

# Hortonworks Data Platform

Reference

(Aug 29, 2014)

## Hortonworks Data Platform : Reference

Copyright © 2012-2014 Hortonworks, Inc. Some rights reserved.

The Hortonworks Data Platform, powered by Apache Hadoop, is a massively scalable and 100% open source platform for storing, processing and analyzing large volumes of data. It is designed to deal with data from many sources and formats in a very quick, easy and cost-effective manner. The Hortonworks Data Platform consists of the essential set of Apache Hadoop projects including MapReduce, Hadoop Distributed File System (HDFS), HCatalog, Pig, Hive, HBase, Zookeeper and Ambari. Hortonworks is the major contributor of code and patches to many of these projects. These projects have been integrated and tested as part of the Hortonworks Data Platform release process and installation and configuration tools have also been included.

Unlike other providers of platforms built using Apache Hadoop, Hortonworks contributes 100% of our code back to the Apache Software Foundation. The Hortonworks Data Platform is Apache-licensed and completely open source. We sell only expert technical support, [training](#) and partner-enablement services. All of our technology is, and will remain free and open source.

Please visit the [Hortonworks Data Platform](#) page for more information on Hortonworks technology. For more information on Hortonworks services, please visit either the [Support](#) or [Training](#) page. Feel free to [Contact Us](#) directly to discuss your specific needs.



Except where otherwise noted, this document is licensed under  
**Creative Commons Attribution ShareAlike 3.0 License.**  
<http://creativecommons.org/licenses/by-sa/3.0/legalcode>

## Table of Contents

1. Hadoop Service Accounts .....	1
2. Configuring Ports .....	2
3. Controlling HDP Services Manually .....	11
3.1. Starting HDP Services .....	11
3.2. Stopping HDP services .....	13
4. Deploying HDP In Production Data Centers with Firewalls .....	16
4.1. Deployment Strategies for Data Centers with Firewalls .....	16
4.1.1. Terminology .....	16
4.1.2. Mirroring or Proxying .....	17
4.1.3. Considerations for choosing a Mirror or Proxy solution .....	18
4.2. Recommendations for Deploying HDP .....	18
4.2.1. RPMs in the HDP repository .....	18
4.3. Detailed Instructions for Creating Mirrors and Proxies .....	19
4.3.1. Option I - Mirror server has no access to the Internet .....	19
4.3.2. Option II - Mirror server has temporary or continuous access to the Internet .....	23
4.3.3. Set up a trusted proxy server .....	28
5. Wire Encryption in Hadoop .....	31
6. Supported Database Matrix for Hortonworks Data Platform .....	32

## List of Tables

2.1. Accumulo Ports .....	2
2.2. Flume Ports .....	3
2.3. HDFS Ports .....	4
2.4. MapReduce Ports .....	4
2.5. YARN Ports .....	5
2.6. Hive Ports .....	6
2.7. HBase Ports .....	6
2.8. Oozie Ports .....	7
2.9. Sqoop Ports .....	8
2.10. Storm Ports .....	8
2.11. ZooKeeper Ports .....	9
2.12. MySQL Ports .....	9
2.13. Kerberos Ports .....	10
4.1. Terminology .....	16
4.2. Comparison - HDP Deployment Strategies .....	18
4.3. Deploying HDP - Option I .....	20
4.4. Options for <i>\$os</i> parameter in repo URL .....	22
4.5. Options for <i>\$os</i> parameter in repo URL .....	23
4.6. Deploying HDP - Option II .....	25
4.7. Options for <i>\$os</i> parameter in repo URL .....	27
6.1. Supported Databases for HDP Stack .....	32

# 1. Hadoop Service Accounts

You can configure service accounts using:

- If you are performing a Manual Install of HDP, refer to the Manual instructions on [Service User and Group accounts](#)
- If you are performing a Ambari Install of HDP, refer to the Ambari instructions in the [Ambari Install Guide](#) .

## 2. Configuring Ports

The tables below specify which ports must be opened for which ecosystem components to communicate with each other. Make sure the appropriate ports are opened before you install HDP.

### Accumulo Ports

The following table lists the default ports used by the various Accumulo services. (**Note:** Neither of these services are used in a standard HDP installation.)

**Table 2.1. Accumulo Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
Master	Master nodes (Active master and any standby)	9999		The Master thrift server	Yes (client API needs)	<code>master.port.client</code> in <code>accumulo-site.xml</code>
TabletServer	Slave nodes	9997		The TabletServer thrift server	Yes (client API needs)	<code>tserver.port.client</code> in <code>accumulo-site.xml</code>
Garbage Collector	GC nodes (Active GC and any standby)	50091		The GarbageCollector thrift server	No	<code>gc.port.client</code> in <code>accumulo-site.xml</code>
Monitor	Monitor nodes (Active Monitor and any standby)	50095	HTTP(S)	Metrics/ Monitoring of an Accumulo instance	Yes	<code>monitor.port.client</code> in <code>accumulo-site.xml</code>
Monitor log aggregation	Monitor nodes (Active Monitor and any standby)	4560		Log4j socket which accepts logs forwarded from other Accumulo services	No	<code>monitor.port.log4j</code> in <code>accumulo-site.xml</code>
Tracer	Tracer nodes	12234		The Tracer thrift server	Yes (if enabled)	<code>trace.port.client</code> in <code>accumulo-site.xml</code>
Thrift Proxy (optional)	Proxy nodes	42424		The Thrift Proxy server	Yes (if enabled)	<code>port</code> in <code>proxy.properties</code>
TabletServer Replication Service	Slave nodes	10002		TabletServer Thrift service supporting multi-instance Accumulo replication	No	<code>replication.receipt.service.port</code> in <code>accumulo-site.xml</code>
Master Replication Service	Master nodes (Active master and any standby)	10001		Master Thrift service supporting multi-instance Accumulo replication	No	<code>master.replication.coordinator.port</code> in <code>accumulo-site.xml</code>

### Flume Ports

The following table lists the default ports used by the various Flume services. (**Note:** Neither of these services are used in a standard HDP installation.)

**Table 2.2. Flume Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
Flume	Flume Agent	41414	TCP	Flume performance metrics in JSON format	Yes (client API needs)	<code>master.port.client</code> in <code>accumulo-site.xml</code>
Flume	HDFS Sink	8020	TCP	Communication from Flume into the Hadoop cluster's NameNode	Yes (client API needs)	<code>tserver.port.client</code> in <code>accumulo-site.xml</code>
Flume	HDFS Sink	9000	TCP	Communication from Flume into the Hadoop cluster's NameNode	No	<code>gc.port.client</code> in <code>accumulo-site.xml</code>
Flume	HDFS Sink	50010	TCP	Communication from Flume into the Hadoop cluster's HDFS DataNode	No	
Flume	HDFS Sink	50020	TCP	Communication from Flume into the Hadoop cluster's HDFS DataNode	No	
Flume	HBase Sink	2181	TCP	Communication from Flume into the Hadoop cluster's Zookeeper	No	
Flume	HBase Sink	60020	TCP	Communication from Flume into the Hadoop cluster's HBase Regionserver	No	
Flume	All Other Sources and Sinks	Variable	Variable	Ports and protocols used by Flume sources and sinks	No	Refer to the flume configuration file(s) for ports actually in use. Ports in use are specified using the port keyword in the Flume configuration file. By default Flume configuration files are located in <code>/etc/flume/conf</code> on Linux and <code>c:\hdp\flume-1.4.0.x.y.z\conf</code> on Windows

### HDFS Ports

The following table lists the default ports used by the various HDFS services.

**Table 2.3. HDFS Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
NameNode WebUI	Master Nodes (NameNode and any back-up NameNodes)	50070	http	Web UI to look at current status of HDFS, explore file system	Yes (Typically admins, Dev/Support teams)	dfs.http.address
		50470	https	Secure http service		dfs.https.address
NameNode metadata service		8020/9000	IPC	File system metadata operations	Yes (All clients who directly need to interact with the HDFS)	Embedded in URI specified by fs.defaultFS
DataNode	All Slave Nodes	50075	http	DataNode WebUI to access the status, logs etc.	Yes (Typically admins, Dev/Support teams)	dfs.datanode.http.address
		50475	https	Secure http service		dfs.datanode.https.address
		50010		Data transfer		dfs.datanode.address
		50020	IPC	Metadata operations	No	dfs.datanode.ipc.address
Secondary NameNode	Secondary NameNode and any backup Secondary NameNode	50090	http	Checkpoint for NameNode metadata	No	dfs.secondary.http.address

**MapReduce Ports:** The following table lists the default ports used by the various MapReduce services.

**Table 2.4. MapReduce Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
MapReduce		10020	http	MapReduce JobHistory server address		mapreduce.jobhistory.address
MapReduce		19888	http	MapReduce JobHistory webapp address		mapreduce.jobhistory.webapp.address
MapReduce		13562	http	MapReduce Shuffle Port		mapreduce.shuffle.port

**YARN Ports:** The following table lists the default ports used by the various YARN services.

Table 2.5. YARN Ports

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
Resource Manager WebUI	Master Nodes (Resource Manager and any back-up Resource Manager node)	8088	http	Web UI for Resource Manager	Yes	<code>yarn.resourcemanager.webapp.address</code>
Resource Manager	Master Nodes (Resource Manager Node)	8050	IPC	For application submissions	Yes (All clients who need to submit the YARN applications including Hive, Hive server, Pig)	Embedded in URI specified by <code>yarn.resourcemanager.address</code>
Resource Manager	Master Nodes (Resource Manager Node)	8025	http	For application submissions	Yes (All clients who need to submit the YARN applications including Hive, Hive server, Pig)	<code>yarn.resourcemanager.resource-tracker.address</code>
Scheduler	Master Nodes (Resource Manager Node)	8030	http	Scheduler Address	Yes (Typically admins, Dev/ Support teams)	<code>yarn.resourcemanager.scheduler.address</code>
Resource Manager	Master Nodes (Resource Manager Node)	8141	http	Scheduler Address	Yes (Typically admins, Dev/ Support teams)	<code>yarn.resourcemanager.admin.address</code>
NodeManager	Master Nodes (NodeManager) and Slave Nodes	45454	http	NodeManager Address	Yes (Typically admins, Dev/ Support teams)	<code>yarn.nodemanager.address</code>
Timeline Server	Master Nodes	10200	http	Timeline Server Address	Yes (Typically admins, Dev/ Support teams)	<code>yarn.timeline-service.address</code>
Timeline Server	Master Nodes	8188	http	Timeline Server Webapp Address	Yes (Typically admins, Dev/ Support teams)	<code>yarn.timeline-service.webapp.address</code>
Timeline Server	Master Nodes	8190	https	Timeline Server	Yes (Typically	<code>yarn.timeline-service.webapp.https.address</code>

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
				Webapp https Address	admins, Dev/ Support teams)	
Job History Service	Master Nodes	19888	https	Job History Service	Yes (Typically admins, Dev/ Support teams)	yarn.log.server.url

## Hive Ports

The following table lists the default ports used by the various Hive services. (**Note:** Neither of these services are used in a standard HDP installation.)

**Table 2.6. Hive Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
Hive Server	Hive Server machine (Usually a utility machine)	10000		Service for programatically (Thrift/JDBC) connecting to Hive	Yes (Clients who need to connect to Hive either programatically or through UI SQL tools that use JDBC)	ENV Variable HIVE_PORT
Hive Web UI	Hive Server machine (Usually a utility machine)	9999	http	Web UI to explore Hive schemas	Yes	hive.hwi.listen.port
Hive Metastore		9933	http		Yes (Clients that run Hive, Pig and potentially M/R jobs that use HCatalog)	hive.metastore.uris

## HBase Ports

The following table lists the default ports used by the various HBase services.

**Table 2.7. HBase Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
HMaster	Master Nodes (HBase Master Node and any back-up HBase Master node)	60000			Yes	hbase.master.port
HMaster Info Web UI	Master Nodes (HBase master Node and	60010	http	The port for the HBaseMaster web UI. Set to	Yes	hbase.master.info.port

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
	back up HBase Master node if any)			-1 if you do not want the info server to run.		
Region Server	All Slave Nodes	60020			Yes (Typically admins, dev/ support teams)	<code>hbase.regionserver.port</code>
Region Server	All Slave Nodes	60030	http		Yes (Typically admins, dev/ support teams)	<code>hbase.regionserver.info.port</code>
HBase REST Server (optional)	All REST Servers	8080	http	The port used by HBase Rest Servers. REST servers are optional, and not installed by default	Yes	<code>hbase.rest.port</code>
HBase REST Server Web UI (optional)	All REST Servers	8085	http	The port used by HBase Rest Servers web UI. REST servers are optional, and not installed by default	Yes (Typically admins, dev/ support teams)	<code>hbase.rest.info.port</code>
HBase Thrift Server (optional)	All Thrift Servers	9090		The port used by HBase Thrift Servers. Thrift servers are optional, and not installed by default	Yes	
HBase Thrift Server Web UI (optional)	All Thrift Servers	9095		The port used by HBase Thrift Servers web UI. Thrift servers are optional, and not installed by default	Yes (Typically admins, dev/ support teams)	<code>hbase.thrift.info.port</code>

**Oozie Ports:** The following table lists the default ports used by Oozie.

**Table 2.8. Oozie Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
Oozie	Oozie Server	11000	TCP	The port Oozie server runs.	Yes	<code>OOZIE_HTTP_PORT</code> in <code>oozie_env.sh</code>
Oozie	Oozie Server	11001	TCP	The admin port Oozie server runs.	No	<code>OOZIE_ADMIN_PORT</code> in <code>oozie_env.sh</code>

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
Oozie	Oozie Server	11443	TCP	The port Oozie server runs when using HTTPS.	Yes	OOZIE_HTTPS_PORT in oozie_env.sh

**Sqoop Ports:** The following table lists the default ports used by Sqoop.

**Table 2.9. Sqoop Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
Sqoop	Metastore	16000	TCP	Connection between Sqoop and the metastore	No	sqoop.metastore.server.port
Sqoop	JDBC Listener	Varies, depends on target database. For example, if moving data from MySQL, TCP port 3306 must be open.	TCP	Outbound port from the Hadoop cluster to the database. Varies depending on Database	No	

### Storm Ports

The following table lists the default ports used by Storm.

**Table 2.10. Storm Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
Zookeeper Port		2181		Port used by localhost to talk to ZooKeeper.		storm.zookeeper.port
DRPC Port		3772				drpc.port
DRPC Invocations Port		3773				drpc.invocations.port
Nimbus Thrift Port		6627				nimbus.thrift.port
Supervisor Slots Ports		6700, 6701,		Defines the amount of		supervisor.slots.ports

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
		6702, 7603		workers that can be run on this machine. Each worker is assigned a port to use for communication.		
Logviewer Port		8000				logviewer.port
UI Port		8080				ui.port
Ambari Port		8744				ambari.port

## ZooKeeper Ports

**Table 2.11. ZooKeeper Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
ZooKeeper Server	All ZooKeeper Nodes	2888		Port used by ZooKeeper peers to talk to each other. See <a href="#">here</a> for more information.	No	hbase.zookeeper.peerport
ZooKeeper Server	All ZooKeeper Nodes	3888		Port used by ZooKeeper peers to talk to each other. See <a href="#">here</a> for more information.	No	hbase.zookeeper.leaderport
ZooKeeper Server	All ZooKeeper Nodes	2181		Property from ZooKeeper's config <code>zoo.cfg</code> . The port at which the clients will connect.	No	hbase.zookeeper.property.clientPort

**MySQL Ports:** The following table lists the default ports used by the various MySQL services.

**Table 2.12. MySQL Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
MySQL	MySQL database server	3306				

**Kerberos Ports:** The following table lists the default port used by the designated Kerberos KDC.

**Table 2.13. Kerberos Ports**

Service	Servers	Default Ports Used	Protocol	Description	Need End User Access?	Configuration Parameters
KDC	Kerberos KDC server	88		Port used by the designated KDC		

## 3. Controlling HDP Services Manually

In this document:

- [Starting HDP Services](#)
- [Stopping HDP Services](#)

### 3.1. Starting HDP Services

Start the Hadoop services in the following order:

- Knox
- ZooKeeper
- HDFS
- YARN
- HBase
- Hive Metastore
- HiveServer2
- WebHCat
- Oozie
- Storm

#### Instructions

1. Start Knox. When starting the gateway with the script below, the process runs in the background. The log output is written to `/var/log/knox` and a PID (process ID) is written to `/var/run/knox`. Execute this command on the Knox host machine.

```
cd $GATEWAY_HOME su -l Knox -c "bin/gateway.sh start"
```

where `$GATEWAY_HOME` is the directory where Knox is installed. For example, `/usr/lib/knox`.



#### Note

If Knox has been stopped without using `gateway.sh stop`, you must start the service using `gateway.sh clean`. The clean option removes all log files in `/var/log/knox`.

2. Start ZooKeeper. Execute this command on the ZooKeeper host machine(s):

```
su - zookeeper -c "export ZOOCFGDIR=/etc/zookeeper/conf ; export ZOOCFG=zoo.cfg ; source /etc/zookeeper/conf/zookeeper-env.sh ; /usr/lib/zookeeper/bin/zkServer.sh start"
```

### 3. Start HDFS

- a. If you are running NameNode HA (High Availability), start the JournalNodes by executing these commands on the JournalNode host machines:

```
su $HDFS_USER  
/usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/hadoop/conf start  
journalnode
```

where `$HDFS_USER` is the HDFS user. For example, `hdfs`.

- b. Execute this command on the NameNode host machine(s):

```
su -l hdfs -c "/usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/  
hadoop/conf start namenode"
```

- c. If you are running NameNode HA, start the Zookeeper Failover Controller (ZKFC) by executing the following command on all NameNode machines. The starting sequence of the ZKFCs determines which NameNode will become Active.

```
su -l hdfs -c "/usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/  
hadoop/conf start zkfc"
```

- d. If you are not running NameNode HA, execute the following command on the Secondary NameNode host machine. If you are running NameNode HA, the Standby NameNode takes on the role of the Secondary NameNode.

```
su -l hdfs -c "/usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/  
hadoop/conf start secondarynamenode"
```

- e. Execute these commands on all DataNodes:

```
su -l hdfs -c "/usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/  
hadoop/conf start datanode"
```

### 4. Start YARN

- a. Execute this command on the ResourceManager host machine(s):

```
su -l yarn -c "export HADOOP_LIBEXEC_DIR=/usr/lib/hadoop/libexec && /  
usr/lib/hadoop-yarn/sbin/yarn-daemon.sh --config /etc/hadoop/conf start  
resourcemanager"
```

- b. Execute this command on the History Server host machine:

```
su -l mapred -c "export HADOOP_LIBEXEC_DIR=/usr/lib/hadoop/libexec && /  
usr/lib/hadoop-mapreduce/sbin/mr-jobhistory-daemon.sh --config /etc/  
hadoop/conf start historyserver"
```

- c. Execute this command on all NodeManagers:

```
su -l yarn -c "export HADOOP_LIBEXEC_DIR=/usr/lib/hadoop/libexec && /  
usr/lib/hadoop-yarn/sbin/yarn-daemon.sh --config /etc/hadoop/conf start  
nodemanager"
```

## 5. Start HBase

- a. Execute this command on the HBase Master host machine:

```
su -l hbase -c "/usr/lib/hbase/bin/hbase-daemon.sh --config /etc/hbase/conf start master; sleep 25"
```

- b. Execute this command on all RegionServers:

```
su -l hbase -c "/usr/lib/hbase/bin/hbase-daemon.sh --config /etc/hbase/conf start regionserver"
```

6. Start the Hive Metastore. On the Hive Metastore host machine, execute the following command:

```
su $HIVE_USER  
nohup hive --service metastore>/var/log/hive/hive.out 2>/var/log/hive/hive.log &
```

where, `$HIVE_USER` is the Hive user. For example, `hive`.

7. Start HiveServer2. On the Hive Server2 host machine, execute the following command:

```
su $HIVE_USER  
nohup /usr/lib/hive/bin/hiveserver2 -hiveconf hive.metastore.uris=" " >>/tmp/hiveserver2HD.out 2>> /tmp/hiveserver2HD.log &
```

where, `$HIVE_USER` is the Hive user. For example, `hive`.

8. Start WebHCat. On the WebHCat host machine, execute the following command:

```
su -l hcat -c "/usr/lib/hcatalog/sbin/webhcat_server.sh start"
```

9. Start Oozie. Execute these commands on the Oozie host machine.

```
su $OOZIE_USER  
/usr/lib/oozie/bin/oozie-start.sh
```

where `$OOZIE_USER` is the Oozie user. For example, `oozie`.

- 10 Start Storm using a process controller, such as supervisor.

```
su $STORM_USER  
/usr/bin/supervisord start
```

where `$STORM_USER` is the Storm user. For example, `storm`.

## 3.2. Stopping HDP services

Before performing any upgrades or uninstalling software, stop all of the Hadoop services in the following order:

- Knox
- Oozie
- WebHCat

- HiveServer2
- Hive Metastore
- HBase
- YARN
- HDFS
- Zookeeper
- Storm

### Instructions

1. Stop Knox. Execute this command on the Knox host machine.

```
cd $GATEWAY_HOME su -l knox -c "bin/gateway.sh stop"
```

where `$GATEWAY_HOME` is the directory where Knox is installed. For example, `/usr/lib/knox`.

2. Stop Oozie. Execute these commands on the Oozie host machine.

```
su $OOZIE_USER  
/usr/lib/oozie/bin/oozie-stop.sh
```

where `$OOZIE_USER` is the Oozie user. For example, `oozie`.

3. Stop WebHCat. On the WebHCat host machine, execute the following command:

```
su -l hcat -c "/usr/lib/hcatalog/sbin/webhcat_server.sh stop"
```

4. Stop Hive. Execute this command on the Hive Metastore and Hive Server2 host machine.

```
ps aux | awk '{print $1,$2}' | grep hive | awk '{print $2}' | xargs kill >/dev/null 2>&1
```

5. Stop HBase

- a. Execute this command on all RegionServers:

```
su -l hbase -c "/usr/lib/hbase/bin/hbase-daemon.sh --config /etc/hbase/conf stop regionserver"
```

- b. Execute this command on the HBase Master host machine:

```
su -l hbase -c "/usr/lib/hbase/bin/hbase-daemon.sh --config /etc/hbase/conf stop master"
```

6. Stop YARN

- a. Execute this command on all NodeManagers:

```
su -l yarn -c "export HADOOP_LIBEXEC_DIR=/usr/lib/hadoop/libexec && /usr/lib/hadoop-yarn/sbin/yarn-daemon.sh --config /etc/hadoop/conf stop nodemanager"
```

- b. Execute this command on the History Server host machine:

```
su -l mapred -c "export HADOOP_LIBEXEC_DIR=/usr/lib/hadoop/libexec && /usr/lib/hadoop-mapreduce/sbin/mr-jobhistory-daemon.sh --config /etc/hadoop/conf stop historyserver"
```

- c. Execute this command on the ResourceManager host machine(s):

```
su -l yarn -c "export HADOOP_LIBEXEC_DIR=/usr/lib/hadoop/libexec && /usr/lib/hadoop-yarn/sbin/yarn-daemon.sh --config /etc/hadoop/conf stop resourcemanager"
```

## 7. Stop HDFS

- a. Execute this command on all DataNodes:

```
su -l hdfs -c "/usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/hadoop/conf stop datanode"
```

- b. If you are not running NameNode HA (High Availability), stop the Secondary NameNode by executing this command on the Secondary NameNode host machine:

```
su -l hdfs -c "/usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/hadoop/conf stop secondarynamenode"
```

- c. Execute this command on the NameNode host machine(s):

```
su -l hdfs -c "/usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/hadoop/conf stop namenode"
```

- d. If you are running NameNode HA, stop the Zookeeper Failover Controllers (ZKFC) by executing this command on the NameNode host machines:

```
su -l hdfs -c "/usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/hadoop/conf stop zkfc"
```

- e. If you are running NameNode HA, stop the JournalNodes by executing these commands on the JournalNode host machines:

```
su $HDFS_USER /usr/lib/hadoop/sbin/hadoop-daemon.sh --config /etc/hadoop/conf stop journalnode
```

where `$HDFS_USER` is the HDFS user. For example, `hdfs`.

8. Stop ZooKeeper. Execute this command on the ZooKeeper host machine(s):

```
su - zookeeper -c "export ZOOCFGDIR=/etc/zookeeper/conf ; export ZOOCFG=zoo.cfg ; source /etc/zookeeper/conf/zookeeper-env.sh ; /usr/lib/zookeeper/bin/zkServer.sh stop"
```

9. Stop Storm using a process controller, such as supervisord.

```
su $STORM_USER /usr/bin/supervisord stop
```

where `$STORM_USER` is the Storm user. For example, `storm`.

## 4. Deploying HDP In Production Data Centers with Firewalls

In this document:

- [Deployment strategies for data centers with firewall](#)
- [Recommendations for deploying HDP](#)
- [Detailed instructions for creating mirrors and proxies](#)

### 4.1. Deployment Strategies for Data Centers with Firewalls

A typical Hortonworks Data Platform (HDP) install requires access to the Internet in order to fetch software packages from a remote repository. Because corporate networks typically have various levels of firewalls, these firewalls may limit or restrict Internet access, making it impossible for your cluster nodes to access the HDP repository during the install process.

The solution for this is to either:

- Create a local mirror repository inside your firewall hosted on a local mirror server inside your firewall; or
- Provide a trusted proxy server inside your firewall that can access the hosted repositories.



#### Note

Many of the descriptions in this section assume you are using RHEL/Centos/Oracle Linux. If you are using SLES, please adjust the commands and directories accordingly.

This document will cover these two options in detail, discuss the trade-offs, provide configuration guidelines, and will also provide recommendations for your deployment strategy.

In general, before installing Hortonworks Data Platform in a production data center, it is best to ensure that both the Data Center Security team and the Data Center Networking team are informed and engaged to assist with these aspects of the deployment.

#### 4.1.1. Terminology

**Table 4.1. Terminology**

Item	Description
Yum Package Manager (yum)	A package management tool that fetches and installs software packages and performs automatic dependency resolution. See <a href="http://yum.baseurl.org/">http://yum.baseurl.org/</a> for more information.

Item	Description
Local Mirror Repository	The yum repository hosted on your Local Mirror Server that will serve the HDP software.
Local Mirror Server	The server in your network that will host the Local Mirror Repository. This server must be accessible from all hosts in your cluster where you will install HDP.
HDP Repositories	A set of repositories hosted by Hortonworks that contains the HDP software packages. HDP software packages include the HDP Repository and the HDP-UTILS Repository.
HDP Repository Tarball	A tarball image that contains the complete contents of the HDP Repositories.

## 4.1.2. Mirroring or Proxying

HDP uses yum or zypper to install software, and this software is obtained from the HDP Repositories. If your firewall prevents Internet access, you must mirror or proxy the HDP Repositories in your Data Center.

Mirroring a repository involves copying the entire repository and all its contents onto a local server and enabling an HTTPD service on that server to serve the repository locally. Once the local mirror server setup is complete, the `*.repo` configuration files on every cluster node must be updated, so that the given package names are associated with the local mirror server instead of the remote repository server.

There are three options for creating a local mirror server. Each of these options is explained in detail in a later section.

- **Mirror server has no access to Internet at all:** Use a web browser on your workstation to download the HDP Repository Tarball, move the tarball to the selected mirror server using scp or an USB drive, and extract it to create the repository on the local mirror server.
- **Mirror server has temporary access to Internet:** Temporarily configure a server to have Internet access, download a copy of the HDP Repository to this server using the **reposync** command, then reconfigure the server so that it is back behind the firewall.



### Note

Option I is probably the least effort, and in some respects, is the most secure deployment option.

Option III is best if you want to be able to update your Hadoop installation periodically from the Hortonworks Repositories.

- **Trusted proxy server:** Proxying a repository involves setting up a standard HTTP proxy on a local server to forward repository access requests to the remote repository server and route responses back to the original requestor. Effectively, the proxy server makes the repository server accessible to all clients, by acting as an intermediary.

Once the proxy is configured, change the `/etc/yum.conf` file on every cluster node, so that when the client attempts to access the repository during installation, the request goes through the local proxy server instead of going directly to the remote repository server.

### 4.1.3. Considerations for choosing a Mirror or Proxy solution

The following table lists some benefits provided by these alternative deployment strategies:

**Table 4.2. Comparison - HDP Deployment Strategies**

Advantages of repository mirroring (Options I and II)	Advantages of creating a proxy
<ul style="list-style-type: none"> <li>Minimizes network access (after the initial investment of copying the repository to local storage). The install process is therefore faster, reliable, and more cost effective (reduced WAN bandwidth minimizes the data center costs).</li> <li>Allows security-conscious data centers to qualify a fixed set of repository files. It also ensures that the remote server will not change these repository files.</li> <li>Large data centers may already have existing repository mirror servers for the purpose of OS upgrades and software maintenance. You can easily add the HDP Repositories to these existing servers.</li> </ul>	<ul style="list-style-type: none"> <li>Avoids the need for long term management of the repository files (including periodic updates for upgrades, new versions, and bug fixes).</li> <li>Almost all data centers already have a setup of well-known proxies. In such cases, you can simply add the local proxy server to the existing proxies' configurations. This approach is easier compared to creating local mirror servers in data centers with no mirror server setup.</li> <li>The network access is same as that required when using a mirror repository, but the source repository handles file management.</li> </ul>

However, each of the above approaches are also known to have the following disadvantages:

- Mirrors have to be managed for updates, upgrades, new versions, and bug fixes.
- Proxy servers rely on the repository provider to not change the underlying files without notice.
- Caching proxies are necessary, because non-caching proxies do not decrease WAN traffic and do not speed up the install process.

## 4.2. Recommendations for Deploying HDP

This section provides information on the various components of the Apache Hadoop ecosystem.

In many data centers, using a mirror for the HDP Repositories can be the best deployment strategy. The HDP Repositories are small and easily mirrored, allowing you secure control over the contents of the Hadoop packages accepted for use in your data center.



### Note

The installer pulls many packages from the base OS repositories (repos). If you do not have a complete base OS available to all your machines at the time of installation, you may run into issues. If you encounter problems with base OS repos being unavailable, please contact your system administrator to arrange for these additional repos to be proxied or mirrored.

### 4.2.1. RPMs in the HDP repository

In the HDP repository, you will find two different source RPM for each component.

For example, for Hadoop, you should find the following two RPMs:

- `hadoop-x.x.x.x.el6.src.rpm`
- `hadoop-source-x.x.x.x.el6.i386.rpm`

The `src` and `source` are two different packages that serve the following purpose:

- The `src` package is used to re-create the binary in a given environment. You can use the `src` package of a particular component if you want to rebuild RPM for that component.
- The `source` package on the other hand, is used for reference or debugging purpose. The `source` package is particularly useful when you want to examine the source code of a particular component in a deployed cluster.

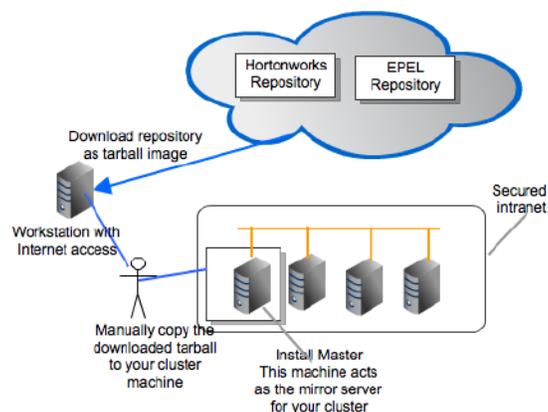
## 4.3. Detailed Instructions for Creating Mirrors and Proxies

In this section:

- [Option I: Mirror server has no access to the Internet](#)
- [Option II - Mirror server has temporary or continuous access to the Internet](#)
- [Trusted proxy server](#)

### 4.3.1. Option I - Mirror server has no access to the Internet

The local mirror setup for Option I is shown in the following illustration:



Complete the following instructions to set up a mirror server that has no access to the Internet:

1. [Check Your Prerequisites](#)
2. [Install the Repos](#)

### 4.3.1.1. Check Your Prerequisites

Select a mirror server host with the following characteristics:

- This server runs on either CentOS (v5.x, v6.x), RHEL (v5.x, v6.x), Oracle Linux(v5.x, v6.x), SLES 11, or Ubuntu 12, and has several GB of storage available.
- This server and the cluster nodes are all running the same OS.



#### Note

To support repository mirroring for heterogeneous clusters requires a more complex procedure than the one documented here.

- The firewall lets all cluster nodes (the servers on which you want to install HDP) to access this server.

### 4.3.1.2. Install the Repos

1. Use a workstation with access to the Internet and download the tarball image of the appropriate Hortonworks yum repository.

**Table 4.3. Deploying HDP - Option I**

Cluster OS	HDP Repository Tarballs
RHEL/ CentOS/ Oracle Linux 5.x	<pre>wget http://public-repo-1.hortonworks.com/HDP/centos5/HDP-2.1.5.0- centos5-rpm.tar.gz</pre> <pre>wget http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.19/repos/ centos5/HDP-UTILS-1.1.0.19-centos5.tar.gz</pre>
RHEL/CentOS/ Oracle Linux 6.x	<pre>wget http://public-repo-1.hortonworks.com/HDP/centos6/HDP-2.1.5.0- centos6-rpm.tar.gz</pre> <pre>wget http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.19/repos/ centos6/HDP-UTILS-1.1.0.19-centos6.tar.gz</pre>
SLES 11 SP1	<pre>wget http://public-repo-1.hortonworks.com/HDP/sles11sp1/HDP-2.1.5.0- sles11sp1-rpm.tar.gz</pre> <pre>wget http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.19/repos/ sles11sp1/HDP-UTILS-1.1.0.19-sles11sp1.tar.gz</pre>
SLES 11 SP3	<pre>wget http://public-repo-1.hortonworks.com/HDP/susel1sp3/HDP-2.1.5.0- susel1sp3-rpm.tar.gz</pre> <pre>wget http://public-repo-1.hortonworks.com/HDP-UTILS-1.1.0.19/repos/ susel1sp3/HDP-UTILS-1.1.0.19-susel1sp3.tar.gz</pre>
Ubuntu 12.04	<pre>wget http://public-repo-1.hortonworks.com/HDP/ubuntu12/HDP-2.1.5.0- ubuntu12-tars-tarball.tar.gz</pre> <pre>wget http://public-repo-1.hortonworks.com/HDP/tools/2.1.5.0/ hdp_manual_install_rpm_helper_files-2.1.5.695.tar.gz</pre>
Debian 6	<pre>wget http://public-repo-1.hortonworks.com/HDP/debian6/HDP-2.1.5.0- debian6-tars-tarball.tar.gz</pre> <pre>wget http://public-repo-1.hortonworks.com/HDP/tools/2.1.5.0/ hdp_manual_install_rpm_helper_files-2.1.5.695.tar.gz</pre>

2. Create an HTTP server.

- a. On the mirror server, install an HTTP server (such as Apache httpd) using the instructions provided [here](#).
- b. Activate this web server.
- c. Ensure that the firewall settings (if any) allow inbound HTTP access from your cluster nodes to your mirror server.



### Note

If you are using EC2, make sure that SELinux is disabled.

3. On your mirror server, create a directory for your web server.
  - For example, from a shell window, type:
    - **For RHEL/CentOS/Oracle:**

```
mkdir -p /var/www/html/hdp/
```
    - **For SLES:**

```
mkdir -p /srv/www/htdocs/rpms
```
    - **For Ubuntu:**

```
mkdir -p /var/www/html/hdp/
```
  - If you are using a symlink, enable the **followsymlinks** on your web server.
4. Copy the HDP Repository Tarball to the directory created in step 3, and untar it.
5. Verify the configuration.
  - The configuration is successful, if you can access the above directory through your web browser.

To test this out, browse to the following location: `http://$yourwebserver/hdp/$os/HDP-2.1.3.0/`.

You should see directory listing for all the HDP components along with the RPMs at: `$os/HDP-2.1.5.0`.



### Note

If you are installing a 2.x.0 release, use: `http://$yourwebserver/hdp/$os/2.x/GA`

If you are installing a 2.x.x release, use: `http://$yourwebserver/hdp/$os/2.x/updates`

where

- `$os` can be `centos5`, `centos6`, `suse11`, or `ubnuntu12`. Use the following options table for `$os` parameter:

**Table 4.4. Options for `$os` parameter in repo URL**

Operating System	Value
CentOS 5	centos5
RHEL 5	
Oracle Linux 5	
CentOS 6	centos6
RHEL 6	
Oracle Linux 6	
SLES 11	suse11
Ubuntu 12	ubuntu12

6. Configure the **yum** clients on all the nodes in your cluster.

- a. Fetch the yum configuration file from your mirror server.

```
http://<$yourwebserver>/hdp/$os/2.x/updates/2.1.5.0/hdp.repo
```

- b. Store the `hdp.repo` file to a temporary location.

- c. Edit `hdp.repo` file changing the value of the **baseurl** property to point to your local repositories based on your cluster OS.

```
[HDP-2.x]
name=Hortonworks Data Platform Version - HDP-2.x
baseurl=http://$yourwebserver/HDP/$os/2.x/GA
gpgcheck=1
gpgkey=http://public-repo-1.hortonworks.com/HDP/$os/RPM-GPG-KEY/RPM-GPG-KEY-Jenkins
enabled=1
priority=1

[HDP-UTILS-1.1.0.19]
name=Hortonworks Data Platform Utils Version - HDP-UTILS-1.1.0.19
baseurl=http://$yourwebserver/HDP-UTILS-1.1.0.19/repos/$os
gpgcheck=1
gpgkey=http://public-repo-1.hortonworks.com/HDP/$os/RPM-GPG-KEY/RPM-GPG-KEY-Jenkins
enabled=1
priority=1

[HDP-2.1.5.0]
name=Hortonworks Data Platform HDP-2.1.5.0
baseurl=http://$yourwebserver/HDP/$os/2.x/updates/2.1.5.0
gpgcheck=1
gpgkey=http://public-repo-1.hortonworks.com/HDP/$os/RPM-GPG-KEY/RPM-GPG-KEY-Jenkins
enabled=1
priority=1
```

where

- `$yourwebserver` is FQDN of your local mirror server.

- `$os` can be `centos5`, `centos6`, `suse11`, or `ubuntu12`. Use the following options table for `$os` parameter:

**Table 4.5. Options for `$os` parameter in repo URL**

Operating System	Value
CentOS 5	centos5
RHEL 5	
Oracle Linux 5	
CentOS 6	centos6
RHEL 6	
Oracle Linux 6	
SLES 11	suse11
Ubuntu 12	ubuntu12

- Use `scp` or `pdsh` to copy the client yum configuration file to `/etc/yum.repos.d/` directory on every node in the cluster.
- d. [Conditional]: If you have multiple repositories configured in your environment, deploy the following plugin on all the nodes in your cluster.
- i. Install the plugin.

- **For RHEL and CentOS v5.x**

```
yum install yum-priorities
```

- **For RHEL and CentOS v6.x**

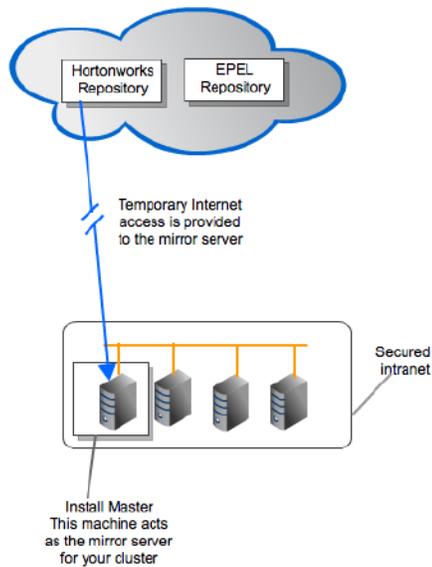
```
yum install yum-plugin-priorities
```

- ii. Edit the `/etc/yum/pluginconf.d/priorities.conf` file to add the following:

```
[main]
enabled=1
gpgcheck=0
```

### 4.3.2. Option II - Mirror server has temporary or continuous access to the Internet

The local mirror setup for Option II is shown in the following illustration:



Complete the following instructions to set up a mirror server that has temporary access to the Internet:

1. [Check Your Prerequisites](#)
2. [Install the Repos](#)

### 4.3.2.1. Check Your Prerequisites

Select a local mirror server host with the following characteristics:

- This server runs on either CentOS/RHEL/Oracle Linux 5.x or 6.x, SLES 11, or Ubuntu 12, and has several GB of storage available.
- The local mirror server and the cluster nodes must have the same OS. If they are not running CentOS or RHEL, the mirror server must not be a member of the Hadoop cluster.



#### Note

To support repository mirroring for heterogeneous clusters requires a more complex procedure than the one documented here.

- The firewall allows all cluster nodes (the servers on which you want to install HDP) to access this server.
- Ensure that the mirror server has **yum** installed.
- Add the **yum-utils** and **createrepo** packages on the mirror server.

```
yum install yum-utils createrepo
```

### 4.3.2.2. Install the Repos

- Temporarily reconfigure your firewall to allow Internet access from your mirror server host.

- Execute the following command to download the appropriate Hortonworks yum client configuration file and save it in `/etc/yum.repos.d/` directory on the mirror server host.

**Table 4.6. Deploying HDP - Option II**

Cluster OS	HDP Repository Tarballs
RHEL/CentOS/ Oracle Linux 5.x	<code>wget http://public-repo-1.hortonworks.com/HDP/centos5/2.x/updates/2.1.5.0/hdp.repo -O /etc/yum.repos.d/hdp.repo</code>
RHEL/CentOS/ Oracle Linux 6.x	<code>wget http://public-repo-1.hortonworks.com/HDP/centos6/2.x/updates/2.1.5.0/hdp.repo -O /etc/yum.repos.d/hdp.repo</code>
SLES 11 SP1	<code>wget http://public-repo-1.hortonworks.com/HDP/sles11sp1/2.x/updates/2.1.5.0/hdp.repo -O /etc/zypp/repos.d/hdp.repo</code>
SLES 11 SP3	<code>wget http://public-repo-1.hortonworks.com/HDP/susellsp3/2.x/updates/2.1.5.0/hdp.repo -O /etc/zypp/repos.d/hdp.repo</code>
Ubuntu 12.04	<code>wget http://public-repo-1.hortonworks.com/HDP/ubuntu12/2.1.5.0/hdp.list -O /etc/apt/sources.list.d/hdp.list</code>
Debian 6	<code>wget http://public-repo-1.hortonworks.com/HDP/deb6/2.1.5.0/hdp.list -O /etc/apt/sources.list.d/hdp.list</code>

- Create an HTTP server.
  1. On the mirror server, install an HTTP server (such as Apache `httpd`) using the instructions provided <http://httpd.apache.org/download.cgi>
  2. Activate this web server.
  3. Ensure that the firewall settings (if any) allow inbound HTTP access from your cluster nodes to your mirror server.



### Note

If you are using EC2, make sure that SELinux is disabled.

4. Optional - If your mirror server uses SLES, modify the `default-server.conf` file to enable the docs root folder listing.

```
sed -e "s/Options None/Options Indexes MultiViews/ig" /etc/apache2/default-server.conf > /tmp/tempfile.tmp
mv /tmp/tempfile.tmp /etc/apache2/default-server.conf
```

- On your mirror server, create a directory for your web server.

- For example, from a shell window, type:

- For RHEL/CentOS/Oracle:

```
mkdir -p /var/www/html/hdp/
```

- For SLES:

```
mkdir -p /srv/www/htdocs/rpms
```

- For Ubuntu and Debian:

```
mkdir -p /var/www/html/hdp/
```

- If you are using a symlink, enable the `followsymlinks` on your web server.
- Copy the contents of entire HDP repository for your desired OS from the remote yum server to your local mirror server.
- Continuing the previous example, from a shell window, type:

- For RHEL/CentOS/Oracle/Ubuntu:

```
cd /var/www/html/hdp
```

- For SLES:

```
cd /srv/www/htdocs/rpms
```

Then for all hosts, type:

- HDP Repository

```
reposync -r HDP
reposync -r HDP-2.1.5.0
reposync -r HDP-UTILS-1.1.0.19
```

You should see both an `HDP-2.1.5.0` directory and an `HDP-UTILS-1.1.0.19` directory, each with several subdirectories.

- Generate appropriate metadata.

This step defines each directory as a yum repository. From a shell window, type:

- For RHEL/CentOS/Oracle:

- HDP Repository:

```
createrepo /var/www/html/hdp/HDP-2.1.5.0
createrepo /var/www/html/hdp/HDP-UTILS-1.1.0.19
```

- For SLES:

- HDP Repository:

```
createrepo /srv/www/htdocs/rpms/hdp/HDP
```

You should see a new folder called `repodata` inside both HDP directories.

- Verify the configuration.
  - The configuration is successful, if you can access the above directory through your web browser.

To test this out, browse to the following location:

- HDP: [http://\\$yourwebserver/hdp/HDP-2.1.5.0/](http://$yourwebserver/hdp/HDP-2.1.5.0/)
- You should now see directory listing for all the HDP components.

- At this point, you can disable external Internet access for the mirror server, so that the mirror server is again entirely within your data center firewall.
- Depending on your cluster OS, configure the **yum** clients on all the nodes in your cluster
  1. Edit the repo files, changing the value of the `baseurl` property to the local mirror URL.
    - Edit the `/etc/yum.repos.d/hdp.repo` file, changing the value of the `baseurl` property to point to your local repositories based on your cluster OS.

```
[HDP-2.x]
name=Hortonworks Data Platform Version - HDP-2.x
baseurl=http://$yourwebserver/HDP/$os/2.x/GA
gpgcheck=1
gpgkey=http://public-repo-1.hortonworks.com/HDP/$os/RPM-GPG-KEY/RPM-GPG-KEY-Jenkins
enabled=1
priority=1

[HDP-UTILS-1.1.0.19]
name=Hortonworks Data Platform Utils Version - HDP-UTILS-1.1.0.19
baseurl=http://$yourwebserver/HDP-UTILS-1.1.0.19/repos/$os
gpgcheck=1
gpgkey=http://public-repo-1.hortonworks.com/HDP/$os/RPM-GPG-KEY/RPM-GPG-KEY-Jenkins
enabled=1
priority=1

[HDP-2.1.5.0]
name=Hortonworks Data Platform HDP-2.1.5.0
baseurl=http://$yourwebserver/HDP/$os/2.x/updates/2.1.5.0
gpgcheck=1
gpgkey=http://public-repo-1.hortonworks.com/HDP/$os/RPM-GPG-KEY/RPM-GPG-KEY-Jenkins
enabled=1
priority=1
```

where

- `$yourwebserver` is FQDN of your local mirror server.
- `$os` can be `centos5`, `centos6`, or `suse11`. Use the following options table for `$os` parameter:

**Table 4.7. Options for `$os` parameter in repo URL**

Operating System	Value
CentOS 5	centos5
RHEL 5	
Oracle Linux 5	
CentOS 6	centos6
RHEL 6	
Oracle Linux 6	
SLES 11	suse11
Ubuntu 12	ubuntu12

2. Copy the yum/zypper client configuration file to all nodes in your cluster.

- RHEL/CentOS/Oracle Linux:

Use **scp** or **pdsh** to copy the client yum configuration file to `/etc/yum.repos.d/` directory on every node in the cluster.

- For SLES:

On every node, invoke the following command:

- HDP Repository: **zypper addrepo -r http://\$yourwebserver/hdp/HDP/suse11/2.x/updates/2.1.5.0/hdp.repo**

- For Ubuntu:

On every node, invoke the following command:

- HDP Repository: **sudo add-apt-repository 'deb http://\$yourwebserver/hdp/HDP/ubuntu12/2.x/hdp.list'**

- Optional - Ambari Repository: **sudo add-apt-repository 'deb http://\$yourwebserver/hdp/ambari/ubuntu12/1.x/updates/1.5.1/ambari.list'**

- If using Ambari, verify the configuration by deploying Ambari server on one of the cluster nodes. **yum install ambari-server**

- If your cluster runs CentOS, Oracle, or RHEL and if you have multiple repositories configured in your environment, deploy the following plugin on all the nodes in your cluster.

1. Install the plugin.

- **For RHEL and CentOS v5.x**

```
yum install yum-priorities
```

- **For RHEL and CentOS v6.x**

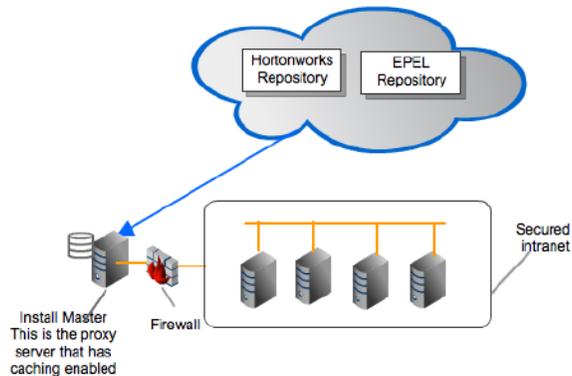
```
yum install yum-plugin-priorities
```

2. Edit the `/etc/yum/pluginconf.d/priorities.conf` file to add the following:

```
[main]
enabled=1
gpgcheck=0
```

### 4.3.3. Set up a trusted proxy server

A trusted proxy server is set up in the following illustration:



Complete the following instructions to set up a trusted proxy server:

1. [Check Your Prerequisites](#)
2. [Install the Repos](#)

### 4.3.3.1. Check Your Prerequisites

Select a mirror server host with the following characteristics:

- This server runs on either CentOS/RHEL/Oracle Linux (5.x or 6.x), SLES 11, or Ubuntu 12, and has several GB of storage available.
- The firewall allows all cluster nodes (the servers on which you want to install HDP) to access this server, and allows this server to access the Internet (at least those Internet servers for the repositories to be proxied).

### 4.3.3.2. Install the Repos

1. Create a caching HTTP PROXY server on the selected host.
  - a. It is beyond the scope of this document to show how to set up an HTTP PROXY server, given the many variations that may be required, depending on your data center's network security policy. If you choose to use the Apache HTTPD server, it starts by installing `httpd`, using the instructions provided [here](#), and then adding the `mod_proxy` and `mod_cache` modules, as stated [here](#).  
  
Please engage your network security specialists to correctly set up the proxy server.
  - b. Activate this proxy server and configure its cache storage location.
  - c. Ensure that the firewall settings (if any) allow inbound HTTP access from your cluster nodes to your mirror server, and outbound access to the desired repo sites, including `public-repo-1.hortonworks.com`.



#### Note

If you are using EC2, make sure that SELinux is disabled.

2. Depending on your cluster OS, configure the `yum` clients on all the nodes in your cluster.



## Note

The following description is taken from the CentOS documentation [here](#).

- a. On each cluster node, add the following lines to the `/etc/yum.conf` file.

(As an example, the settings below will enable **yum** to use the proxy server **mycache.mydomain.com**, connecting to port **3128**, with the following credentials **yum-user/qwerty**.)

```
# proxy server:port number
proxy=http://mycache.mydomain.com:3128
```

```
# account details for secure yum proxy connections
proxy_username=yum-user
proxy_password=qwerty
```

- b. Once all nodes have their `/etc/yum.conf` file updated with appropriate configuration info, you can proceed with the HDP installation just as though the nodes had direct access to the Internet repositories.
- c. If this proxy configuration does not seem to work, try adding a `/` at the end of the proxy URL. For example:

```
proxy=http://mycache.mydomain.com:3128/
```

## 5. Wire Encryption in Hadoop

Information on how over-the-wire encryption works in Hadoop, has been moved to the [Hadoop Security Guide](#).

## 6. Supported Database Matrix for Hortonworks Data Platform

This page contains certification information on supported databases for Hortonworks Data Platform (HDP).

The following table identifies the supported databases for HDP.

**Table 6.1. Supported Databases for HDP Stack**

Operating System	Component	Database				
		PostgreSQL 8.x	PostgreSQL 9.x	MySQL 5.x	Oracle 11gr2	Other
RHEL/Centos/ Oracle Linux 5.x  RHEL/CentOS/ Oracle Linux 6.x  SLES 11  Ubuntu 12	Hive / HCatalog	Supported. For instructions on configuring this database for Hive metastore, see <a href="#">Instructions for Manual Install</a> .	Supported. For instructions on configuring this database for Hive metastore, see <a href="#">Instructions for Manual Install</a> .	Default. For instructions on configuring this database for Hive metastore, either see <a href="#">Instructions for Manual Install</a> or see <a href="#">Using Ambari with MySQL</a> .	Default. For instructions on configuring this database for Hive metastore, either see <a href="#">Instructions for Manual Install</a> or see <a href="#">Using Ambari with Oracle</a> .	
	Oozie	Supported. For instructions on configuring this database for Oozie metastore, see <a href="#">Instructions for Manual Install</a> .	Supported. For instructions on configuring this database for Oozie metastore, see <a href="#">Instructions for Manual Install</a> .	Supported. For instructions on configuring this database for Oozie metastore, either see <a href="#">Instructions for Manual Install</a> or see <a href="#">Using Oozie with MySQL</a> .	Supported. For instructions on configuring this database for Oozie metastore, either see <a href="#">Instructions for Manual Install</a> or see <a href="#">Instructions for Oozie</a> .	Derby (default).
	Hue <sup>a</sup>	Supported. For instructions on configuring this database for Hue, see <a href="#">Instructions for Manual Install</a> .	Supported. For instructions on configuring this database for Hue, see <a href="#">Instructions for Manual Install</a> .	Supported. For instructions on configuring this database for Hue, see <a href="#">Instructions for Manual Install</a> .	Supported. For instructions on configuring this database for Hue, see <a href="#">Instructions for Manual Install</a> .	SQLite (default)
	Ambari <sup>b</sup>	Default. For more information, see <a href="#">Database Requirements</a> .	Supported. For more information, see <a href="#">Using Ambari with PostgreSQL</a> .	Supported. For more information, see <a href="#">Using Ambari with MySQL</a> .	Supported. For more information, see <a href="#">Using Ambari with Oracle</a> .	

<sup>a</sup>Hue does not currently support Ubuntu 12.

<sup>b</sup>Ambari does not currently support Ubuntu 12.