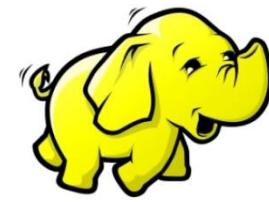


edureka!

Big Data & Hadoop

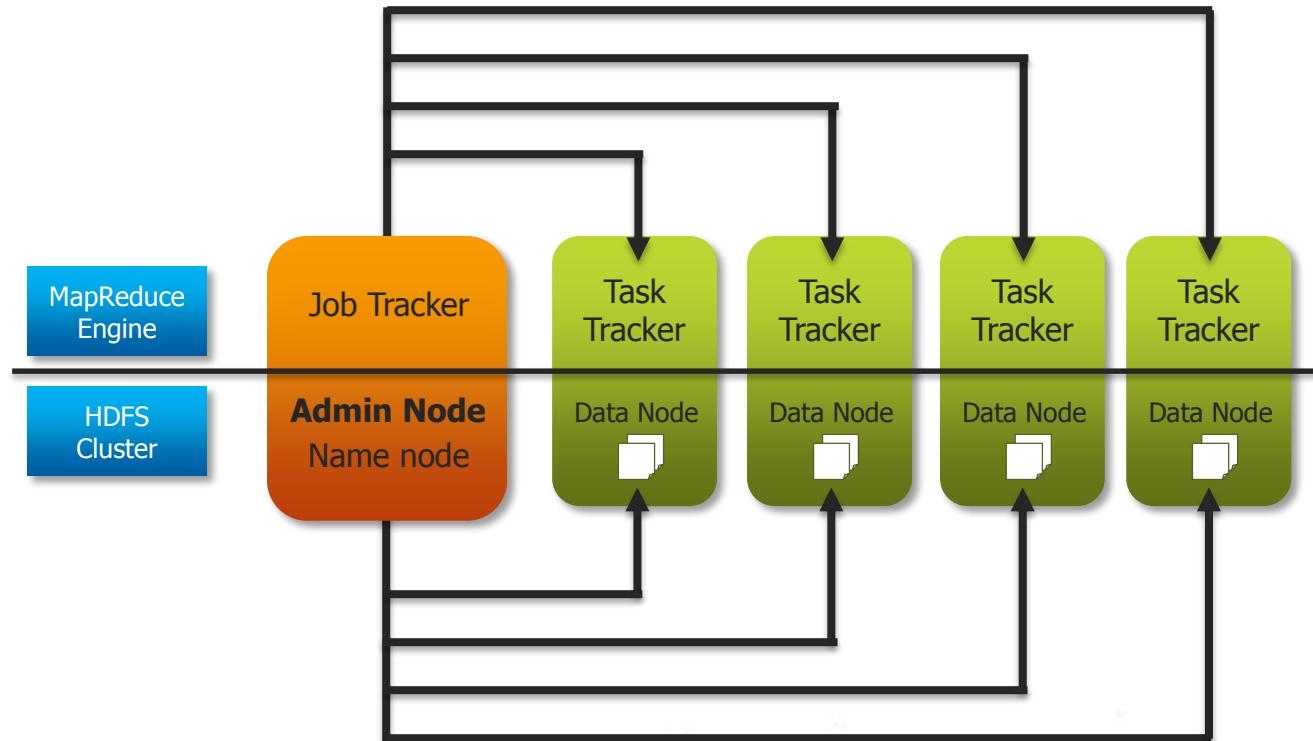


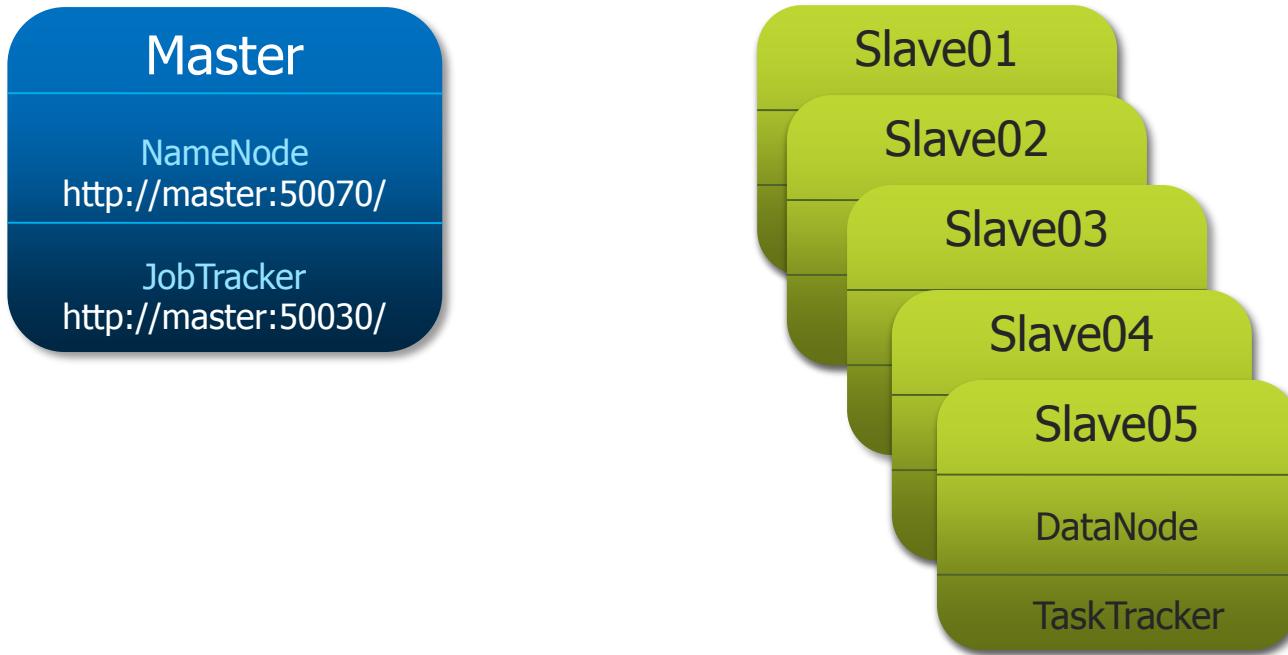
- ✓ **Module 1**
 - ✓ Understanding Big Data
 - ✓ Hadoop Architecture
- ✓ **Module 2**
 - ✓ **Hadoop Cluster Configuration**
 - ✓ **Data loading Techniques**
 - ✓ **Hadoop Project Environment**
- ✓ **Module 3**
 - ✓ Hadoop MapReduce framework
 - ✓ Programming in Map Reduce
- ✓ **Module 4**
 - ✓ Advance MapReduce
 - ✓ MRUnit testing framework
- ✓ **Module 5**
 - ✓ Analytics using Pig
 - ✓ Understanding Pig Latin
- ✓ **Module 6**
 - ✓ Analytics using Hive
 - ✓ Understanding HIVE QL
- ✓ **Module 7**
 - ✓ Advance Hive
 - ✓ NoSQL Databases and HBASE
- ✓ **Module 8**
 - ✓ Advance HBASE
 - ✓ Zookeeper Service
- ✓ **Module 9**
 - ✓ Hadoop 2.0 – New Features
 - ✓ Programming in MRv2
- ✓ **Module 10**
 - ✓ Apache Oozie
 - ✓ Real world Datasets and Analysis
 - ✓ Project Discussion

- ✓ Revision
- ✓ A Typical Hadoop Cluster
- ✓ Hadoop Cluster Configuration
- ✓ Hadoop Configuration Files
- ✓ Hadoop Cluster Modes
- ✓ Sample example list in Hadoop
- ✓ Running Teragen Example
- ✓ Dump of MR job
- ✓ Data Loading Techniques
 - ✓ Using Hadoop Copy Commands
 - ✓ FLUME
 - ✓ SQOOP
- ✓ Course Project Problem Statement

Lets's Revise

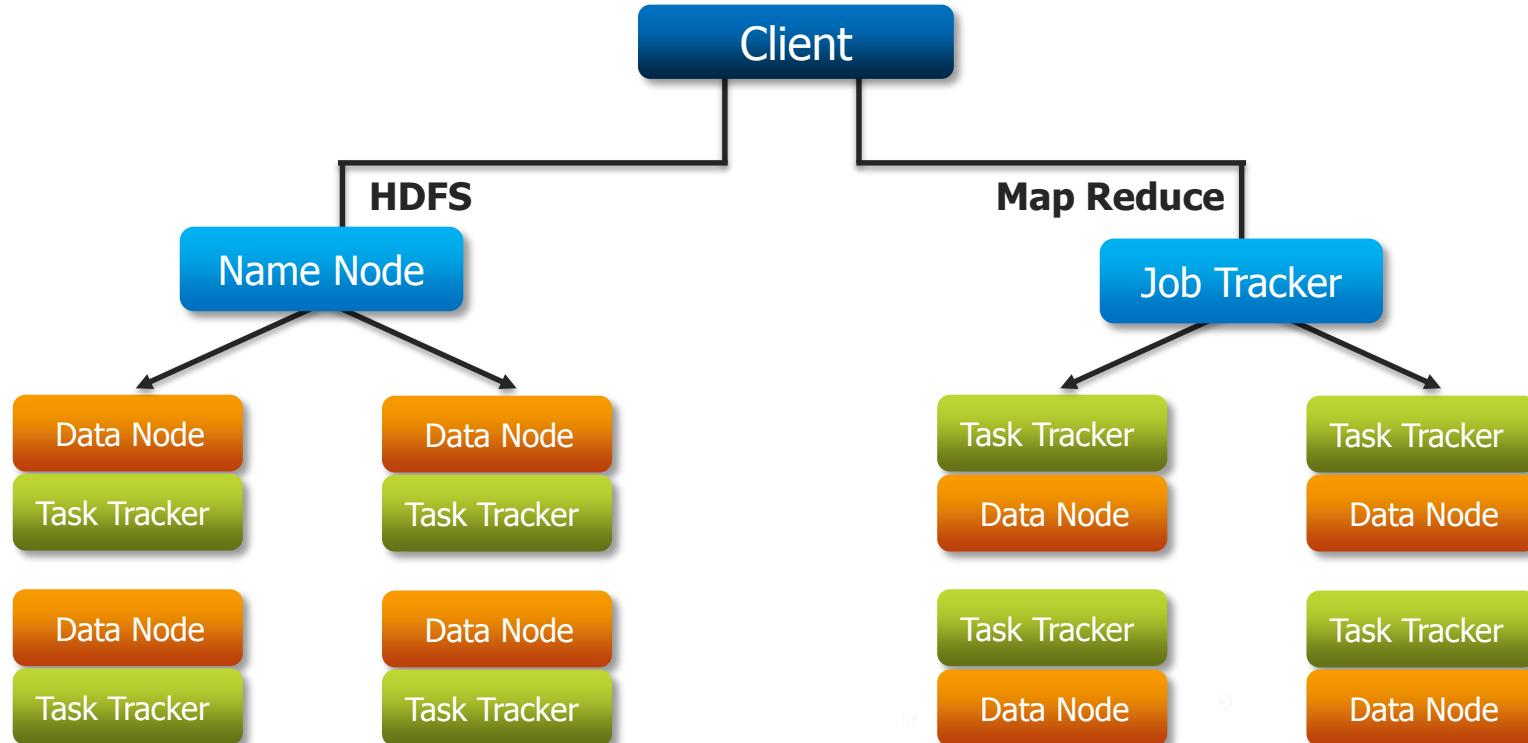
1. Hadoop Core Components
2. HDFS Architecture
3. What is HDFS?
4. Hadoop Vs. traditional systems
5. Namenode and Secondary Namenode





Hadoop Cluster Architecture (Contd.)

edureka!



Pre-Class Questions



The default replication factor is:

- 2
- 4
- 5
- 3

The answer is 3. It means if you move a file to hdfs then by default 3 copies of the file will be stored on different datanodes.



Every Slave node has two daemons running on them that is DataNode & TaskTracker in a MultiNode Cluster.

- TRUE
- FALSE



TRUE. DataNode service for HDFS and TaskTracker for processing (MapReduce) tasks.



A block is replicated in 4 nodes K,L,M, and N. If M, K and N fails. A client can still read the data.

- TRUE
- FALSE

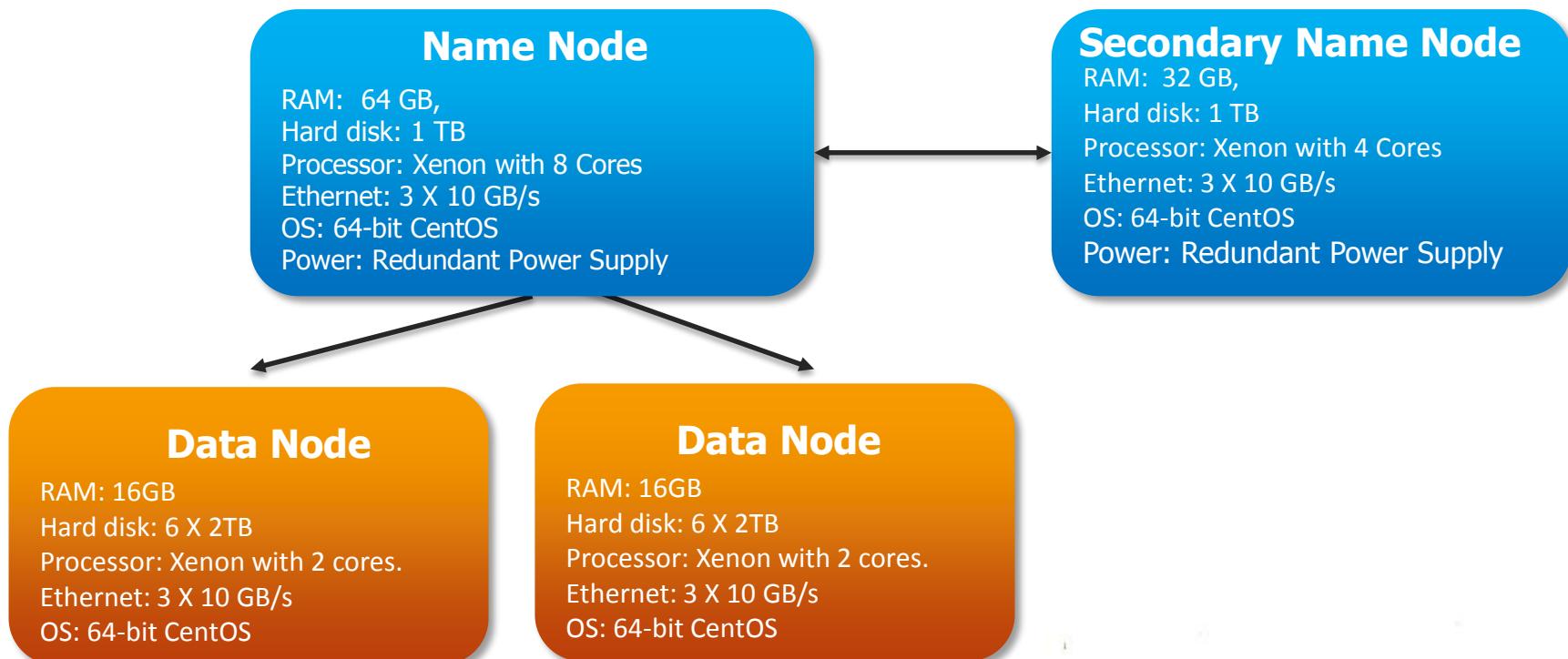


TRUE. As the remaining node 'L' will contain the block in question.



Hadoop Cluster: A Typical Use Case

edureka!





Facebook

- We use Hadoop to store copies of internal log and dimension data sources and use it as a source for reporting/analytics and machine learning.
- Currently we have 2 major clusters:
 - A 1100-machine cluster with 8800 cores and about 12 PB raw storage.
 - A 300-machine cluster with 2400 cores and about 3 PB raw storage.
 - Each (commodity) node has 8 cores and 12 TB of storage.
 - We are heavy users of both streaming as well as the Java APIs. We have built a higher level data warehousing framework using these features called Hive (see the <http://hadoop.apache.org/hive/>). We have also developed a FUSE implementation over HDFS.

<http://wiki.apache.org/hadoop/PoweredBy>

Hadoop can run in any of the following three modes:

Standalone (or Local) Mode

- ✓ No daemons, everything runs in a single JVM.
- ✓ Suitable for running MapReduce programs during development.
- ✓ Has no DFS.

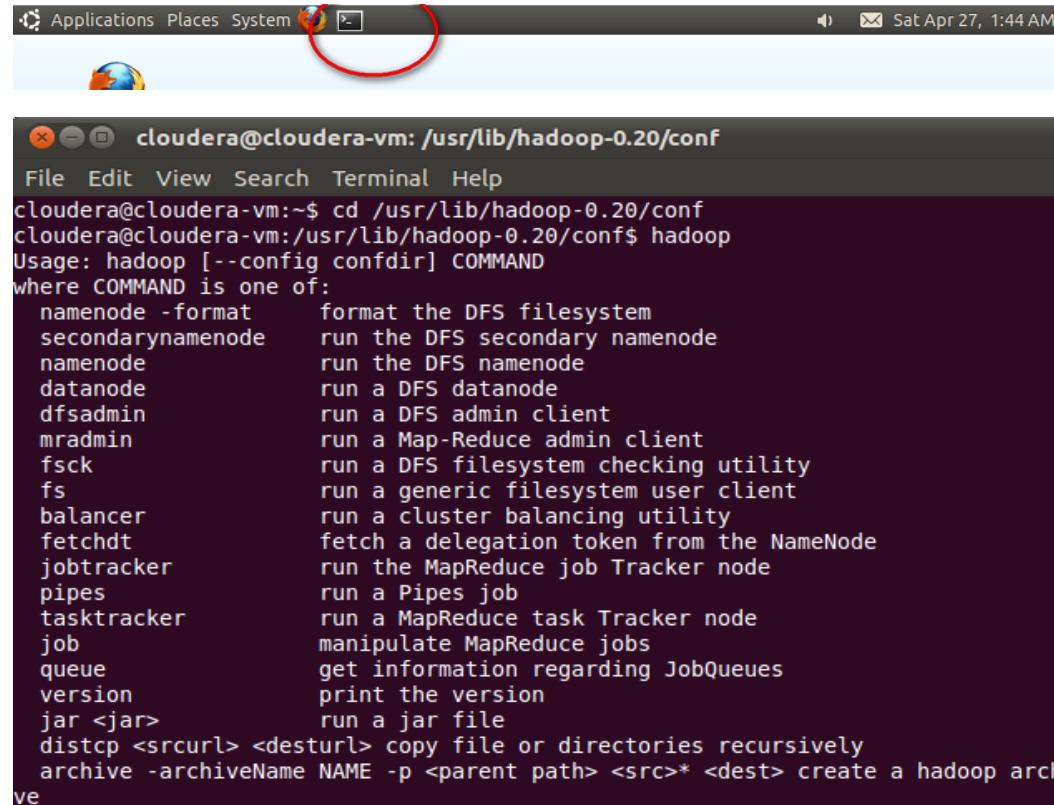
Pseudo-Distributed Mode

- ✓ Hadoop daemons run on the local machine.

Fully-Distributed Mode

- ✓ Hadoop daemons run on a cluster of machines.

Terminal Commands



The screenshot shows a Linux desktop interface with a terminal window open. The terminal window has a red circle drawn around its title bar. The title bar displays "cloudera@cloudera-vm: /usr/lib/hadoop-0.20/conf". The terminal content shows the usage information for the Hadoop command-line interface:

```
File Edit View Search Terminal Help
cloudera@cloudera-vm:~$ cd /usr/lib/hadoop-0.20/conf
cloudera@cloudera-vm:/usr/lib/hadoop-0.20/conf$ hadoop
Usage: hadoop [--config confdir] COMMAND
where COMMAND is one of:
  namenode -format      format the DFS filesystem
  secondarynamenode     run the DFS secondary namenode
  namenode              run the DFS namenode
  datanode              run a DFS datanode
  dfsadmin              run a DFS admin client
  mradmin               run a Map-Reduce admin client
  fsck                  run a DFS filesystem checking utility
  fs                   run a generic filesystem user client
  balancer              run a cluster balancing utility
  fetchdt              fetch a delegation token from the NameNode
  jobtracker            run the MapReduce job Tracker node
  pipes                 run a Pipes job
  tasktracker           run a MapReduce task Tracker node
  job                  manipulate MapReduce jobs
  queue                get information regarding JobQueues
  version               print the version
  jar <jar>              run a jar file
  distcp <srcurl> <desturl> copy file or directories recursively
  archive -archiveName NAME -p <parent path> <src>* <dest> create a hadoop archive
```

Listing of files present on HDFS

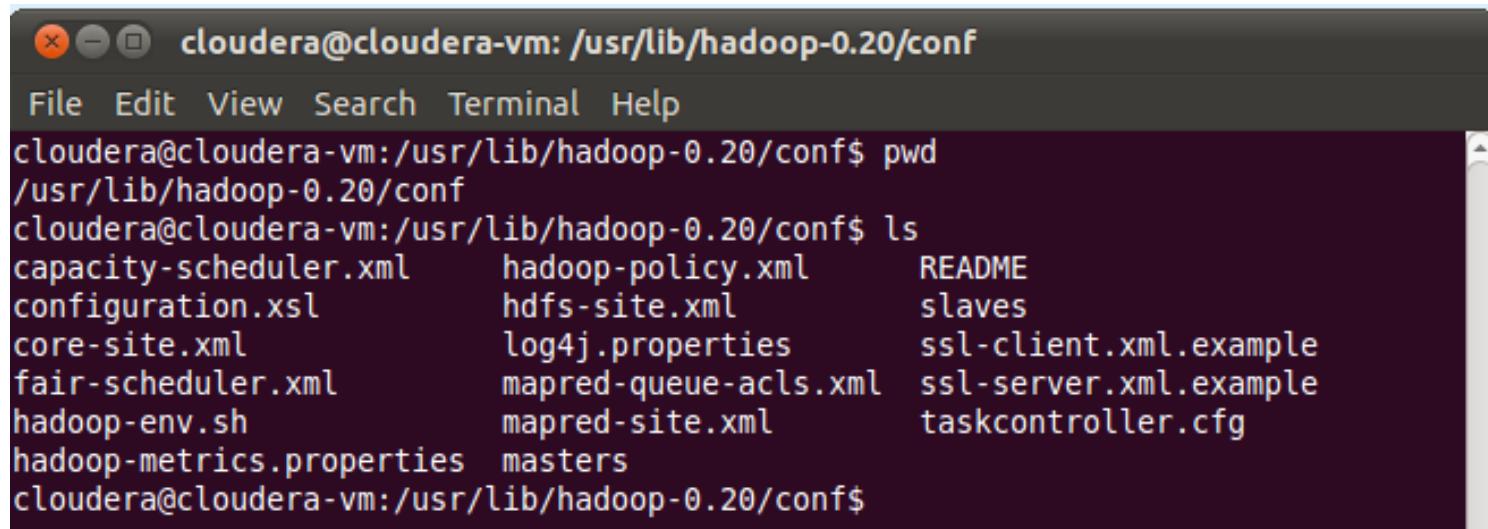
```
cloudera@cloudera-vm:/usr/lib/hadoop-0.20/conf$ hadoop dfs -ls /
Found 3 items
drwxrwxrwx  - hue    supergroup          0 2013-04-20 03:01 /tmp
drwxr-xr-x  - hue    supergroup          0 2013-04-20 03:22 /user
drwxr-xr-x  - mapred supergroup          0 2011-04-12 06:20 /var
cloudera@cloudera-vm:/usr/lib/hadoop-0.20/conf$ █
```

Listing of files present in bin Directory

```
cloudera@cloudera-vm:~$ cd /usr/lib/hadoop-0.20/bin/
cloudera@cloudera-vm:/usr/lib/hadoop-0.20/bin$ ls
fuse_dfs_wrapper.sh  hadoop-daemons.sh  start-all.sh      stop-all.sh
hadoop              rcc                  start-balancer.sh  stop-balancer.sh
hadoop-config.sh    README                start-dfs.sh      stop-dfs.sh
hadoop-daemon.sh   slaves.sh             start-mapred.sh  stop-mapred.sh
cloudera@cloudera-vm:/usr/lib/hadoop-0.20/bin$
```

Hadoop Configuration Files

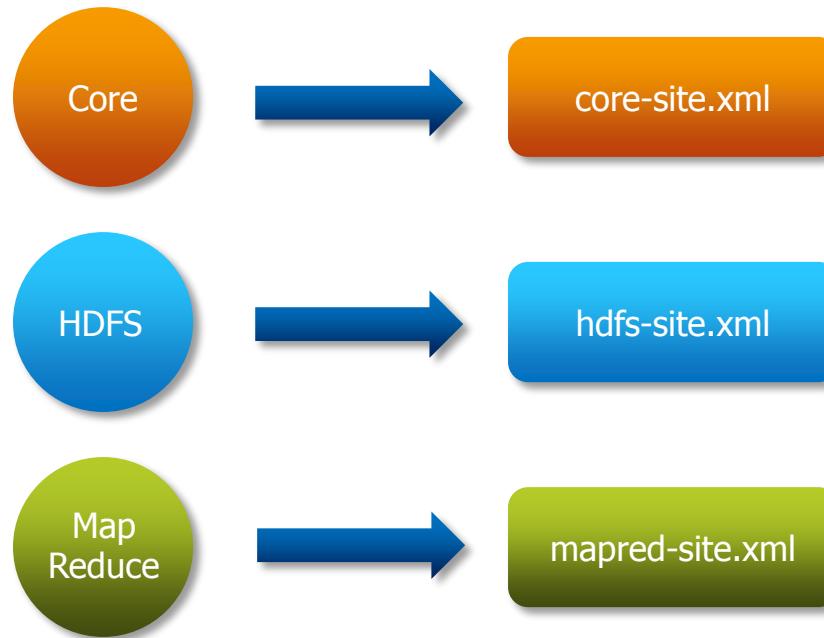
Configuration Filenames	Description of Log Files
hadoop-env.sh	Environment variables that are used in the scripts to run Hadoop.
core-site.xml	Configuration settings for Hadoop Core such as I/O settings that are common to HDFS and MapReduce.
hdfs-site.xml	Configuration settings for HDFS daemons, the namenode, the secondary namenode and the data nodes.
mapred-site.xml	Configuration settings for MapReduce daemons : the job-tracker and the task-trackers.
masters	A list of machines (one per line) that each run a secondary namenode.
slaves	A list of machines (one per line) that each run a datanode and a task-tracker.



A screenshot of a terminal window titled "cloudera@cloudera-vm: /usr/lib/hadoop-0.20/conf". The window shows a list of configuration files in the "/usr/lib/hadoop-0.20/conf" directory. The files listed are:

- capacity-scheduler.xml
- configuration.xsl
- core-site.xml
- fair-scheduler.xml
- hadoop-env.sh
- hadoop-metrics.properties
- hadoop-policy.xml
- hdfs-site.xml
- log4j.properties
- mapred-queue-acls.xml
- mapred-site.xml
- masters
- README
- slaves
- ssl-client.xml.example
- ssl-server.xml.example
- taskcontroller.cfg

```
cloudera@cloudera-vm: /usr/lib/hadoop-0.20/conf$ pwd
/usr/lib/hadoop-0.20/conf
cloudera@cloudera-vm: /usr/lib/hadoop-0.20/conf$ ls
capacity-scheduler.xml      hadoop-policy.xml      README
configuration.xsl           hdfs-site.xml        slaves
core-site.xml                log4j.properties    ssl-client.xml.example
fair-scheduler.xml          mapred-queue-acls.xml  ssl-server.xml.example
hadoop-env.sh                 mapred-site.xml   taskcontroller.cfg
hadoop-metrics.properties   masters
cloudera@cloudera-vm: /usr/lib/hadoop-0.20/conf$
```



hdfs-site.xml	core-site.xml
<?xml version - "1.0"?>	<?xml version ="1.0"?>
<!--hdfs-site.xml-->	<!--core-site.xml-->
<configuration>	<configuration>
<property>	<property>
<name> dfs.replication </name>	<name> fs.default.name </name>
<value> 1 </value>	<value> hdfs://localhost:8020/ </value>
</property>	</property>
</configuration>	</configuration>

Defining HDFS Details In hdfs-site.xml

Property	Value	Description
dfs.data.dir	<value> /disk1/hdfs/data, /disk2/hdfs/data </value>	A list of directories where the datanode stores blocks. Each block is stored in only one of these directories. \${hadoop.tmp.dir}/dfs/data
fs.checkpoint.dir	<value> /disk1/hdfs/namesecondary, /disk2/hdfs/namesecondary </value>	A list of directories where the secondary namenode stores checkpoints. It stores a copy of the checkpoint in each directory in the list \${hadoop.tmp.dir}/dfs/namesecondary

mapred-site.xml

```
<?xml version="1.0"?>
<configuration>
<property>
    <name>mapred.job.tracker</name>
    <value>localhost:8021</value>
<property>
</configuration>
```

Defining mapred-site.xml

Property	Value	Description
mapred.job.tracker	<value> localhost:8021 <td>The hostname and the port that the jobtracker RPC server runs on. If set to the default value of local, then the jobtracker runs in-process on demand when you run a MapReduce job.</td>	The hostname and the port that the jobtracker RPC server runs on. If set to the default value of local, then the jobtracker runs in-process on demand when you run a MapReduce job.
mapred.local.dir	\${hadoop.tmp.dir}/mapred/local	A list of directories where MapReduce stores intermediate data for jobs. The data is cleared out when the job ends.
mapred.system.dir	\${hadoop.tmp.dir}/mapred/system	The directory relative to fs.default.name where shared files are stored, during a job run.
mapred.tasktracker.map.tasks.maximum	2	The number of map tasks that may be run on a tasktracker at any one time
mapred.tasktracker.reduce.tasks.maximum	2	The number of reduce tasks tat may be run on a tasktracker at any one time.

1. <http://hadoop.apache.org/docs/r1.1.2/core-default.html>
2. <http://hadoop.apache.org/docs/r1.1.2/mapred-default.html>
3. <http://hadoop.apache.org/docs/r1.1.2/hdfs-default.html>

Two files are used by the startup and shutdown commands:

Slaves

- ✓ Contains a list of hosts, one per line, that are to host **DataNode** and **TaskTracker** servers.

Masters

- ✓ Contains a list of hosts, one per line, that are to host **Secondary NameNode** servers.



- ✓ This file also offers a way to provide custom parameters for each of the servers.
- ✓ Hadoop-env.sh is sourced by all of the Hadoop Core scripts provided in the conf/ directory of the installation.
- ✓ **Examples of environment variables that you can specify:**

```
export HADOOP_DATANODE_HEAPSIZE="128"
```

```
export HADOOP_TASKTRACKER_HEAPSIZE="512"
```

- ✓ NameNode status: <http://localhost:50070/dfshealth.jsp>
- ✓ JobTracker status: <http://localhost:50030/jobtracker.jsp>
- ✓ TaskTracker status: <http://localhost:50060/tasktracker.jsp>
- ✓ DataBlock Scanner Report: <http://localhost:50075/blockScannerReport>



Which of the following file is used to specify the NameNode's heap size?

- .bashrc
- hadoop-env.sh
- hdfs-site.sh
- core-site.xml



hadoop-env.sh. This file specifies environment variables that affect the JDK used by Hadoop Daemon (bin/hadoop).

It is necessary to define all the properties in core-site.xml, hdfs-site.xml & mapred-site.xml.

- TRUE
- FALSE



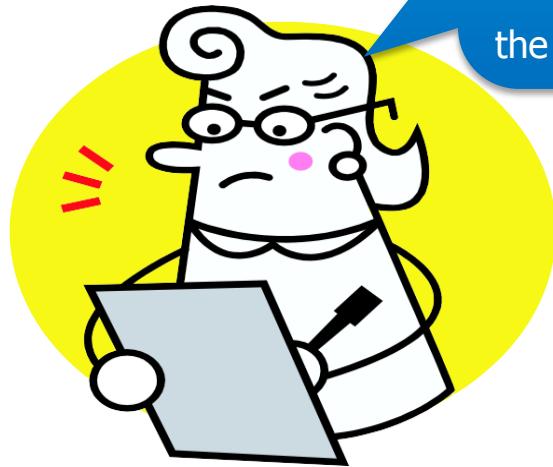


False. Detailed answer is after the next question.

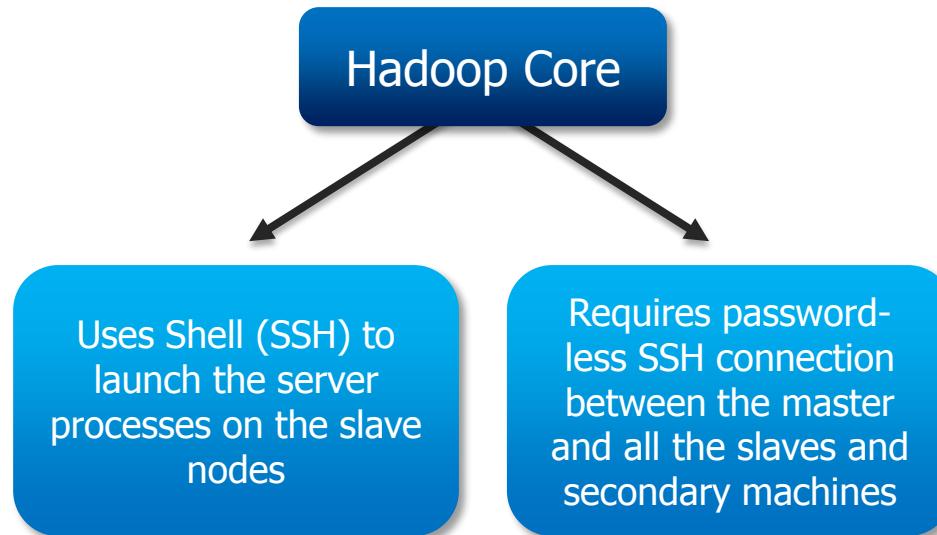
Standalone Mode uses default configuration?

- TRUE
- FALSE





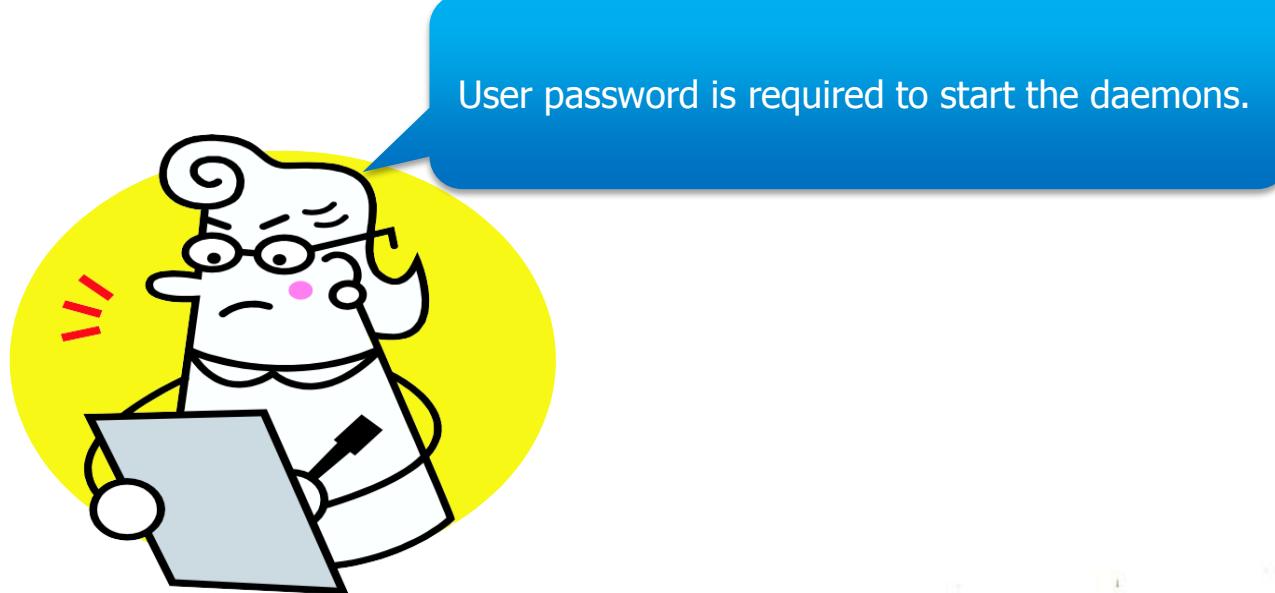
True. In Standalone mode Hadoop runs with default configuration (Empty configuration files i.e. no configuration settings in core-site.xml, hdfs-site.xml, and mapred-site.xml). If properties are not defined in the configuration files, hadoop runs with default values for the corresponding properties.





If password-less SSH Login is not set-up,

- hadoop deamons won't start
- only NameNode will start in Master
- user's password has to be entered to start every daemon
- None of these



Sample Examples List

```
cloudera@cloudera-vm:/usr/lib/hadoop-0.20$ ls
bin
build.xml
CHANGES.txt
cloudera
cloudera-pom.xml
conf
contrib
example-conf
hadoop-0.20.2-cdh3u0-ant.jar
hadoop-0.20.2-cdh3u0-core.jar
hadoop-0.20.2-cdh3u0-examples.jar
hadoop-0.20.2-cdh3u0-test.jar
hadoop-0.20.2-cdh3u0-tools.jar
hadoop-ant-0.20.2-cdh3u0.jar
hadoop-ant.jar
hadoop-core-0.20.2-cdh3u0.jar
cloudera@cloudera-vm:/usr/lib/hadoop-0.20$ hadoop jar hadoop-0.20.2-cdh3u0-examples.jar
An example program must be given as the first argument.
Valid program names are:
  aggregatewordcount: An Aggregate based map/reduce program that counts the word
  s in the input files.
```

Running the Teragen Example

```
cloudera@cloudera-vm:/usr/lib/hadoop-0.20$ hadoop jar hadoop-0.20.2-cdh3u0-examples.jar teragen 1000000 /user/teragen-inputTest
Generating 1000000 using 2 maps with step of 500000
13/04/27 02:20:58 INFO mapred.JobClient: Running job: job_201304270136_0001
13/04/27 02:20:59 INFO mapred.JobClient: map 0% reduce 0%
```

```
Generating 1000000 using 2 maps with step of 500000
13/04/27 02:20:58 INFO mapred.JobClient: Running job: job_201304270136_0001
13/04/27 02:20:59 INFO mapred.JobClient: map 0% reduce 0%
13/04/27 02:21:15 INFO mapred.JobClient: map 49% reduce 0%
13/04/27 02:21:18 INFO mapred.JobClient: map 81% reduce 0%
13/04/27 02:21:21 INFO mapred.JobClient: map 100% reduce 0%
13/04/27 02:21:23 INFO mapred.JobClient: Job complete: job_201304270136_0001
13/04/27 02:21:23 INFO mapred.JobClient: Counters: 13
13/04/27 02:21:23 INFO mapred.JobClient: Job Counters
13/04/27 02:21:23 INFO mapred.JobClient: SLOTS_MILLIS_MAPS=37341
13/04/27 02:21:23 INFO mapred.JobClient: Total time spent by all reduces waiting after reserving slots (ms)=0
13/04/27 02:21:23 INFO mapred.JobClient: Total time spent by all maps waiting after reserving slots (ms)=0
13/04/27 02:21:23 INFO mapred.JobClient: Launched map tasks=2
13/04/27 02:21:23 INFO mapred.JobClient: SLOTS_MILLIS_REDUCES=0
13/04/27 02:21:23 INFO mapred.JobClient: FileSystemCounters
13/04/27 02:21:23 INFO mapred.JobClient: HDFS_BYTES_READ=167
13/04/27 02:21:23 INFO mapred.JobClient: FILE_BYTES_WRITTEN=105170
13/04/27 02:21:23 INFO mapred.JobClient: HDFS_BYTES_WRITTEN=100000000
13/04/27 02:21:23 INFO mapred.JobClient: Map-Reduce Framework
13/04/27 02:21:23 INFO mapred.JobClient: Map input records=1000000
13/04/27 02:21:23 INFO mapred.JobClient: Spilled Records=0
13/04/27 02:21:23 INFO mapred.JobClient: Map input bytes=1000000
13/04/27 02:21:23 INFO mapred.JobClient: Map output records=1000000
13/04/27 02:21:23 INFO mapred.JobClient: SPLIT_RAW_BYTES=167
```

Checking the Output

```
cloudera@cloudera-vm:/usr/lib/hadoop-0.20/conf$ hadoop dfs -ls /
Found 3 items
drwxrwxrwx  - hue    supergroup      0 2013-04-27 01:59 /tmp
drwxr-xr-x  - hue    supergroup      0 2013-04-27 02:20 /user
drwxr-xr-x  - mapred supergroup      0 2011-04-12 06:20 /var
cloudera@cloudera-vm:/usr/lib/hadoop-0.20/conf$ hadoop dfs -ls /user
Found 5 items
drwxr-xr-x  - cloudera supergroup    0 2011-04-08 16:48 /user/cloudera
drwxr-xr-x  - hue    supergroup      0 2011-04-01 17:25 /user/hive
drwxr-xr-x  - cloudera supergroup    0 2013-04-20 03:19 /user/teragen-input14
drwxr-xr-x  - cloudera supergroup    0 2013-04-20 03:22 /user/teragen-input15
drwxr-xr-x  - cloudera supergroup    0 2013-04-27 02:21 /user/teragen-inputTest
```

Checking the Output

```
cloudera@cloudera-vm:/usr/lib/hadoop-0.20$ hadoop fs -get /user/teragen-inputTest/part-00001 .
get: Permission denied
cloudera@cloudera-vm:/usr/lib/hadoop-0.20$ sudo su
[sudo] password for cloudera:
root@cloudera-vm:/usr/lib/hadoop-0.20# hadoop fs -get /user/teragen-inputTest/part-00001 .
root@cloudera-vm:/usr/lib/hadoop-0.20# ls
bin         .hadoop-0.20.2-cdh3u0-ant.jar      hadoop-core.jar           ivy.xml
build.xml    .hadoop-0.20.2-cdh3u0-core.jar    hadoop-examples-0.20.2-cdh3u0.jar lib
CHANGES.txt   .hadoop-0.20.2-cdh3u0-examples.jar hadoop-examples.jar        logs
cloudera     .hadoop-0.20.2-cdh3u0-test.jar    hadoop-test-0.20.2-cdh3u0.jar NOTICE.txt
cloudera-pom.xml .hadoop-0.20.2-cdh3u0-tools.jar hadoop-test.jar         part-00001
conf          .hadoop-ant-0.20.2-cdh3u0.jar    hadoop-tools-0.20.2-cdh3u0.jar pids
contrib       .hadoop-ant.jar                  hadoop-tools.jar        README.txt
example-conf .hadoop-core-0.20.2-cdh3u0.jar    ivy                         webapps
root@cloudera-vm:/usr/lib/hadoop-0.20# mkdir test
root@cloudera-vm:/usr/lib/hadoop-0.20# cd test
root@cloudera-vm:/usr/lib/hadoop-0.20/test# hadoop fs -get /user/teragen-inputTest/part-00001 .
root@cloudera-vm:/usr/lib/hadoop-0.20/test# ls
part-00001
root@cloudera-vm:/usr/lib/hadoop-0.20/test# vi part-00001
root@cloudera-vm:/usr/lib/hadoop-0.20/test# vi part-00001
root@cloudera-vm:/usr/lib/hadoop-0.20/test#
```

Dump of a MR Job

edureka!

```
13/08/03 00:58:40 WARN mapred.JobClient: Use GenericOptionsParser for parsing th
e arguments. Applications should implement Tool for the same.
13/08/03 00:58:40 INFO mapred.FileInputFormat: Total input paths to process : 1
13/08/03 00:58:40 INFO mapred.JobClient: Running job: job_201308022025_0003
13/08/03 00:58:41 INFO mapred.JobClient: map 0% reduce 0%
13/08/03 00:58:44 INFO mapred.JobClient: map 100% reduce 0%
13/08/03 00:58:51 INFO mapred.JobClient: map 100% reduce 11%
13/08/03 00:58:52 INFO mapred.JobClient: map 100% reduce 66%
13/08/03 00:58:59 INFO mapred.JobClient: map 100% reduce 100%
13/08/03 00:58:59 INFO mapred.JobClient: Job complete: job_201308022025_0003
13/08/03 00:58:59 INFO mapred.JobClient: Counters: 23
13/08/03 00:58:59 INFO mapred.JobClient: Job Counters
13/08/03 00:58:59 INFO mapred.JobClient: Launched reduce tasks=3
13/08/03 00:58:59 INFO mapred.JobClient: SLOTS_MILLIS_MAPS=4053
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all reduces wai
ting after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all maps waitin
g after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Launched map tasks=2
13/08/03 00:58:59 INFO mapred.JobClient: Data-local map tasks=2
13/08/03 00:58:59 INFO mapred.JobClient: SLOTS_MILLIS_REDUCES=23684
13/08/03 00:58:59 INFO mapred.JobClient: FileSystemCounters
13/08/03 00:58:59 INFO mapred.JobClient: FILE_BYTES_READ=81770
13/08/03 00:58:59 INFO mapred.JobClient: HDFS_BYTES_READ=136111
13/08/03 00:58:59 INFO mapred.JobClient: FILE_BYTES_WRITTEN=429317
13/08/03 00:58:59 INFO mapred.JobClient: HDFS_BYTES_WRITTEN=61194
13/08/03 00:58:59 INFO mapred.JobClient: Map-Reduce Framework
13/08/03 00:58:59 INFO mapred.JobClient: Reduce input groups=3586
13/08/03 00:58:59 INFO mapred.JobClient: Combine output records=4027
13/08/03 00:58:59 INFO mapred.JobClient: Map input records=2403
13/08/03 00:58:59 INFO mapred.JobClient: Reduce shuffle bytes=81788
13/08/03 00:58:59 INFO mapred.JobClient: Reduce output records=3586
13/08/03 00:58:59 INFO mapred.JobClient: Spilled Records=8054
13/08/03 00:58:59 INFO mapred.JobClient: Map output bytes=151013
13/08/03 00:58:59 INFO mapred.JobClient: Map input bytes=132663
13/08/03 00:58:59 INFO mapred.JobClient: Combine input records=11037
13/08/03 00:58:59 INFO mapred.JobClient: Map output records=11037
13/08/03 00:58:59 INFO mapred.JobClient: SPLIT_RAW_BYTES=146
13/08/03 00:58:59 INFO mapred.JobClient: Reduce input records=4027
```



Dump of a MR Job

```
13/08/03 00:58:40 WARN mapred.JobClient: Use GenericOptionsParser for parsing th
e arguments. Applications should implement Tool for the same.
13/08/03 00:58:40 INFO mapred.FileInputFormat: Total input paths to process : 1
13/08/03 00:58:40 INFO mapred.JobClient: Running job: job_201308022025_0003
13/08/03 00:58:41 INFO mapred.JobClient: map 0% reduce 0%
13/08/03 00:58:44 INFO mapred.JobClient: map 100% reduce 0%
13/08/03 00:58:51 INFO mapred.JobClient: map 100% reduce 11%
13/08/03 00:58:52 INFO mapred.JobClient: map 100% reduce 66%
13/08/03 00:58:59 INFO mapred.JobClient: map 100% reduce 100%
13/08/03 00:58:59 INFO mapred.JobClient: Job complete: job_201308022025_0003
13/08/03 00:58:59 INFO mapred.JobClient: Counters: 23
13/08/03 00:58:59 INFO mapred.JobClient: Job Counters
13/08/03 00:58:59 INFO mapred.JobClient: Launched reduce tasks=3
13/08/03 00:58:59 INFO mapred.JobClient: SLOTS_MILLIS_MAPS=4053
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all reduces wait
ing after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all maps waitin
g after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Launched map tasks=2
13/08/03 00:58:59 INFO mapred.JobClient: Data-local map tasks=2
13/08/03 00:58:59 INFO mapred.JobClient: SLOTS_MILLIS_REDUCES=23684
13/08/03 00:58:59 INFO mapred.JobClient: FileSystemCounters
13/08/03 00:58:59 INFO mapred.JobClient: FILE_BYTES_READ=81770
13/08/03 00:58:59 INFO mapred.JobClient: HDFS_BYTES_READ=136111
13/08/03 00:58:59 INFO mapred.JobClient: FILE_BYTES_WRITTEN=429317
13/08/03 00:58:59 INFO mapred.JobClient: HDFS_BYTES_WRITTEN=61194
13/08/03 00:58:59 INFO mapred.JobClient: Map-Reduce Framework
13/08/03 00:58:59 INFO mapred.JobClient: Reduce input groups=3586
13/08/03 00:58:59 INFO mapred.JobClient: Combine output records=4027 ↗ 2
13/08/03 00:58:59 INFO mapred.JobClient: Map input records=2403
13/08/03 00:58:59 INFO mapred.JobClient: Reduce shuffle bytes=81788
13/08/03 00:58:59 INFO mapred.JobClient: Reduce output records=3586 ↘ 1
13/08/03 00:58:59 INFO mapred.JobClient: Spilled Records=8054
13/08/03 00:58:59 INFO mapred.JobClient: Map output bytes=151013
13/08/03 00:58:59 INFO mapred.JobClient: Map input bytes=132663
13/08/03 00:58:59 INFO mapred.JobClient: Combine input records=11037 ↗ 2
13/08/03 00:58:59 INFO mapred.JobClient: Map output records=11037 ↘ 1
13/08/03 00:58:59 INFO mapred.JobClient: SPLIT_RAW_BYTES=146
13/08/03 00:58:59 INFO mapred.JobClient: Reduce input records=4027
```

Dump of a MR Job

edureka!

```
13/08/03 00:58:40 WARN mapred.JobClient: Use GenericOptionsParser for parsing th
e arguments. Applications should implement Tool for the same.
13/08/03 00:58:40 INFO mapred.FileInputFormat: Total input paths to process : 1
13/08/03 00:58:40 INFO mapred.JobClient: Running job: job_201308022025_0003
13/08/03 00:58:41 INFO mapred.JobClient: map 0% reduce 0%
13/08/03 00:58:44 INFO mapred.JobClient: map 100% reduce 0%
13/08/03 00:58:51 INFO mapred.JobClient: map 100% reduce 11%
13/08/03 00:58:52 INFO mapred.JobClient: map 100% reduce 66%
13/08/03 00:58:59 INFO mapred.JobClient: map 100% reduce 100%
13/08/03 00:58:59 INFO mapred.JobClient: Job complete: job_201308022025_0003
13/08/03 00:58:59 INFO mapred.JobClient: Counters: 23
13/08/03 00:58:59 INFO mapred.JobClient: Job Counters
13/08/03 00:58:59 INFO mapred.JobClient: Launched reduce tasks=3
13/08/03 00:58:59 INFO mapred.JobClient: SLOTS_MILLIS_MAPS=4053
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all reduces wai
ting after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all maps waitin
g after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Launched map tasks=2
13/08/03 00:58:59 INFO mapred.JobClient: Data-local map tasks=2
13/08/03 00:58:59 INFO mapred.JobClient: SLOTS_MILLIS_REDUCES=23684
13/08/03 00:58:59 INFO mapred.JobClient: FileSystemCounters
13/08/03 00:58:59 INFO mapred.JobClient: FILE_BYTES_READ=81770
13/08/03 00:58:59 INFO mapred.JobClient: HDFS_BYTES_READ=136111
13/08/03 00:58:59 INFO mapred.JobClient: FILE_BYTES_WRITTEN=429317
13/08/03 00:58:59 INFO mapred.JobClient: HDFS_BYTES_WRITTEN=61194
13/08/03 00:58:59 INFO mapred.JobClient: Map-Reduce Framework
13/08/03 00:58:59 INFO mapred.JobClient: Reduce input groups=3586
13/08/03 00:58:59 INFO mapred.JobClient: Combine output records=4027 2
13/08/03 00:58:59 INFO mapred.JobClient: Map input records=2403
13/08/03 00:58:59 INFO mapred.JobClient: Reduce shuffle bytes=81788
13/08/03 00:58:59 INFO mapred.JobClient: Reduce output records=3586 2
13/08/03 00:58:59 INFO mapred.JobClient: Spilled Records=8054
13/08/03 00:58:59 INFO mapred.JobClient: Map output bytes=151013
13/08/03 00:58:59 INFO mapred.JobClient: Map input bytes=132663
13/08/03 00:58:59 INFO mapred.JobClient: Combine input records=11037 1
13/08/03 00:58:59 INFO mapred.JobClient: Map output records=11037
13/08/03 00:58:59 INFO mapred.JobClient: SPLIT_RAW_BYTES=146
13/08/03 00:58:59 INFO mapred.JobClient: Reduce input records=4027 3
```

Dump of a MR Job

edureka!

```
13/08/03 00:58:40 WARN mapred.JobClient: Use GenericOptionsParser for parsing th
e arguments. Applications should implement Tool for the same.
13/08/03 00:58:40 INFO mapred.FileInputFormat: Total input paths to process : 1
13/08/03 00:58:40 INFO mapred.JobClient: Running job: job_201308022025_0003
13/08/03 00:58:41 INFO mapred.JobClient: map 0% reduce 0%
13/08/03 00:58:44 INFO mapred.JobClient: map 100% reduce 0%
13/08/03 00:58:51 INFO mapred.JobClient: map 100% reduce 11%
13/08/03 00:58:52 INFO mapred.JobClient: map 100% reduce 66%
13/08/03 00:58:59 INFO mapred.JobClient: map 100% reduce 100%
13/08/03 00:58:59 INFO mapred.JobClient: Job complete: job_201308022025_0003
13/08/03 00:58:59 INFO mapred.JobClient: Counters: 23
13/08/03 00:58:59 INFO mapred.JobClient: Job Counters
13/08/03 00:58:59 INFO mapred.JobClient: Launched reduce tasks=3
13/08/03 00:58:59 INFO mapred.JobClient: SLOTS_MILLIS_MAPS=4053
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all reduces wai
ting after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all maps waitin
g after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Launched map tasks=2
13/08/03 00:58:59 INFO mapred.JobClient: Data-local map tasks=2
13/08/03 00:58:59 INFO mapred.JobClient: SLOTS_MILLIS_REDUCES=23684
13/08/03 00:58:59 INFO mapred.JobClient: FileSystemCounters
13/08/03 00:58:59 INFO mapred.JobClient: FILE_BYTES_READ=81770
13/08/03 00:58:59 INFO mapred.JobClient: HDFS_BYTES_READ=136111
13/08/03 00:58:59 INFO mapred.JobClient: FILE_BYTES_WRITTEN=429317
13/08/03 00:58:59 INFO mapred.JobClient: HDFS_BYTES_WRITTEN=61194
13/08/03 00:58:59 INFO mapred.JobClient: Map-Reduce Framework
13/08/03 00:58:59 INFO mapred.JobClient: Reduce input groups=3586
13/08/03 00:58:59 INFO mapred.JobClient: Combine output records=4027 4
13/08/03 00:58:59 INFO mapred.JobClient: Map input records=2403
13/08/03 00:58:59 INFO mapred.JobClient: Reduce shuffle bytes=81788
13/08/03 00:58:59 INFO mapred.JobClient: Reduce output records=3586 2
13/08/03 00:58:59 INFO mapred.JobClient: Spilled Records=8054
13/08/03 00:58:59 INFO mapred.JobClient: Map output bytes=151013
13/08/03 00:58:59 INFO mapred.JobClient: Map input bytes=132663
13/08/03 00:58:59 INFO mapred.JobClient: Combine input records=11037 1
13/08/03 00:58:59 INFO mapred.JobClient: Map output records=11037
13/08/03 00:58:59 INFO mapred.JobClient: SPLIT_RAW_BYTES=146
13/08/03 00:58:59 INFO mapred.JobClient: Reduce input records=4027 3
```

The output of a MR job will be stored on HDFS:

- TRUE
- FALSE



True. It is stored in different part files for eg – part-m-00000, part-m-00001 and so on. The part files are created on the basis of the block size.



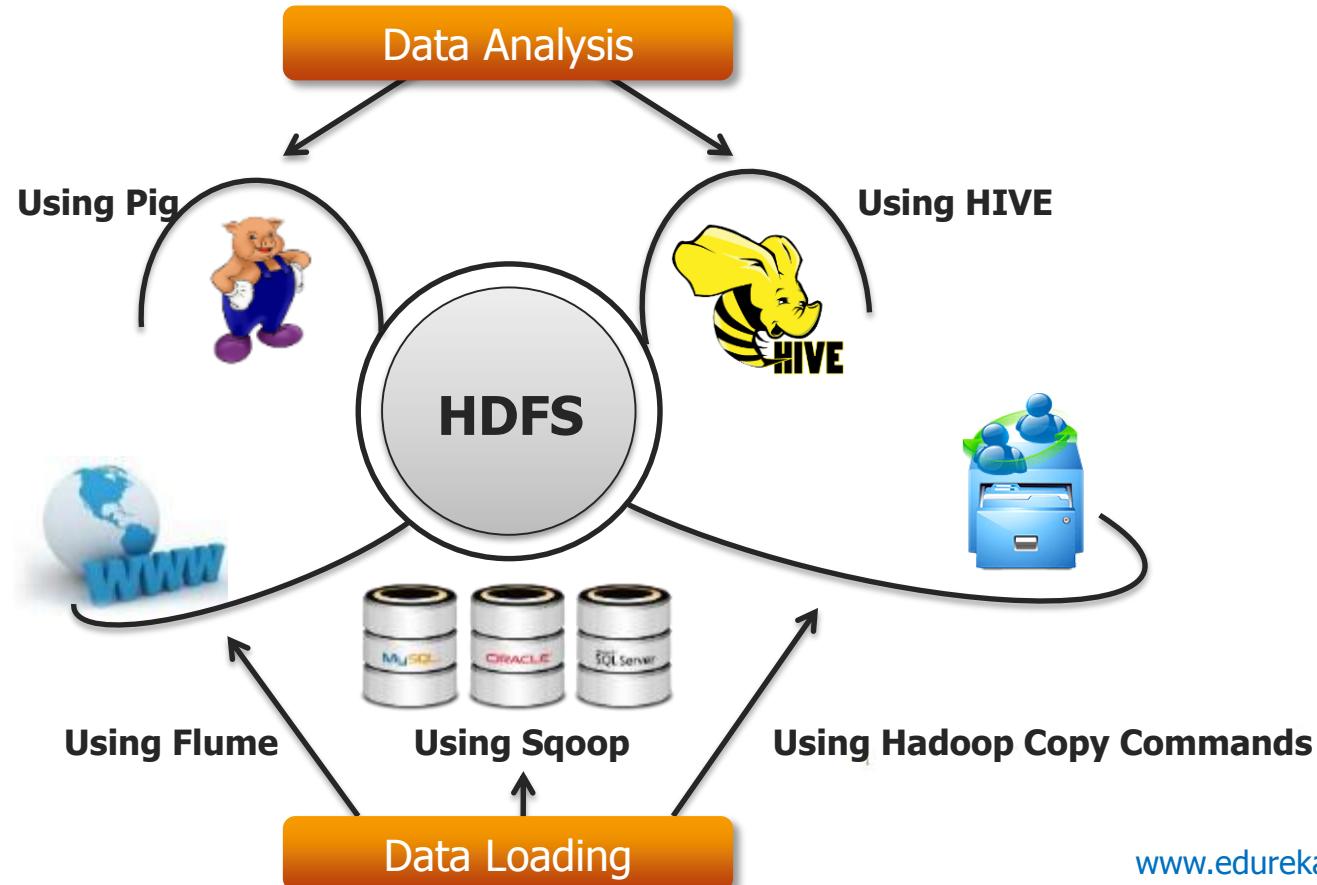
To run MR job data should be present on HDFS:

- TRUE
- FALSE





True. In order to process data in parallel it is necessary that it is present on HDFS so that MR can work on chunks of data in parallel.



put: Copy file(s) from local file system to destination file system. It can also read from “stdin” and writes to destination file system.

```
hadoop dfs –put weather.txt hdfs://<target Namenode>
```

copyFromLocal: Similar to “put” command, except that the source is restricted to a local file reference.

```
hadoop dfs –copyFromLocal weather.txt hdfs://<target Namenode>
```

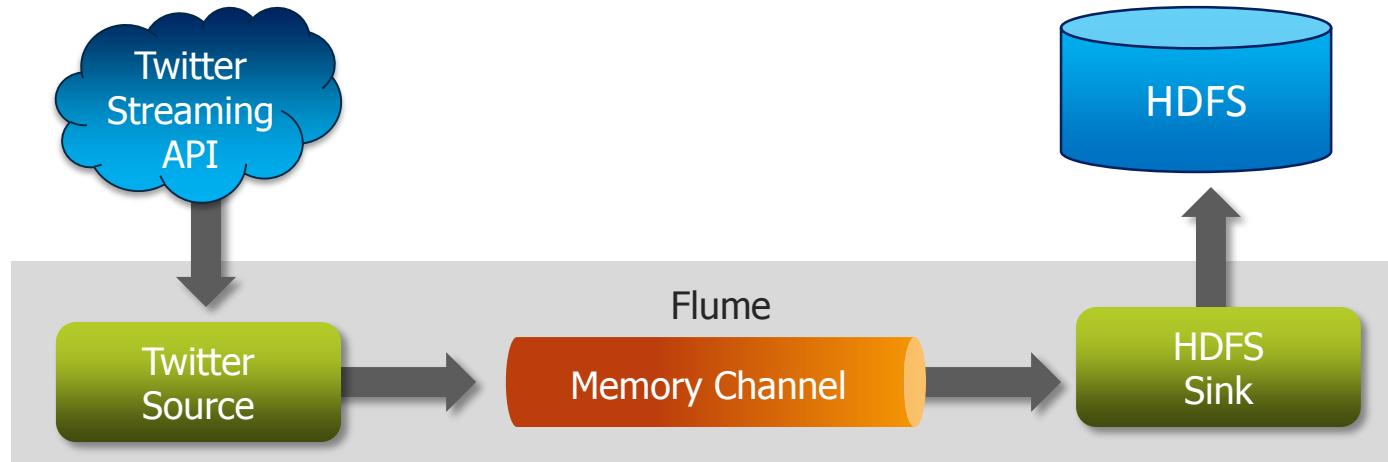
distcp: Distributed Copy to move data between clusters, used for backup and recovery

```
hadoop distcp hdfs://<source NN> hdfs://<target NN>
```

Data Loading Using Flume

edureka!

Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of streaming event data.



Apache Sqoop (TM) is a tool designed for efficiently transferring bulk data between [Apache Hadoop](#) and structured data stores such as relational databases.

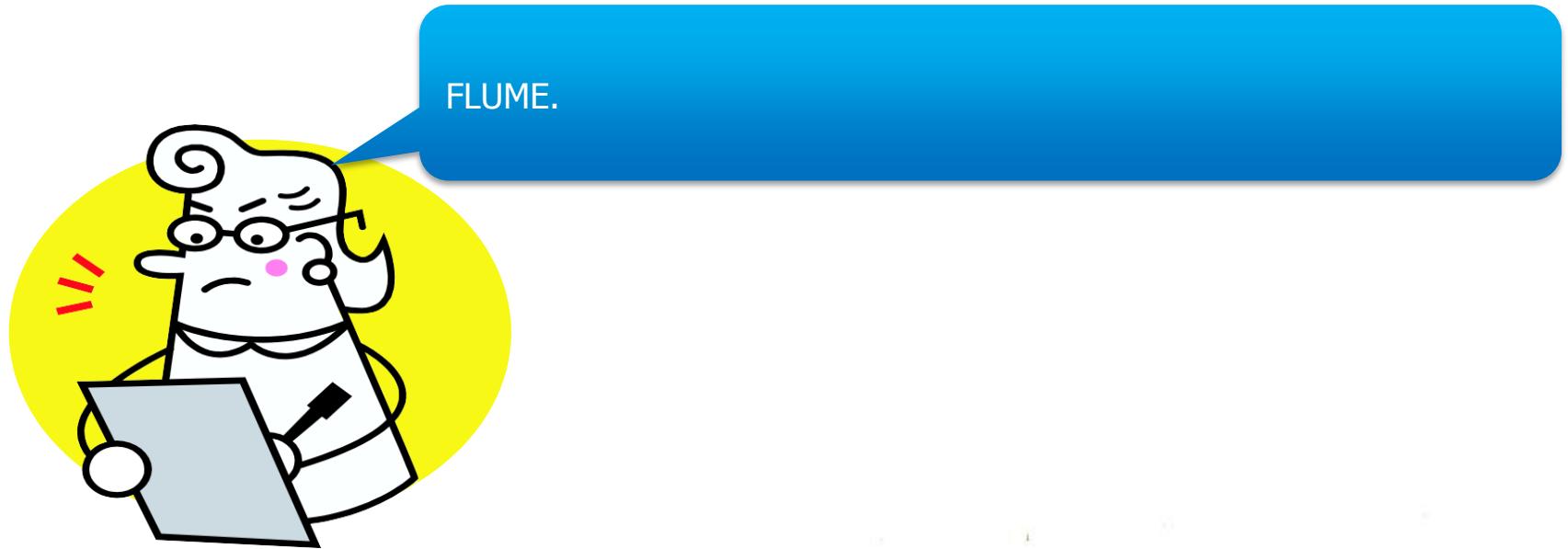
- ✓ Imports individual tables or entire databases to HDFS.
- ✓ Generates Java classes to allow you to interact with your imported data.
- ✓ Provides the ability to import from SQL databases straight into your Hive data warehouse.



Your website is hosting a group of more than 300 sub-websites. You want to have an analytics on the shopping patterns of different visitors? What is the best way to collect those information from the weblogs?

- SQOOP
- FLUME



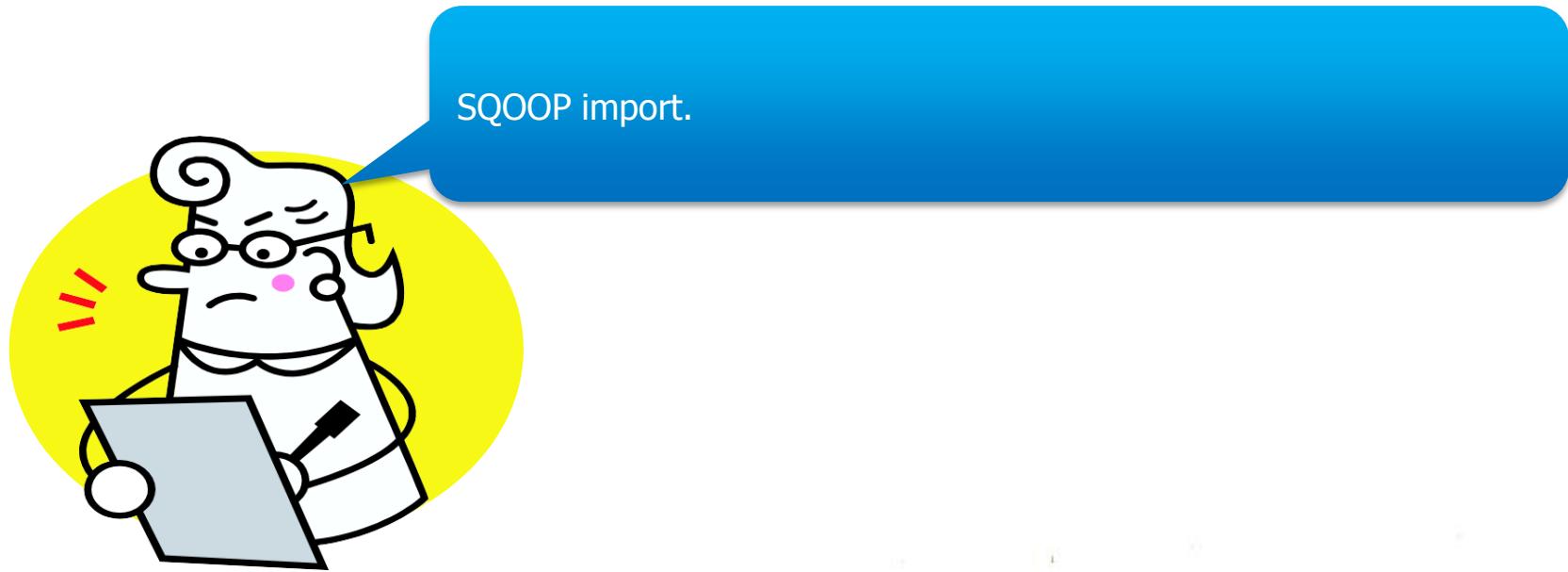




You want to join data collected from two sources. One source of data collected from a big database of call records is already available in HDFS. The another source of data is available in a database table.

The best way to move that data in HDFS is:

- SQOOP import
- PIG script
- Hive Query



Use Case – How do I find out the best?

MY SUBREDDITS ▾ FRONT - ALL - RANDOM | PICS - FUNNY - GAMING - ASKREDDIT - WORLDNEWS - NEWS - VIDEOS - IAMA - TODAYILEARNED - WTF - AWW - TECHNOLOGY - ADVICEANIMALS - SCIENCE - MUSIC - MOVIES - BES* MORE »

want to join? login or register in seconds | English

reddit hot new rising controversial top gilded wiki

were gamers saving kids! Will you join us? (extra-life.org)
submitted 9 hours ago by ExtraLife4Kids to funny
1 comment share

Help me break this thing, bro. (i.imgur.com)
submitted 9 hours ago by tuberjFAGmod to funny
2905 522 comments share

And I always make sure I say she's hot... (livememe.com)
submitted 7 hours ago by jerrytodd to AdviceAnimals
2014 411 comments share

Facebook no longer lets users hide from search (sfgate.com)
submitted 3 hours ago by GainSeAya to technology
800 179 comments share

TIL that in California and 3 other US states, "Ladies' Night" are against the law because they are gender discrimination (en.wikipedia.org)
submitted 9 hours ago by Gertonification to todayilearned
2465 1028 comments share

Boy, 15, kills himself after 'facing expulsion and being put on sex offender registry' for STREAKING at high school football game Misleading Title! (engineeringevil.com)
submitted 10 hours ago by gn3xu5 to news
3122 3045 comments share

Proud Owner of an Award-Winning Butthole (self.IAmA)
submitted 9 hours ago by AsaAkira1 to IAmA
2118 2772 comments share

A crowd gathering in Union Square caught my attention. The kid won. (imgur.com)
submitted 11 hours ago by worldsarmy to pics
3079 1262 comments share

I'm tiny and new and cuter than you. (imgur.com)
submitted 10 hours ago by hamuchannic to aww
2386 273 comments share

sponsored link what's this?

search reddit

username password

remember me reset password login

Submit a new link

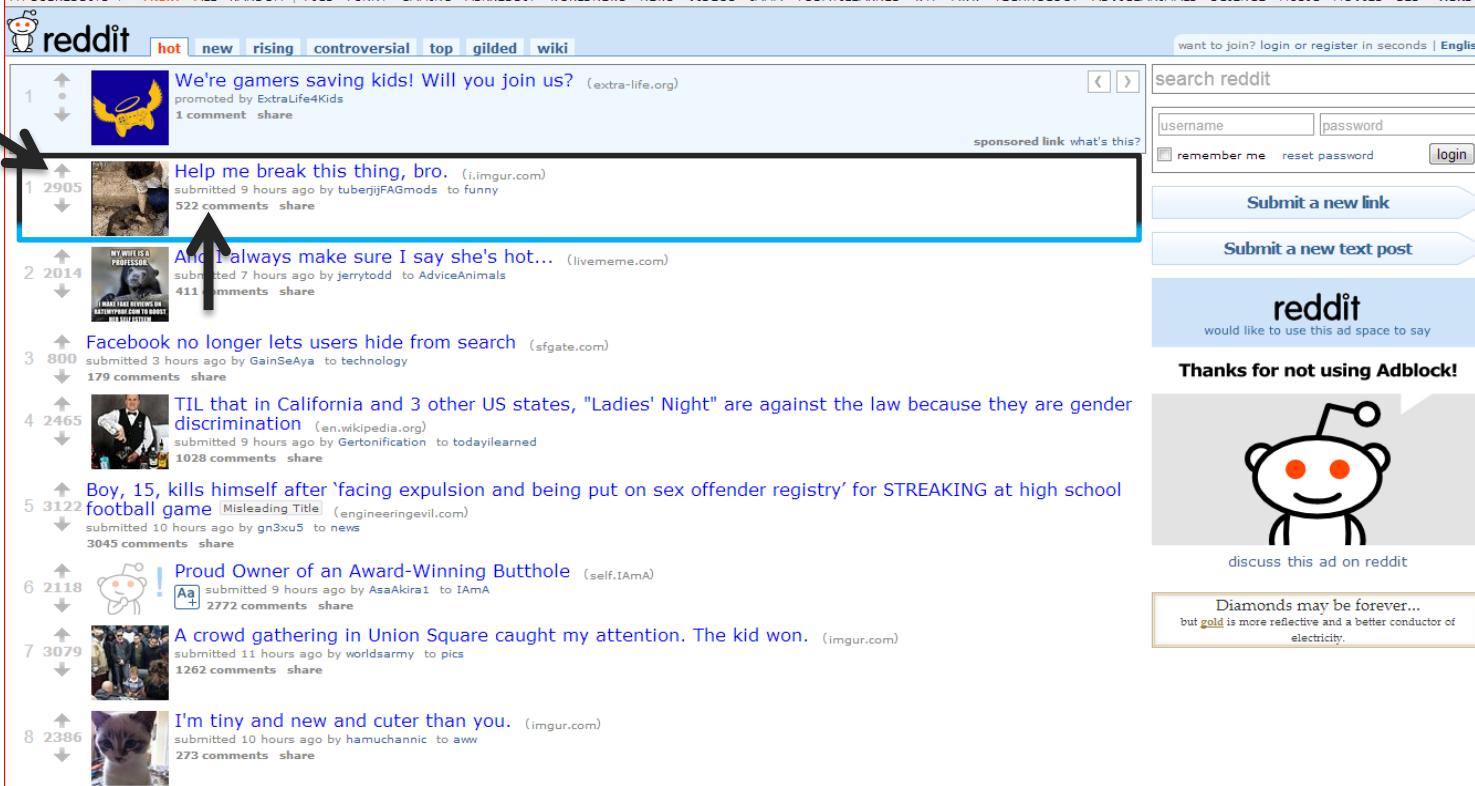
Submit a new text post

reddit would like to use this ad space to say

Thanks for not using Adblock!

discuss this ad on reddit

Diamonds may be forever... but gold is more reflective and a better conductor of electricity.

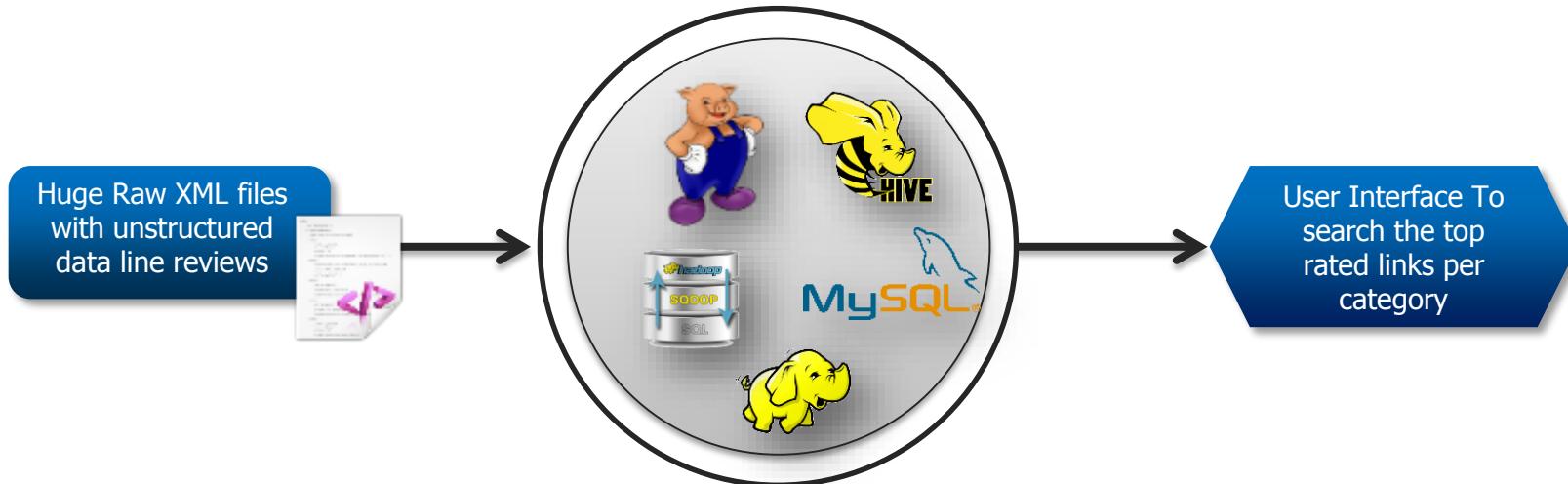


Use Case – Data to Analyse

```
<?xml version="1.0" encoding="UTF-8"?>
<documents>
  <document>
    <hash>0c539040d03f4ba7edb386b3d5b7c8c8</hash>
    <url>http://www.bemboszoo.com/</url>
    <category>Top/Reference/Education/Early_Childhood/Subjects/Language_Arts/Alphabet</category>
    <usercount>724</usercount>
      <tag>
        <name><! [CDATA[fun]]></name>
        <count>114</count>
      </tag>
      <tag>
        <name><! [CDATA[kids]]></name>
        <count>69</count>
      </tag>
      <tag>
        <name><! [CDATA[alphabet]]></name>
        <count>66</count>
      </tag>
      <tag>
        <name><! [CDATA[children]]></name>
        <count>55</count>
      </tag>
    </tags>
    <reviews>
      <review><! [CDATA["This high concept abcdeRARY, the picture book debut for deVicq Cumptich, should delight collectors of Publishers Weekly!]]></review>
      <review><! [CDATA[Really neat flash site making animal art out of the letters of the alphabet.]]></review>
      <review><! [CDATA[I love this site!]]]></review>
      <review><! [CDATA[Since I'm forever young LOL I LOVE this and it's even better in the flash version!!!!!! :)]]></review>
      <review><! [CDATA[Cool once you get past the illegible splash page.]]></review>
      <review><! [CDATA[Silly, yet clever and funny. Probably helps if you have some interest in typography, and know that E]]></review>
      <review><! [CDATA[Clever!]]></review>
    </reviews>
  </document>
</documents>
```

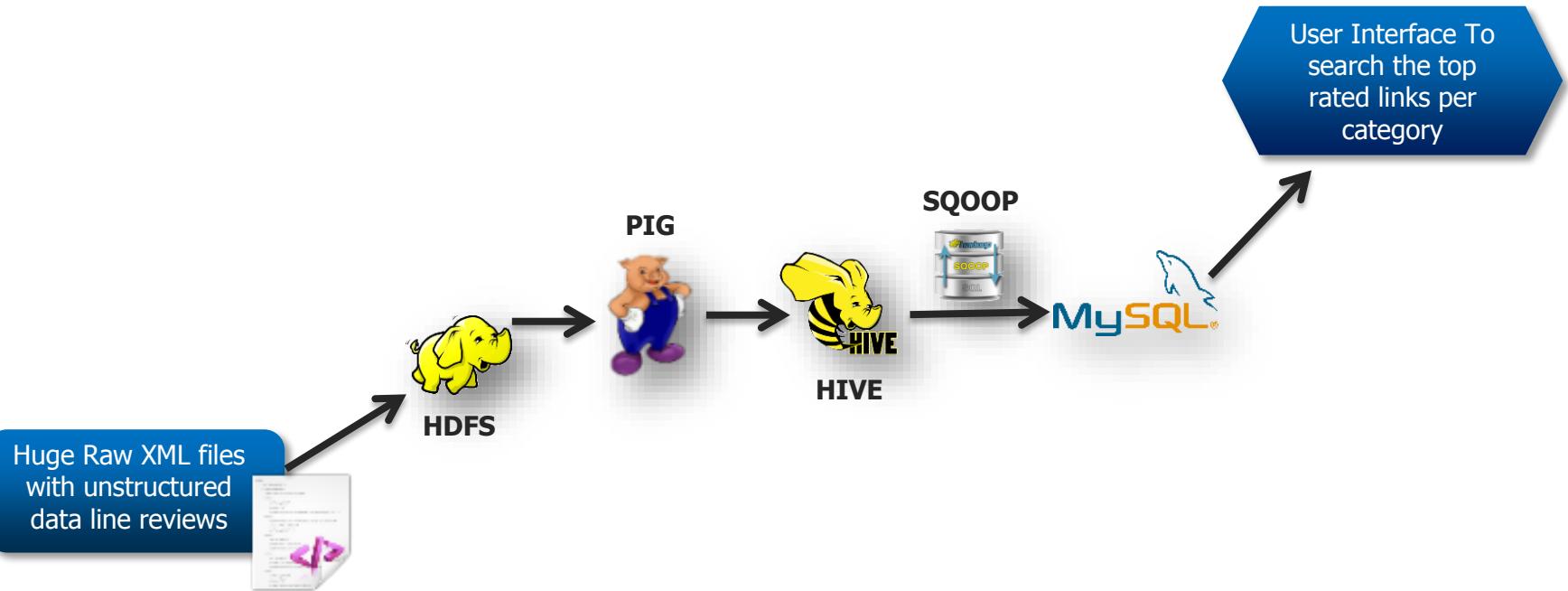
Abstract Flow Diagram

edureka!



Flow Diagram

edureka!



Demo for Data Load

- ✓ Hadoop Cluster Setup
http://hadoop.apache.org/docs/r0.19.1/cluster_setup.html
- ✓ Hadoop on Amazon AWS ec2
<http://www.edureka.in/blog/install-apache-hadoop-cluster/>
- ✓ Hadoop Hardware Selection
<http://blog.cloudera.com/blog/2013/08/how-to-select-the-right-hardware-for-your-new-hadoop-cluster/>
http://docs.hortonworks.com/HDPDocuments/HDP1/HDP-Win-1.3.0/bk_cluster-planning-guide/content/ch_hardware-recommendations.html
- ✓ Hadoop Cluster Configuration
<http://www.edureka.in/blog/hadoop-cluster-configuration-files/>

- ✓ MapReduce Job execution

<http://www.edureka.in/blog/anatomy-of-a-mapreduce-job-in-apache-hadoop/>

- ✓ Add/Remove Nodes in a Cluster

<http://www.edureka.in/blog/commissioning-and-decommissioning-nodes-in-a-hadoop-cluster/>

- ✓ Secondary Namenode

https://hadoop.apache.org/docs/r1.2.1/hdfs_user_guide.html#Secondary+NameNode

- ✓ Setup the Hadoop development environment using the documents present in the LMS.
 - ✓ Flume Set-up on Cloudera
 - ✓ SQOOP Set-up on Cloudera
 - ✓ Refresh your Java Skills using Java Essential for Hadoop Tutorial
- ✓ Attempt the Module-2 Assignments present in the LMS.
- ✓ Review the Interview Questions for setting up hadoop cluster
<http://www.edureka.in/blog/hadoop-interview-questions-hadoop-cluster/>

edureka!

Thank You

See You in Class Next Week