

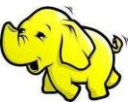
A word cloud visualization featuring various terms related to Big Data technologies. The most prominent words are "Big Data" and "PetaBytes". Other significant words include "MapReduce", "HBase", "HiveQL", "Fume", "MRUnit", "Zookeeper", "YARN", "HDFS", "TaskTracker", "Mappers", "Partitioners", "MongoDB", "JobTracker", "Localization", "Oozie", "Hive", "Sqoop", "Reducers", "File System", "Cassandra", "Combiners", "Distributed", "API", and "Pig". The words are arranged in a dense, overlapping manner, with colors ranging from blue and green to red and yellow.

Course Content

- **Module1**
 - BigData Introduction
- **Module2**
 - HDFS Installation and Commands
- **Module 3**
 - Deployment Modes
- **Module 4**
 - Mapreduce
- **Module 5**
 - Advance Mapreduce - Part1
 - Join and Counters
- **Module 6**
 - Advance MapReduce - Part 2
 - Custom Input Formats and MRUnit
- **Module 7**
 - Pig and Pig Latin
- **Module 8**
 - Hive, Hive QL and Hlve Architecture
- **Module 9**
 - Advance Hive QL, Rollups and Custom Functions

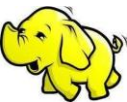


- **Module 10**
 - Flume
- **Module 11**
 - Sqoop
- **Module 12**
 - Oozie
- **Module 13**
 - NoSQL Databases
- **Module 14**
 - MongoDB and Cassandra
- **Module 15**
 - Hbase and Advanced Hbase
- **Module 16**
 - Zookeeper
- **Module 17**
 - Hadoop 2.0
 - HA and YARN and MRV2
- **Module 18**
 - Project

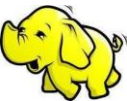
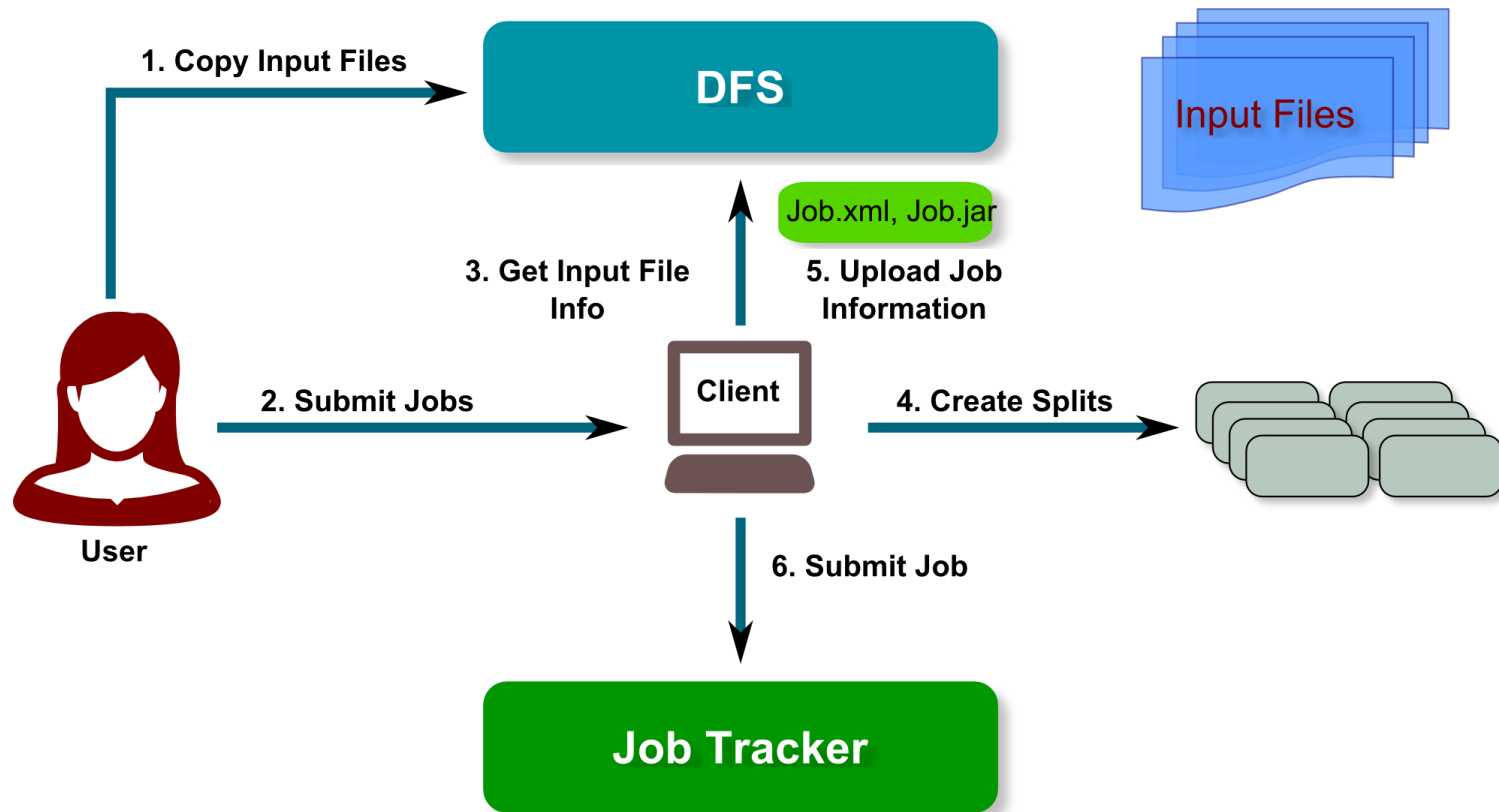


Agenda for the day

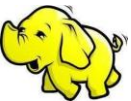
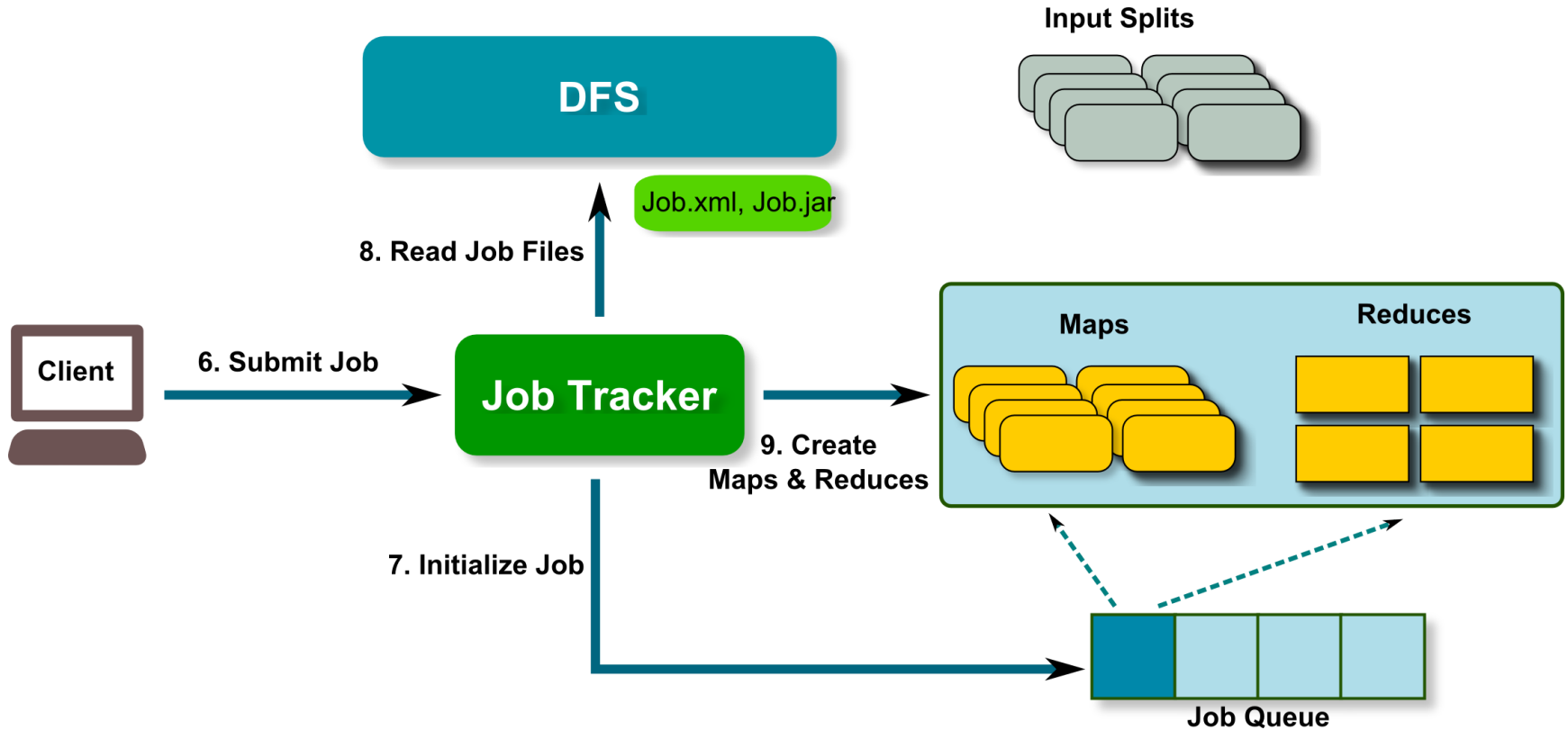
1. Understanding Job Tracker
2. Hadoop Modes
3. Hadoop Terminal Commands
4. Production Hadoop Clusters
5. Cluster Configuration
6. Hadoop Configuration Files
7. Recovery
8. MapReduce in Action



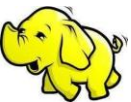
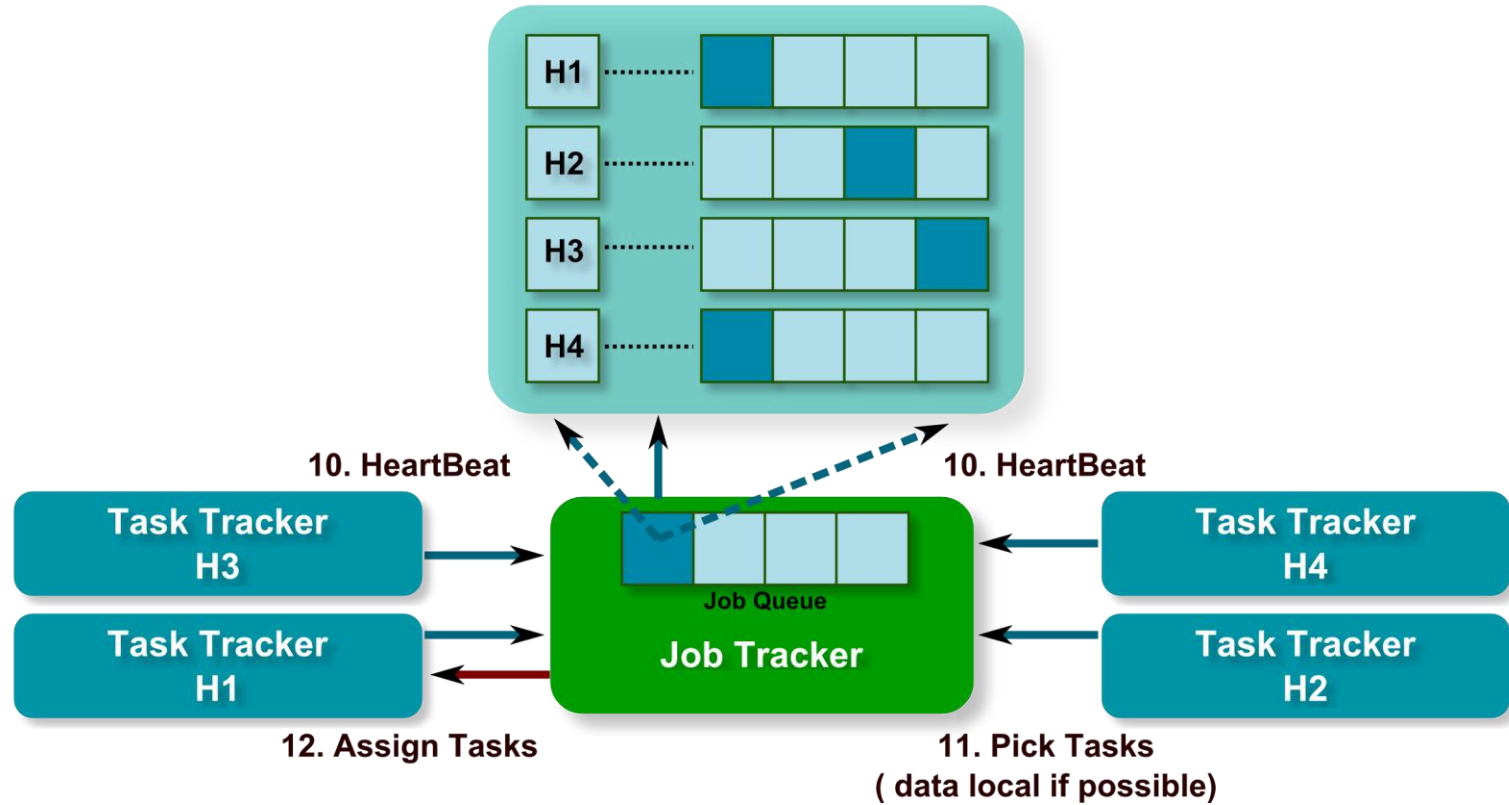
Job Tracker



Job Tracker

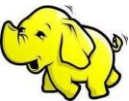


Job Tracker



Job Tracker

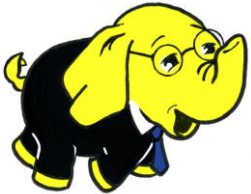
- MapReduce Master : delegating jobs to task tracker
 - Client submits jobs to job tracker , jobs are kept in queue
 - FIFO Scheduler
- Capacity Scheduler
- Job Tracker Determines the location of data through Name Node
- Job Tracker determines available task tracker (prefers the slots near to the data)
- Job Tracker submits the work to Task Tracker
- Task Tracker monitors it and send update to Job Tracker
- After completion Job Tracker Updates its Status
- Job Tracker is a single point of failure



Hadoop Wants to Know?

Which of these is responsible to assign a task to TaskTracker ?

- a) Namenode
- b) Jobtracker
- c) Secondary NameNode
- d) Data Node



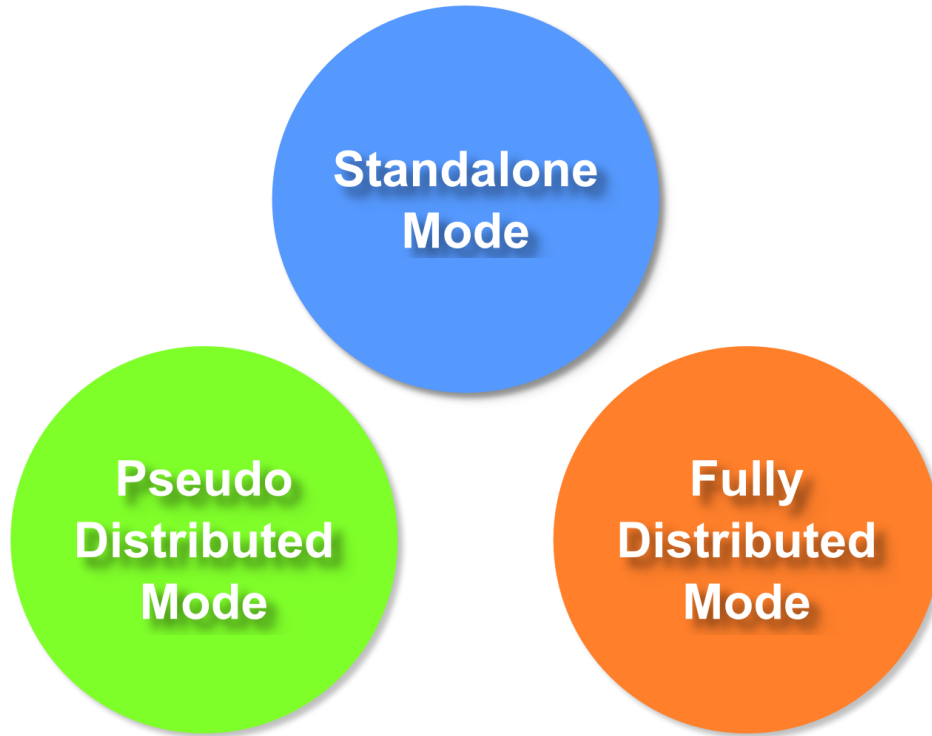
Hadoop Wants to Know?

Which of these is responsible for executing a task ?

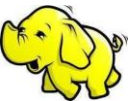
- a) Namenode
- b) Jobtracker
- c) TaskTracker
- d) Data Node



Hadoop 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

```
hadoopjob@ubuntu:~$ hadoop
Warning: $HADOOP_HOME is deprecated.

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
  oiv                   apply the offline fsimage viewer to an fsimage
  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
  historyserver         run job history servers as a standalone daemon
  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
  distcp2 <srcurl> <desturl> DistCp version 2
  archive -archiveName NAME -p <parent path> <src>* <dest> create a hadoop archi
ve
  classpath             prints the class path needed to get the
                        Hadoop jar and the required libraries
  daemonlog             get/set the log level for each daemon
  or
  CLASSNAME             run the class named CLASSNAME
Most commands print help when invoked w/o parameters.
hadoopjob@ubuntu:~$
```



Terminal Commands

Listing of Files present on HDFS

```
hadoopjob@ubuntu:~$ hadoop fs -ls /  
Warning: $HADOOP_HOME is deprecated.  
  
Found 4 items  
drwxr-xr-x - stratapps supergroup          0 2014-02-12 09:56 /data  
drwxr-xr-x - stratapps supergroup          0 2014-02-12 10:19 /datasets  
drwxr-xr-x - stratapps supergroup          0 2014-02-12 16:47 /tmp  
drwxr-xr-x - stratapps supergroup          0 2014-02-12 15:53 /user
```

Listing of files present in bin Directory

```
hadoopjob@ubuntu:~$ cd /usr/local/hadoop/bin  
hadoopjob@ubuntu:/usr/local/hadoop/bin$ ls  
hadoop          hadoop-daemons.sh  start-all.sh      start-jobhistoryserver.sh  stop-balancer.sh  stop-mapred.sh  
hadoop-config.sh rcc                 start-balancer.sh  start-mapred.sh           stop-dfs.sh       task-controller  
hadoop-daemon.sh slaves.sh           start-dfs.sh       stop-all.sh              stop-jobhistoryserver.sh  
hadoopjob@ubuntu:/usr/local/hadoop/bin$
```



Web UI URL's

Name Node Status : <http://localhost:50070/dfshealth.jsp>

Job Tracker Status : <http://localhost:50030/jobtracker.jsp>

Task Tracker Status : <http://localhost:50060/tasktracker.jsp>

Data Block Scanner Report : <http://localhost:50075/blockScannerReport>



Hadoop Wants to Know?

Which of these is used for uploading data from local to HDFS?

- a) Hadoop dfs -ls /input
- b) Hadoop dfs -get /input/data.json
- c) hadoop dfs -mkdir /input/data
- d) Hadoop dfs -copyFromLocal /input/data/data1.json /input/data/



Real Life Hadoop Implementations

- **EBay**

- 532 nodes cluster (8 * 532 cores, 5.3PB).
- Heavy usage of Java MapReduce, Apache Pig, Apache Hive, Apache HBase
- Using it for Search optimization and Research.

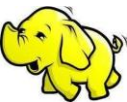


- **Facebook**

- Uses Apache Hadoop to store copies of internal log and dimension data sources and use it as source for reporting/analytics and machine learning.
- Currently FB 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.
- FB are heavy users of both streaming as well as the Java APIs. Have built a higher level data warehousing framework using these features called Hive



(see the <http://hadoop.apache.org/hive/>). Have also developed a FUSE implementation over HDFS.



Real Life Hadoop Implementations



- **Spotify**

- Uses Apache Hadoop for content generation, data aggregation, reporting and analysis
- 690 node cluster = 8280 physical cores, 38TB RAM, 28 PB storage (read more about Hadoop issues while growing fast: Hadoop Adventures At Spotify)
- 7,500 daily Hadoop jobs (scheduled by Luigi, our home-grown and recently open-sourced job scheduler - <https://vimeo.com/63435580>)

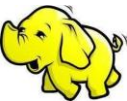
<http://www.slideshare.net/AdamKawa/hadoop-adventures-at-spotify-strata-conference-hadoop-world-2013>

<http://files.meetup.com/5139282/SHUG%201%20-%20Hadoop%20at%20Spotify.pdf>

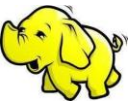
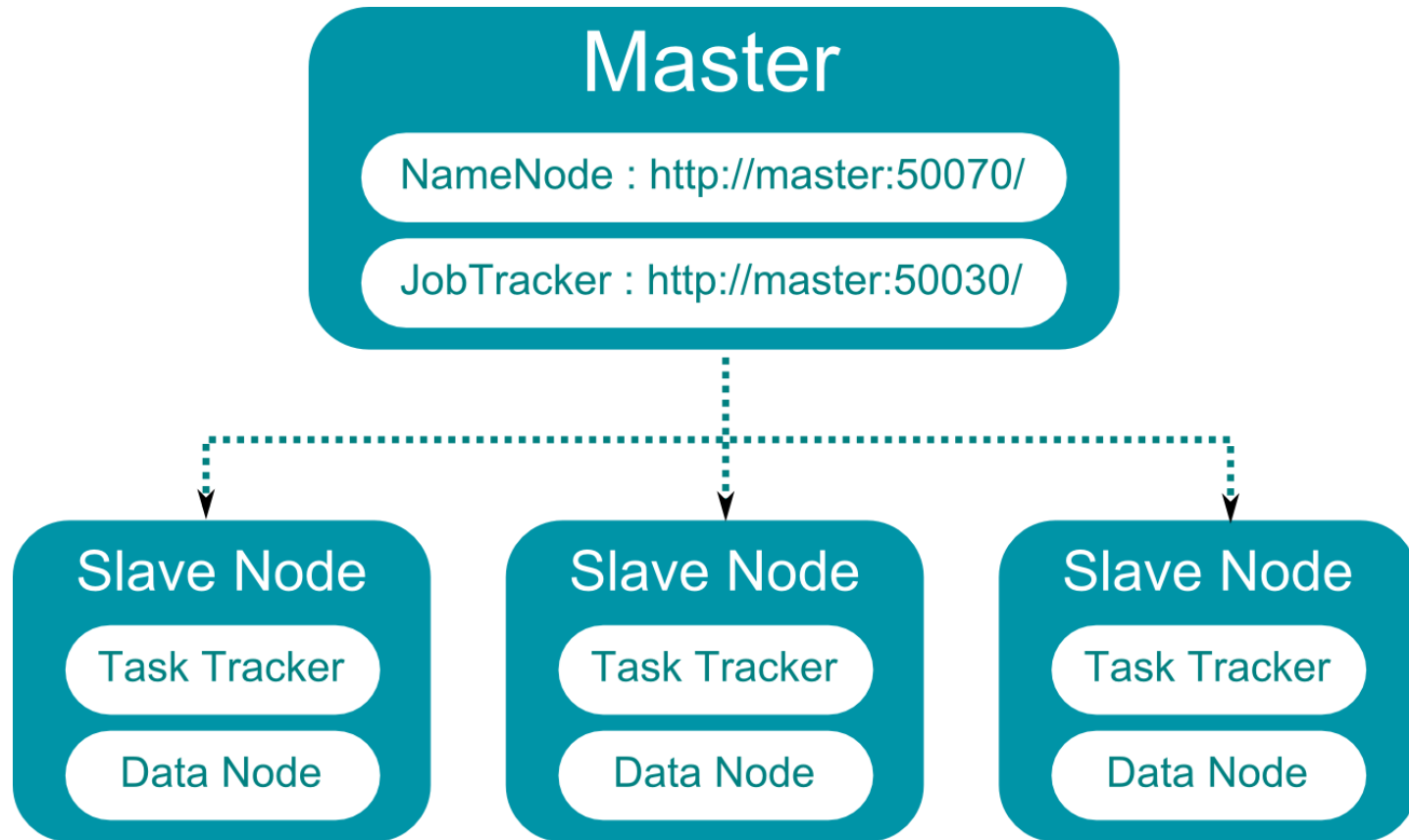


- **Last.fm**

- 100 nodes
- Dual quad-core Xeon L5520 @ 2.27GHz & L5630 @ 2.13GHz , 24GB RAM, 8TB(4x2TB)/node
- Used for charts calculation, royalty reporting, log analysis, A/B testing, dataset merging
- Also used for large scale audio feature analysis over millions of tracks

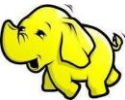


Sample Cluster Configuration



Hadoop Configuration Files

Configuration Files	Format	Description
hadoop-env.sh	Base Script	Environment variables that are used in the scripts to run Hadoop.
core-site.xml	Hadoop Configuration XML	Configuration settings for Hadoop Core such as I/O settings that are common to HDFS and MapReduce.
hdfs-site.xml	Hadoop Configuration XML	Configuration settings for HDFS daemons, the namenode, the secondary namenode and the data nodes.
mapred-site.xml	Hadoop Configuration XML	Configuration settings for MapReduce daemons : the job-tracker and the task-trackers.
masters	Plain Text	A list of machines (one per line) that each run a secondary namenode.
slave	Plain Text	A list of machines (one per line) that each run a datanode and a task-tracker.
hadoop-metric.properties	Java Properties	Properties for controlling how metrics are published in Hadoop.
log4j.properties	Java Properties	Properties for system log files, the namenode audit log and the task log for the task-tracker child process.



Hadoop Configuration Files



capacity-scheduler.xml



configuration.xsl



core-site.xml



fair-scheduler.xml



hadoop-env.sh



hadoop-metrics2.properties



hadoop-policy.xml



hdfs-site.xml



log4j.properties



mapred-queue-acls.xml



mapred-site.xml



masters



slaves



ssl-client.xml.example



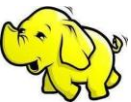
ssl-server.xml.example



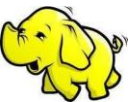
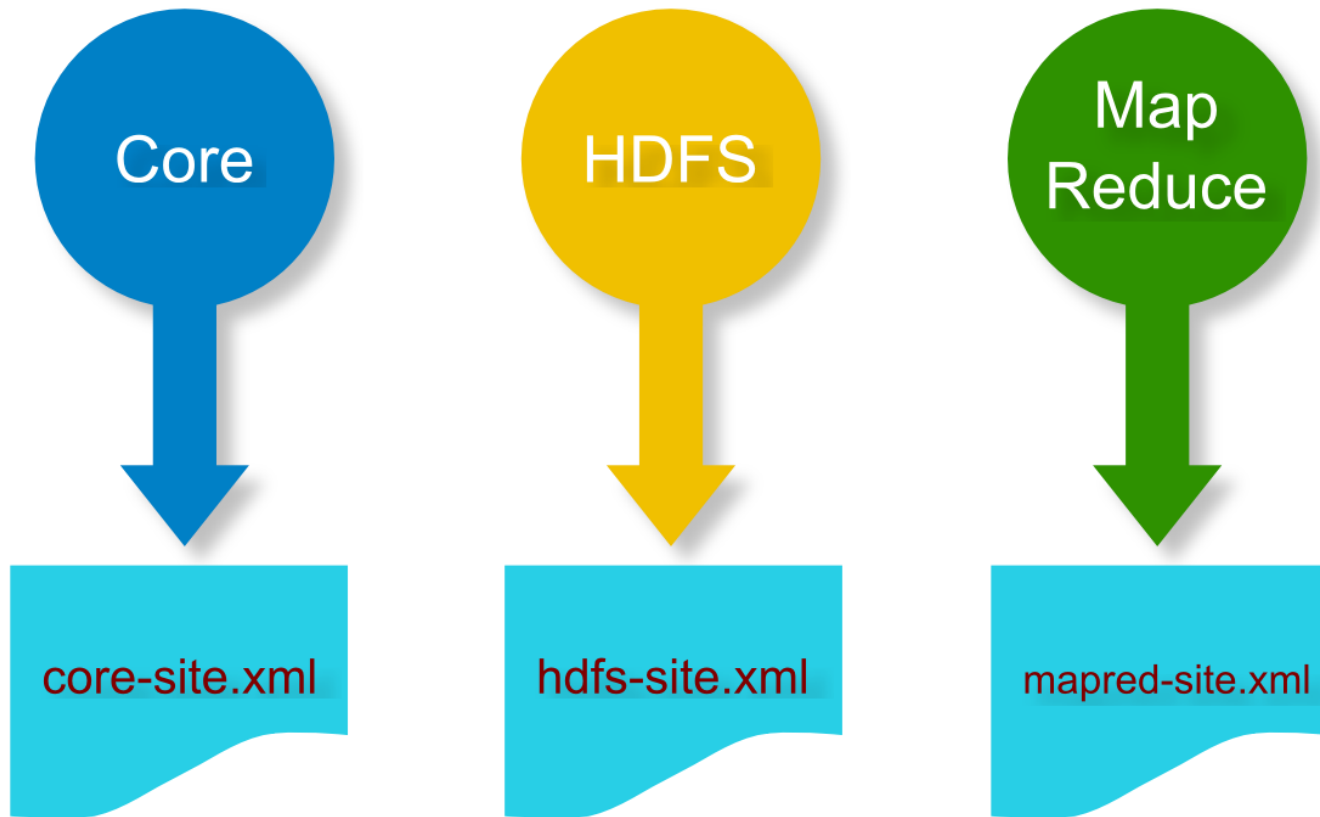
task-log4j.properties



taskcontroller.cfg

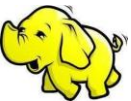


DD for each Component



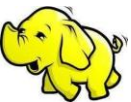
Core-site.xml and hdfs-site.xml

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>http://localhost:8020</value>
</property>	</property>
</configuration>	</configuration>



Defining HDFS details in hdfs-site.xml

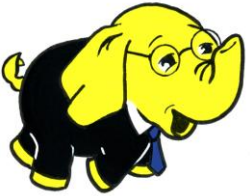
Property	Value	Description
dfs.data.dir	<code><value> /disk1/hdfs/data, /disk2/hdfs/data </value></code>	A list of directories where the datanode stores blocks. Each block is stored in only one of these directories. <code>\${hadoop.tmp.dir}/dfs/data</code>
fs.checkpoint.dir	<code><value> /disk1/hdfs/namesecondary,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



Hadoop Wants to Know?

Can you provide two different data directory locations for a single data node?

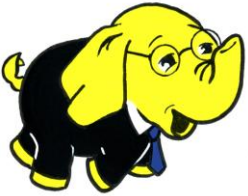
- a) True
- b) False



Hadoop Wants to Know?

Default Replication Factor in HDFS ?

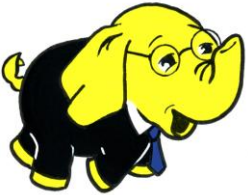
- a) 2
- b) 3
- c) 4
- d) None



Hadoop Wants to Know?

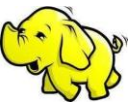
Which of following daemons would be running on Slave Machine ?

- a) DataNode
- b) Secondary NameNode
- c) NameNode
- d) TaskTracker
- e) JobTracker



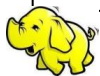
mapred-site.xml

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



Defining mapred-site.xml

Property	Value	Description
mapred.job.tracker	<code><value> localhost:8021 </value></code>	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	<code>\${hadoop.tmp.dir}/mapred/local</code>	A list of directories where MapReduce stores intermediate data for jobs. The data is cleared out when the job ends.
mapred.system.dir	<code>\${hadoop.tmp.dir}/mapred/system</code>	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.tasktracker. map/reducer .tasks.maximum	2	The number of map/reducer tasks that may be run on a tasktracker at any one time



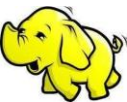
All Properties

<http://hadoop.apache.org/docs/r1.1.2/core-default.html>

<http://hadoop.apache.org/docs/r1.1.2/mapred-default.htm>

|

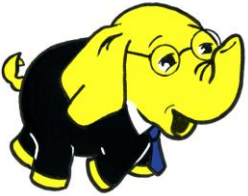
<http://hadoop.apache.org/docs/r1.1.2/hdfs-default.html>



Hadoop Wants to Know?

Which configuration file contains information about task trackers?

- a) /conf/masters
- b) /conf/slaves
- c) /conf/mapred-site.xml
- d) /conf/hdfs-site.xml



Slaves and Masters

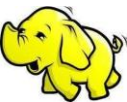
Two files are used by the startup and shutdown commands:

- **Slaves**

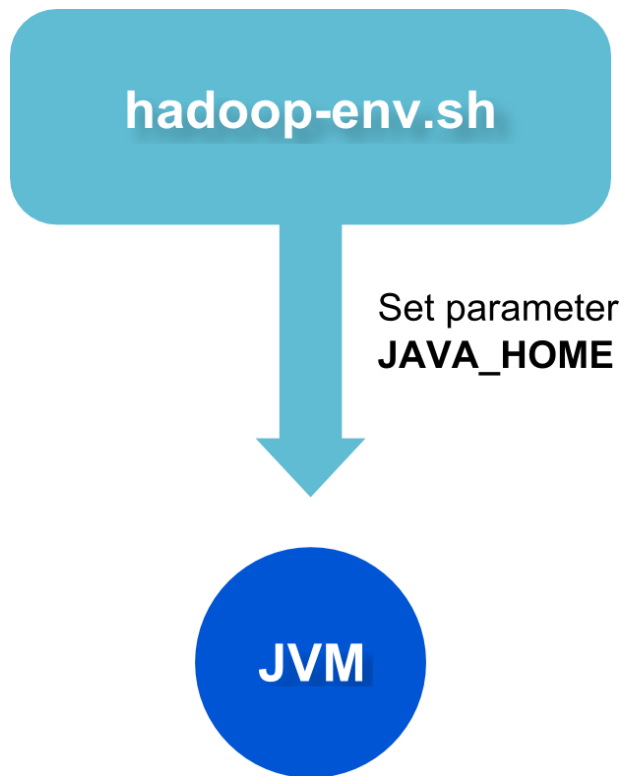
- Contains a list of hosts, one per line, that are to host DataNode and TaskTracker servers.

- **Masters**

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



Per Process Run-time Environment



- 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.

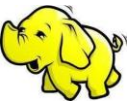
Examples of environment variables that you can specify:

Export:

```
HADOOP_DATANODE_HEAPSIZE="128"
```

Export :

```
HADOOP_TASKTRACKER_HEAPSIZE="512"
```

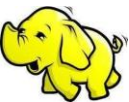


hadoop.env.sh - Sample

```
# Set Hadoop-specific environment variables here.
# The only required environment variable is JAVA_HOME. All others are
# optional. When running a distributed configuration it is best to
# set JAVA_HOME in this file, so that it is correctly defined on
# remote nodes.

# The java implementation to use. Required.
export JAVA_HOME=/usr/lib/jvm/java-6-sun-1.6.0.45

# Extra Java runtime options. Empty by default.
export HADOOP_OPTS="-Djava.net.preferIPv4Stack=true ${HADOOP_OPTS}"
...
...
# A string representing this instance of hadoop. $USER by default.
export HADOOP_IDENT_STRING=$USER
```



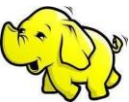
Critical Properties

fs.default.name

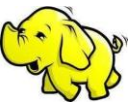
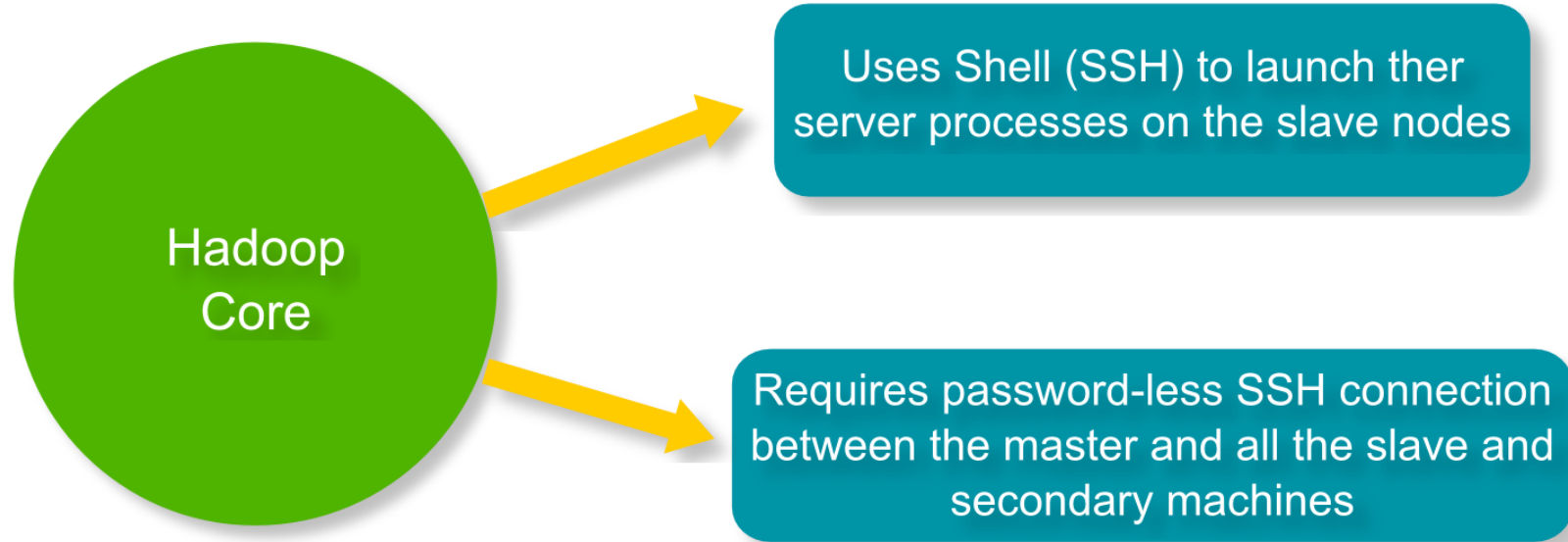
hadoop.tmp.dir

mapred.job.tracker

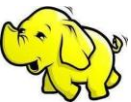
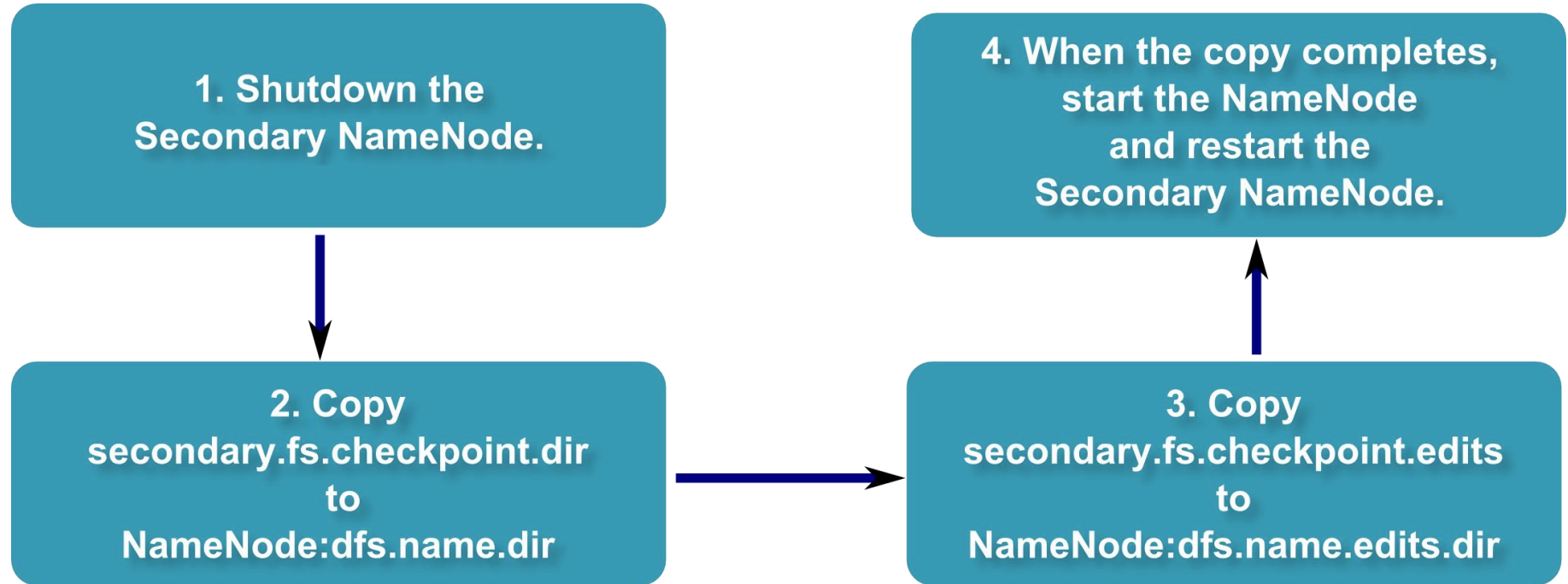
- **fs.default.name:**
 - It points to the default URI for all file system requests in Hadoop.
- **Hadoop.tmp.dir**
 - hadoop.tmp.dir is used as the base for temporary directories locally, and also in HDFS
- **Mapred.job.tracker**
 - The host and port of the MapReduce job tracker where it runs. If "local", then jobs are run in-process as a single map and reduce task



Network Requirements

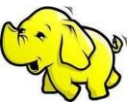


Name Node Recovery



Sample Examples List

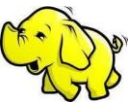
```
hadoopjob@ubuntu:/usr/local/hadoop$ ls
bin          hadoop-ant-1.2.1.jar      ivy          README.txt
build.xml    hadoop-client-1.2.1.jar   ivy.xml      sbin
c++          hadoop-core-1.2.1.jar     lib          share
CHANGES.txt hadoop-examples-1.2.1.jar libexec      src
conf         hadoop-minicluster-1.2.1.jar LICENSE.txt  webapps
contrib      hadoop-test-1.2.1.jar     logs
docs         hadoop-tools-1.2.1.jar    NOTICE.txt
```



Running the Example

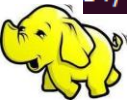
```
hadoopjob@ubuntu:/usr/local/hadoop$ hadoop jar hadoop-examples-1.2.1.jar teragen 100000 /user/teragen/teragen-inputTest
Warning: $HADOOP_HOME is deprecated.

Generating 100000 using 2 maps with step of 50000
14/02/14 14:31:35 INFO mapred.JobClient: Running job: job_201402141140_0001
14/02/14 14:31:36 INFO mapred.JobClient:  map 0% reduce 0%
```



Running the Example

```
14/02/14 14:32:04 INFO mapred.JobClient: map 100% reduce 0%
14/02/14 14:32:07 INFO mapred.JobClient: Job complete: job_201402141140_0001
14/02/14 14:32:07 INFO mapred.JobClient: Counters: 19
14/02/14 14:32:07 INFO mapred.JobClient:   Job Counters
14/02/14 14:32:07 INFO mapred.JobClient:     SLOTS_MILLIS_MAPS=43652
14/02/14 14:32:07 INFO mapred.JobClient:     Total time spent by all reduces waiting after reserving slots (ms)=0
14/02/14 14:32:07 INFO mapred.JobClient:     Total time spent by all maps waiting after reserving slots (ms)=0
14/02/14 14:32:07 INFO mapred.JobClient:     Launched map tasks=2
14/02/14 14:32:07 INFO mapred.JobClient:     SLOTS_MILLIS_REDUCE=0
14/02/14 14:32:07 INFO mapred.JobClient:   File Input Format Counters
14/02/14 14:32:07 INFO mapred.JobClient:     Bytes Read=0
14/02/14 14:32:07 INFO mapred.JobClient:   File Output Format Counters
14/02/14 14:32:07 INFO mapred.JobClient:     Bytes Written=10000000
14/02/14 14:32:07 INFO mapred.JobClient:   FileSystemCounters
14/02/14 14:32:07 INFO mapred.JobClient:     HDFS_BYTES_READ=164
14/02/14 14:32:07 INFO mapred.JobClient:     FILE_BYTES_WRITTEN=107108
14/02/14 14:32:07 INFO mapred.JobClient:     HDFS_BYTES_WRITTEN=10000000
14/02/14 14:32:07 INFO mapred.JobClient:   Map-Reduce Framework
14/02/14 14:32:07 INFO mapred.JobClient:     Map input records=100000
14/02/14 14:32:07 INFO mapred.JobClient:     Physical memory (bytes) snapshot=149934080
14/02/14 14:32:07 INFO mapred.JobClient:     Spilled Records=0
14/02/14 14:32:07 INFO mapred.JobClient:     CPU time spent (ms)=2370
14/02/14 14:32:07 INFO mapred.JobClient:     Total committed heap usage (bytes)=63307776
14/02/14 14:32:07 INFO mapred.JobClient:     Virtual memory (bytes) snapshot=1946124288
14/02/14 14:32:07 INFO mapred.JobClient:     Map input bytes=100000
14/02/14 14:32:07 INFO mapred.JobClient:     Map output records=100000
14/02/14 14:32:07 INFO mapred.JobClient:     SPLIT_RAW_BYTES=164
```



Running the Example

```
hadoopjob@ubuntu:/usr/local/hadoop$ hadoop fs -ls /user/teragen/teragen-inputTest
Warning: $HADOOP_HOME is deprecated.
```

```
Found 4 items
```

```
-rw-r--r--    1 stratapps supergroup          0 2014-02-14 14:32 /user/teragen/teragen-inputTest/_SUCCESS
drwxr-xr-x    - stratapps supergroup          0 2014-02-14 14:31 /user/teragen/teragen-inputTest/_logs
-rw-r--r--    1 stratapps supergroup 5000000 2014-02-14 14:31 /user/teragen/teragen-inputTest/part-00000
-rw-r--r--    1 stratapps supergroup 5000000 2014-02-14 14:31 /user/teragen/teragen-inputTest/part-00001
hadoopjob@ubuntu:/usr/local/hadoop$ hadoop fs -cat /user/teragen/teragen-inputTest/part-00000
```



Running the Example

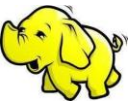
I#wql^@Woe
wyGUMkZHG^
=&z_SmO,bL
HqY:M y%A=
DNm&R>f\$be
(4L|p/"sBM
<AIQFo]_k;
=;WN,D\$xp
KD\$-4a.-]g
_A=w*qK\`4
08a)"V\<p`
xQAr\$"TP&Q
pK0~WIn5k:
55E(X \$+aq
J[OUDpFN]U
hSf3v9T1v_
q|0HybIHYM
~eA6?(=>>R
-ZsQxF9FCn
.J@0^WDE/d
Bg`+BA8=k
4{y*(O%\bL
w`K0q5xJ/I
K!E*n*3\!l
Q@6\$0\\?my
f'+H CzV)H
;o!HJ9Z\i<
+9IQ\$H^yv,

49972AAAAAAAAAABBBBBBBBBBCCCCCCCCCDDDDDDDDDEEEEEEEEEFFFFFFFFFFFGGGGGGGGGHHHHHHHH
49973IIIIIIIIIIJJJJJJJJJJKKKKKKKKKKLLLLLLLLLLMMMMMMMMMMMMNNNNNNNNNOOOOOOOOOPPPPPPP
49974QQQQQQQQQQRRRRRRRRRRSSSSSSSSSTTTTTTTTTUUUUUUUUUVVVVVVVVWWWWWWWWWXXXXXXX
49975YYYYYYYYYYZZZZZZZZZAAAAAAAAAABBBBBBBBBBCCCCCCCCCDDDDDDDDDEEEEEEEEEFFFFFFF
49976GGGGGGGGGGHHHHHHHHHHIIIIIIIIJJJJJJJJJJKKKKKKKKKKLLLLLLLLLLMMMMMMMMMMMMNNNNNNNN
49977000000000PPPPPPPPPPQQQQQQQQQQRRRRRRRRRRSSSSSSSSSTTTTTTTTTUUUUUUUUUVVVVVVVV
49978WWWWWWWWWXXXXXXXXXXYYYYYYYYYYZZZZZZZZZAAAAAAAAAABBBBBBBBBBCCCCCCCCCDDDDDDDD
49979EEEEEEEEEEFFFFFFFFFFGGGGGGGGGGHHHHHHHHHHIIIIIIIIJJJJJJJJJJKKKKKKKKKKLLLLLLLL
49980MMMMMMMMMMMMNNNNNNNNNNNOOOOOOOOOPPPPPPPPPQQQQQQQQQQRRRRRRRRRRSSSSSSSSSTTTTTTT
49981UUUUUUUUUVVVVVVVVWWWWWWWWWXXXXXXXXXXYYYYYYYYYYZZZZZZZZZAAAAAAAAAABBBBBBBB
49982CCCCCCCCCDDDDDDDDDEEEEEEEEEFFFFFFFFFFGGGGGGGGGGHHHHHHHHHHIIIIIIIIJJJJJJJJ
49983KKKKKKKKKKLLLLLLLLLLMMMMMMMMMMMMNNNNNNNNNNNOOOOOOOOOPPPPPPPPPQQQQQQQQQQRRRRRRRR
49984SSSSSSSSSTTTTTTTTTUUUUUUUUUVVVVVVVVWWWWWWWWWXXXXXXXXXXYYYYYYYYYYZZZZZZZZ
49985AAAAAAAAAABBBBBBBBCCCCCCCCCDDDDDDDDDEEEEEEEEEFFFFFFFFFFGGGGGGGGGGHHHHHHHH
49986IIIIIIIIIIJJJJJJJJJJKKKKKKKKKKLLLLLLLLLLMMMMMMMMMMMMNNNNNNNNNNNOOOOOOOOOPPPPPPP
49987QQQQQQQQQQRRRRRRRRRRSSSSSSSSSTTTTTTTTTUUUUUUUUUVVVVVVVVWWWWWWWWWXXXXXXXX
49988YYYYYYYYYYZZZZZZZZZAAAAAAAAAABBBBBBBBBBCCCCCCCCCDDDDDDDDDEEEEEEEEEFFFFFFF
49989GGGGGGGGGGHHHHHHHHHHIIIIIIIIJJJJJJJJJJKKKKKKKKKKLLLLLLLLLLMMMMMMMMMMMMNNNNNNNN
49990000000000PPPPPPPPPPQQQQQQQQQQRRRRRRRRRRSSSSSSSSSTTTTTTTTTUUUUUUUUUVVVVVVVV
49991WWWWWWWWWXXXXXXXXXXYYYYYYYYYYZZZZZZZZZAAAAAAAAAABBBBBBBBCCCCCCCCCDDDDDDDD
49992EEEEEEEEEEFFFFFFFFFFGGGGGGGGGGHHHHHHHHHHIIIIIIIIJJJJJJJJJJKKKKKKKKKKLLLLLLLL
49993MMMMMMMMMMMMNNNNNNNNNNNOOOOOOOOOPPPPPPPPPQQQQQQQQQQRRRRRRRRRRSSSSSSSSSTTTTTTT
49994UUUUUUUUUVVVVVVVVWWWWWWWWWXXXXXXXXXXYYYYYYYYYYZZZZZZZZZAAAAAAAAAABBBBBBBB
49995CCCCCCCCCDDDDDDDDDEEEEEEEEEFFFFFFFFFFGGGGGGGGGGHHHHHHHHHHIIIIIIIIJJJJJJJJ
49996KKKKKKKKKKLLLLLLLLLLMMMMMMMMMMMMNNNNNNNNNNNOOOOOOOOOPPPPPPPPPQQQQQQQQQQRRRRRRRR
49997SSSSSSSSSTTTTTTTTTUUUUUUUUUVVVVVVVVWWWWWWWWWXXXXXXXXXXYYYYYYYYYYZZZZZZZZ
49998AAAAAAAAAABBBBBBBBCCCCCCCCCDDDDDDDDDEEEEEEEEEFFFFFFFFFFGGGGGGGGGGHHHHHHHH
49999IIIIIIIIIIJJJJJJJJJJKKKKKKKKKKLLLLLLLLLLMMMMMMMMMMMMNNNNNNNNNNNOOOOOOOOOPPPPPPP



MapReduce

```
13/08/03 00:58:40 WARN mapred.JobClient: Use GenericOptionsParser for parsing the 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 waiting after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all maps waiting 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
```



MapReduce

```
13/08/03 00:58:40 WARN mapred.JobClient: Use GenericOptionsParser for parsing the 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 waiting after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all maps waiting 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
```

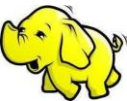
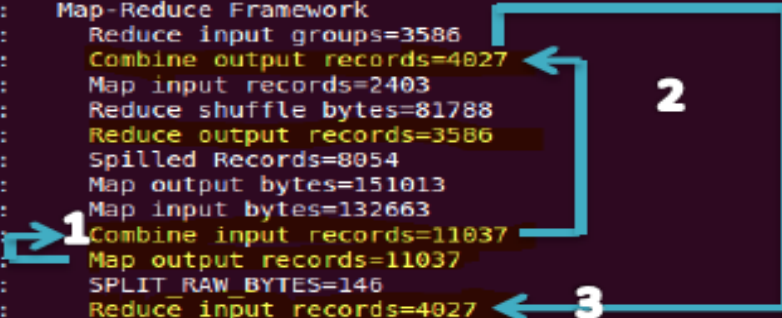
2

1



MapReduce

```
13/08/03 00:58:40 WARN mapred.JobClient: Use GenericOptionsParser for parsing the 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 waiting after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all maps waiting 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
```



MapReduce

```
13/08/03 00:58:40 WARN mapred.JobClient: Use GenericOptionsParser for parsing the 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 waiting after reserving slots (ms)=0
13/08/03 00:58:59 INFO mapred.JobClient: Total time spent by all maps waiting 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
```

4

2

1

3



Further reading and reference...

Companies powered by Hadoop

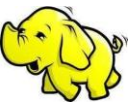
<https://wiki.apache.org/hadoop/PoweredBy>

MapReduce Paper by Google

<http://research.google.com/archive/mapreduce.html>

MapReduce Tutorial

https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html



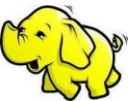
Further reading and reference...

Single Node Cluster Installation on Linux(Ubuntu)

<http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-single-node-cluster/>

Multi-Node Cluster Installation on Linux(Ubuntu)

<http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-multi-node-cluster/>



Thank You

