mount default atime` option.)

Finally, move to the Custom hdfs-site section, click the Add Property link, and add the following property:

[**Click here to view code image**](#)

```
property>
  <name>dfs.nfs3.dump.dir</name>
  <value>/tmp/.hdfs-nfs</value>
</property>
```

The NFSv3 `dump` directory is needed because the NFS client often reorders writes. Sequential writes can arrive at the NFS gateway in random order. This directory is used to temporarily save out-of-order writes before writing to HDFS. Make sure the `dump` directory has enough space. For

example, if the application uploads 10 files, each of size 100MB, it is recommended that this directory have 1GB of space to cover a worst-case write reorder for every file.

Once all the changes have been made, click the green Save button and note the changes you made to the Notes box in the Save confirmation dialog. Then restart all of HDFS by clicking the orange Restart button.

Step 2: Start the Gateway

Log into a DataNode and make sure all NFS services are stopped. In this example, DataNode n0 is used as the gateway.

```
# service rpcbind stop  
# service nfs stop
```

Next, start the HDFS gateway by using the `hadoop-daemon` script to start `portmap` and `nfs3` as follows:

[Click here to view code image](#)

```
# /usr/hdp/2.2.4.2-2/hadoop/sbin/hadoop-daemon.sh start portmap  
# /usr/hdp/2.2.4.2-2/hadoop/sbin/hadoop-daemon.sh start nfs3
```

The `portmap` daemon will write its log to

[Click here to view code image](#)

```
/var/log/hadoop/root/hadoop-root-portmap-n0.log
```

The `nfs3` daemon will write its log to

[Click here to view code image](#)

```
/var/log/hadoop/root/hadoop-root-nfs3-n0.log
```

To confirm the gateway is working, issue the following command. The output should look like the following:

[Click here to view code image](#)

```
# rpcinfo -p n0  
program vers proto port service  
100005 2 tcp 4242 mountd  
100000 2 udp 111 portmapper  
100000 2 tcp 111 portmapper  
100005 1 tcp 4242 mountd  
100003 3 tcp 2049 nfs  
100005 1 udp 4242 mountd  
100005 3 udp 4242 mountd
```

```
100005 3 tcp 4242 mountd  
100005 2 udp 4242 mountd
```

Finally, make sure the mount is available by issuing the following command:

```
# showmount -e n0  
Export list for n0:  
/*
```

If the `rpcinfo` or `showmount` command does not work correctly, check the previously mentioned log files for problems.

Step 3: Mount HDFS

The final step is to mount HDFS on a client node. In this example, the main login node is used. To mount the HDFS files, exit from the gateway node (in this case node `n0`) and create the following directory:

```
# mkdir /mnt/hdfs
```

The mount command is as follows. Note that the name of the gateway node will be different on other clusters, and an IP address can be used instead of the node name.

[Click here to view code image](#)

```
# mount -t nfs -o vers=3,proto=tcp,nolock n0:/ /mnt/hdfs/
```

Once the file system is mounted, the files will be visible to the client users. The following command will list the mounted file system:

[Click here to view code image](#)

```
# ls /mnt/hdfs
```

```
app-logs apps benchmarks hdp mapred mr-history system tmp user var
```

The gateway in the current Hadoop release uses `AUTH_UNIX`-style authentication and requires that the login user name on the client match the user name that NFS passes to HDFS. For example, if the NFS client is user `admin`, the NFS gateway will access HDFS as user `admin` and existing HDFS permissions will prevail.

The system administrator must ensure that the user on the NFS client machine has the same user name and user ID as that on the NFS gateway machine. This is usually not a problem if you use the same user management system, such as LDAP/NIS, to create and deploy users to cluster nodes.

CAPACITY SCHEDULER BACKGROUND

The Capacity scheduler is the default scheduler for YARN that enables multiple groups to securely share a large Hadoop cluster. Developed by the original Hadoop team at Yahoo!, the Capacity scheduler has successfully run many of the largest Hadoop clusters.

To use the Capacity scheduler, one or more queues are configured with a predetermined fraction of the total slot (or processor) capacity. This assignment guarantees a minimum amount of resources for each queue. Administrators can configure soft limits and optional hard limits on the capacity allocated to each queue. Each queue has strict ACLs (Access Control Lists) that control which users can submit applications to individual queues. Also, safeguards are in place to ensure that users cannot view or modify applications from other users.

The Capacity scheduler permits sharing a cluster while giving each user or group certain minimum capacity guarantees. These minimum amounts are not given away in the absence of demand (i.e., a group is always guaranteed a minimum number of resources is available). Excess slots are given to the most starved queues, based on the number of running tasks divided by the queue capacity. Thus, the fullest queues as defined by their initial minimum capacity guarantee get the most needed resources. Idle capacity can be assigned and provides elasticity for the users in a cost-effective manner.

Administrators can change queue definitions and properties, such as capacity and ACLs, at run time without disrupting users. They can also add more queues at run time, but cannot delete queues at run time. In addition, administrators can stop queues at run time to ensure that while existing applications run to completion, no new applications can be submitted.

The Capacity scheduler currently supports memory-intensive applications, where an application can optionally specify higher memory resource requirements than the default. Using information from the NodeManagers, the Capacity scheduler can then place containers on the best-suited nodes.

The Capacity scheduler works best when the workloads are well known, which helps in assigning the minimum capacity. For this scheduler to work most effectively, each queue should be assigned a minimal capacity that is less than the maximal expected workload. Within each queue, multiple jobs are scheduled using hierarchical (first in, first out) FIFO queues similar to the approach used with the stand-alone FIFO scheduler. If there are no queues configured, all jobs are placed in the default queue.

The ResourceManager UI provides a graphical representation of the scheduler queues and their utilization. [Figure 10.6](#) shows two jobs running on a four-node cluster. To select the scheduler view, click the Scheduler option at the bottom of the left-side vertical menu.

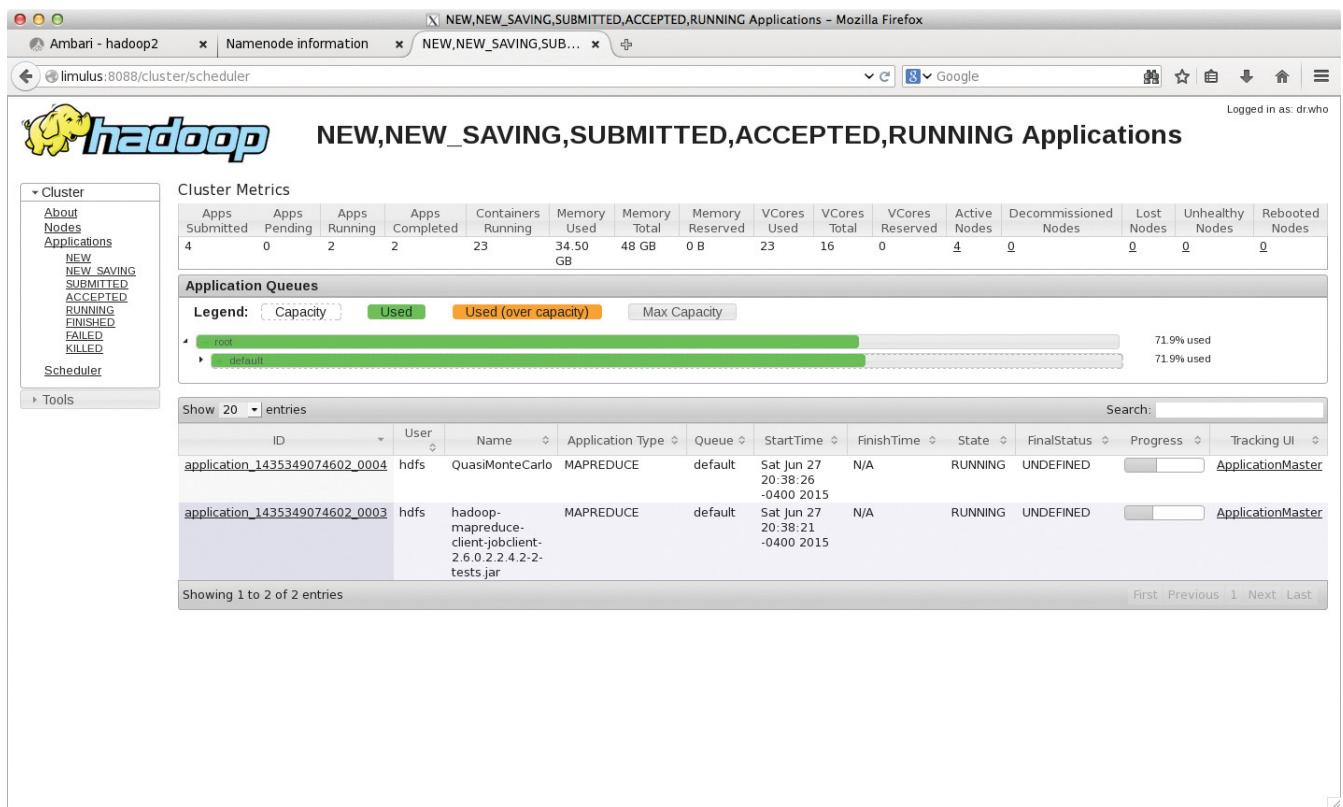


Figure 10.6 Apache YARN ResourceManager web interface showing Capacity scheduler information

Information on configuring the Capacity scheduler can be found at <https://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn-site/CapacityScheduler.html> and from *Apache Hadoop YARN: Moving beyond MapReduce and Batch Processing with Apache Hadoop 2*. (See the list of resources at the end of this chapter.) In addition to the Capacity scheduler, Hadoop YARN offers a Fair scheduler. More information can be found on the Hadoop website.

HADOOP VERSION 2 MAPREDUCE COMPATIBILITY

Hadoop version 1 is essentially a monolithic MapReduce engine. Moving this technology to YARN as a separate application framework was a complex task because MapReduce requires many important processing features, including data locality, fault tolerance, and application priorities. See [Chapter 8, “Hadoop YARN Applications,”](#) for more background on the structure of YARN applications.

To provide data locality, the MapReduce ApplicationMaster is required to locate blocks for processing and then request containers on these blocks. To implement fault tolerance, the capability to handle failed map or reduce tasks and request them again on other nodes was needed. Fault tolerance moved hand-in-hand with the complex intra-application priorities.

The logic to handle complex intra-application priorities for map and reduce tasks had to be built into the ApplicationMaster. There was no need to start idle reducers before mappers finished processing enough data. Reducers were now under control of the ApplicationMaster and were

not fixed, as they had been in Hadoop version 1. This design actually made Hadoop version 2 much more efficient and increased cluster throughput.

The following sections provide basic background on how the MapReduce framework operates under YARN. The new Hadoop version 2 MapReduce (often referred to as MRv2) was designed to provide as much backward compatibility with Hadoop version 1 MapReduce (MRv1) as possible. As with the other topics in this chapter, the following discussion provides an overview of some of the important considerations when administering the Hadoop version 2 MapReduce framework. For more information, consult <https://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>.

Enabling ApplicationMaster Restarts

Should an error occur in a MapReduce job, the ApplicationMaster can be automatically restarted by YARN. To enable ApplicationMaster restarts, set the following properties:

- Inside `yarn-site.xml`, you can tune the property `yarn.resourcemanager.am.max-retries`. The default is 2.
- Inside `mapred-site.xml`, you can more directly tune how many times a MapReduce ApplicationMaster should restart with the property `mapreduce.am.max-attempts`. The default is 2.

Calculating the Capacity of a Node

YARN has removed the hard-partitioned mapper and reducer slots of Hadoop version 1. To determine the MapReduce capacity of a cluster node, new capacity calculations are required. Estimates as to the number of mapper and reducer tasks that can efficiently run on a node help determine the amount of computing resources made available to Hadoop users. There are eight important parameters for calculating a node's MapReduce capacity; they are found in the `mapred-site.xml` and `yarn-site.xml` files.

- `mapred-site.xml`
- `mapreduce.map.memory.mb`
`mapreduce.reduce.memory.mb`

The hard limit enforced by Hadoop on the mapper or reducer task.

- `mapreduce.map.java.opts`
`mapreduce.reduce.java.opts`

The heap size of the `jvm -Xmx` for the mapper or reducer task. Remember to leave room for the JVM Perm Gen and Native Libs used. This value should always be smaller than `mapreduce.[map|reduce].memory.mb`.

- `yarn-site.xml`
- `yarn.scheduler.minimum-allocation-mb`

The smallest container YARN will allow.

- `yarn.scheduler.maximum-allocation-mb`

The largest container YARN will allow.

- `yarn.nodemanager.resource.memory-mb`

The amount of physical memory (RAM) on the compute node for containers. It is important that this value is not equal to the total RAM on the node, as other Hadoop services also require RAM.

- `yarn.nodemanager.vmem-pmem-ratio`

The amount of virtual memory each container is allowed. It is calculated with the following formula:

Click here to view code image

```
containerMemoryRequest * vmem-pmem-ratio.
```

As an example, consider a configuration with the settings in [Table 10.1](#). Using these settings, we have given each map and reduce task a generous 512MB of overhead for the container, as seen with the difference between the `mapreduce.[map|reduce].memory.mb` and the `mapreduce.[map|reduce].java.opts`.

Property	Value
<code>mapreduce.map.memory.mb</code>	1536
<code>mapreduce.reduce.memory.mb</code>	2560
<code>mapreduce.map.java.opts</code>	<code>-Xmx1024m</code>
<code>mapreduce.reduce.java.opts</code>	<code>-Xmx2048m</code>
<code>yarn.scheduler.minimum-allocation-mb</code>	512
<code>yarn.scheduler.maximum-allocation-mb</code>	4096
<code>yarn.nodemanager.resource.memory-mb</code>	36864
<code>yarn.nodemanager.vmem-pmem-ratio</code>	2.1

Table 10.1 Example YARN MapReduce Settings

Next, we have configured YARN to allow a container no smaller than 512MB and no larger than 4GB. Assuming the compute nodes have 36GB of RAM available for containers, and with a virtual memory ratio of 2.1 (the default value), each map can have as much as 3225.6MB of RAM and a reducer can have 5376MB of virtual RAM. Thus our compute node configured for 36GB of container space can support up to 24 maps or 14 reducers, or any combination of mappers and reducers allowed by the available resources on the node.

Running Hadoop Version 1 Applications

To ease the transition from Hadoop version 1 to version 2 with YARN, a major goal of YARN and the MapReduce framework implementation on top of YARN is to ensure that existing MapReduce applications that were programmed and compiled against previous MapReduce APIs (MRv1 applications) can continue to run with little work on top of YARN (MRv2 applications).

Binary Compatibility of org.apache.hadoop.mapred APIs

For the vast majority of users who use the `org.apache.hadoop.mapred` APIs, MapReduce on YARN ensures full binary compatibility. These existing applications can run on YARN directly without recompilation. You can use jar files of your existing application that code against MapReduce APIs and use `bin/hadoop` to submit them directly to YARN.

Source Compatibility of org.apache.hadoop.mapreduce APIs

Unfortunately, it has proved difficult to ensure full binary compatibility of applications that were originally compiled against MRv1 `org.apache.hadoop.mapreduce` APIs. These APIs have gone through lots of changes. For example, many of the classes stopped being abstract classes and changed to interfaces. The YARN community eventually reached a compromise on this issue, supporting source compatibility only for `org.apache.hadoop.mapreduce` APIs. Existing applications using MapReduce APIs are source compatible and can run on YARN either with no changes, with simple recompilation against MRv2 jar files that are shipped with Hadoop version 2, or with minor updates.

Compatibility of Command-Line Scripts

Most of the command-line scripts from Hadoop 1.x should work without any tweaking. The only exception is `mrmadmin`, whose functionality was removed from MRv2 because the JobTracker and TaskTracker no longer exist. The `mrmadmin` functionality has been replaced with `rmadmin`. The suggested method to invoke `rmadmin` is through the command line, even though you can directly invoke the APIs. In YARN, when `mrmadmin` commands are executed, warning messages will appear, reminding users to use YARN commands (i.e., `rmadmin` commands). Conversely, if the user's applications programmatically invoke `mrmadmin`, those applications will break when running on top of YARN. There is no support for either binary or source compatibility under YARN.

Running Apache Pig Scripts on YARN

Pig is one of the two major data process applications in the Hadoop ecosystem, with the other being Hive (see [Chapter 7, “Essential Hadoop Tools”](#)). Thanks to the significant efforts made by the Pig community, Pig scripts of existing users do not need any modifications. Pig on YARN in Hadoop 0.23 has been supported since version 0.10.0, and Pig working with Hadoop 2.x has been supported since version 0.10.1.

Existing Pig scripts that work with Pig 0.10.1 and beyond will work just fine on top of YARN. In contrast, versions earlier than Pig 0.10.x may not run directly on YARN due to some of the incompatible MapReduce APIs and configuration.

Running Apache Hive Queries on YARN

Hive queries of existing users do not need any changes to work on top of YARN, starting with Hive 0.10.0, thanks to the work done by Hive community. Support for Hive to work on YARN in the Hadoop 0.23 and 2.x releases has been in place since version 0.10.0. Queries that work in Hive 0.10.0 and beyond will work without changes on top of YARN. However, as with Pig, earlier versions of Hive may not run directly on YARN, as those Hive releases do not support Hadoop 0.23 and 2.x.

Running Apache Oozie Workflows on YARN

Like the Pig and Hive communities, the Apache Oozie community worked to ensure existing Oozie workflows would run in a completely backward-compatible manner on Hadoop version 2. Support for Hadoop 0.23 and 2.x is available starting with Oozie release 3.2.0. Existing Oozie

workflows can start taking advantage of YARN in versions 0.23 and 2.x with Oozie 3.2.0 and above.

SUMMARY AND ADDITIONAL RESOURCES

Administering an Apache Hadoop cluster may involve many different services and issues. The two core Hadoop services are the YARN ResourceManager and the HDFS. The Apache Ambari management GUI, presented in [Chapter 9](#), is a good tool for making changes to these services. Several basic YARN administration topics were presented in this chapter, including using the `radmin` tool, taking advantage of the YARN Web Proxy, decommissioning YARN nodes, managing YARN applications, and setting important YARN properties.

Important topics and procedures related to basic HDFS administration were described as well. These include exploring the NameNode UI and learning how to add HDFS users, run file system checks, balance DataNodes, create HDFS snapshots, and start an HDFS NFSv3 gateway. In addition, this chapter provided background information on HDFS Safe Mode and the SecondaryNameNode.

The Capacity scheduler offers key functionality within Hadoop version 2, as discussed in this chapter. An example of the Scheduler view on the YARN ResourceManager UI was provided in this chapter as well. Finally, Hadoop version 2 MapReduce compatibility and node capacity are issues that many programmers need to address.

Additional information can be found from the following sources:

■ Apache Hadoop YARN Administration

- <http://hadoop.apache.org/docs/current/> (Scroll down to the lower left-hand corner under “Configuration.”)
- https://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn-site/YarnCommands.html#Administration_Commands (YARN Administration Commands)
- Book: Murthy, A., et al. 2014. *Apache Hadoop YARN: Moving beyond MapReduce and Batch Processing with Apache Hadoop 2*. Boston, MA: Addison-Wesley. <http://www.informit.com/store/apache-hadoop-yarn-moving-beyond-mapreduce-and-batch-9780321934505>

■ HDFS Administration:

- <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsUserGuide.html> (Notice the HDFS topic menu on the left-hand side of the page.)
- <http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsPermissionsGuide.html> (HDFS permissions)
- <http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HDFSCommands.html> (HDFS commands)
- <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsSnapshots.htm> (HDFS snapshots)
- <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsNfsGateway.html> (HDFS NFSv3 gateway)

■ Capacity Scheduler Administration

- <https://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn-site/CapacityScheduler.html>(Capacity scheduler configuration)
- Book: Murthy, A., et al. 2014. *Apache Hadoop YARN: Moving beyond MapReduce and Batch Processing with Apache Hadoop 2*. Boston, MA: Addison-Wesley. <http://www.informit.com/store/apache-hadoop-yarn-moving-beyond-mapreduce-and-batch-9780321934505> (See “Apache Hadoop YARN Administration”)
- **MapReduce Version 2 (MRv2) Administration**
- <https://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html> (MapReduce with YARN)
- Book: Murthy, A., et al. 2014. *Apache Hadoop YARN: Moving beyond MapReduce and Batch Processing with Apache Hadoop 2*. Boston, MA: Addison-Wesley. <http://www.informit.com/store/apache-hadoop-yarn-moving-beyond-mapreduce-and-batch-9780321934505>



Future Vision

FUTURE VISION BIE

By K B Hemanth Raj

Visit : <https://hemanthrajhemu.github.io>

Quick Links for Faster Access.

CSE 8th Semester - <https://hemanthrajhemu.github.io/CSE8/>

ISE 8th Semester - <https://hemanthrajhemu.github.io/ISE8/>

ECE 8th Semester - <https://hemanthrajhemu.github.io/ECE8/>

8th Semester CSE - TEXTBOOK - NOTES - QP - SCANNER & MORE

17CS81 IOT - <https://hemanthrajhemu.github.io/CSE8/17SCHEME/17CS81/>

17CS82 BDA - <https://hemanthrajhemu.github.io/CSE8/17SCHEME/17CS82/>

17CS832 UID - <https://hemanthrajhemu.github.io/CSE8/17SCHEME/17CS832/>

17CS834 SMS - <https://hemanthrajhemu.github.io/CSE8/17SCHEME/17CS834/>

8th Semester Computer Science & Engineering (CSE)

8th Semester CSE Text Books: <https://hemanthrajhemu.github.io/CSE8/17SCHEME/Text-Book.html>

8th Semester CSE Notes: <https://hemanthrajhemu.github.io/CSE8/17SCHEME/Notes.html>

8th Semester CSE Question Paper: <https://hemanthrajhemu.github.io/CSE8/17SCHEME/Question-Paper.html>

8th Semester CSE Scanner: <https://hemanthrajhemu.github.io/CSE8/17SCHEME/Scanner.html>

8th Semester CSE Question Bank: <https://hemanthrajhemu.github.io/CSE8/17SCHEME/Question-Bank.html>

8th Semester CSE Answer Script: <https://hemanthrajhemu.github.io/CSE8/17SCHEME/Answer-Script.html>

Contribution Link:

<https://hemanthrajhemu.github.io/Contribution/>

Stay Connected... get Updated... ask your queries...

Join Telegram to get Instant Updates:

<https://telegram.me/joinchat/AAAAAFTp8kuvCHALxuMaQ>

Contact: MAIL: futurevisionbie@gmail.com

INSTAGRAM: www.instagram.com/futurevisionbie/