



Hadoop Rocks

March 12 2009

Michael Thomas



Assumptions

Rocks is already configured and running

You manage the fuse kernel modules yourself

You attach `hadoop.xml` to your desired node type

★ `/export/rocks/install/site-profiles/5.1/graphs/default/hadoop.xml`

Only one data directory per node

Namenode does not act as a data node



Components I



hadoop-0.9.0-4.el5.x86_64.rpm

- ★ **/opt/hadoop/***
- ★ **/etc/profile.d/hadoop.sh**
- ★ **/etc/init.d/hadoop**
- ★ **/etc/sysconfig/hadoop**
- ★ **/usr/bin/hadoop-to-hostname.sh**
- ★ **/usr/bin/hdfs**
- ★ **/usr/bin/getPoolSize.sh**

hadoop-firstboot-0.1-3.el5.noarch.rpm

- ★ **/etc/init.d/hadoop-firstboot**
- ★ **Requires: hadoop**
- ★ **Modifies /opt/hadoop/conf/***



Components II



`fuse-2.7.4-8_10.el5.x86_64.rpm`

`fuse-libs-2.7.4-8_10.el5.x86_64.rpm`

`(fuse-kmdl-2.6.18-91.1.22.el5-2.7.4-8_10.el5.x86_64.rpm)`

`/export/rocks/install/site-profiles/5.1/nodes/hadoop.xml`

`gridftp-hdfs-1.0.0-1.10.1u.el5.1.x86_64.rpm`

`/export/rocks/install/site-profiles/5.1/nodes/gridftp.xml`



hadoop.xml



```
<package>hadoop-firstboot</package>
```

```
<package>hadoop</package>
```

```
<package>fuse</package>
```

```
<package>fuse-libs</package>
```

```
<post>
```

```
chkconfig hadoop-firstboot on
```

```
chkconfig hadoop on
```

```
# Update a few site-specific configurations
```

```
sed -i -e 's/^HADOOP_NAMENODE=.* /HADOOP_NAMENODE=compute-13-1' \  
    /etc/sysconfig/hadoop
```

```
# Create the fuse mount point and mount the hdfs automatically
```

```
mkdir /mnt/hadoop
```

```
echo "hdfs# /mnt/hadoop fuse server=compute-13-  
    1,port=9000,rdbuffer=131072,allow_other 0 0" >> /etc/fstab
```

```
</post>
```



How to make it work

- **Download files from** <http://ultralight.caltech.edu/~wart/hadoop/>
- **Copy rpms to** `/export/rocks/install/contrib/5.1/RPMS/x86_64`
- **Copy hadoop.xml to** `/export/rocks/install/site-profiles/5.1/nodes/hadoop.xml`
- **Modify namenode hostname, datanode data directory in** `hadoop.xml`
- **Copy gridftp.xml to** `/export/rocks/install/site-profiles/5.1/nodes/hadoop.xml`
- **Modify namenode hostname, datanode data directory in** `gridftp.xml`
- **Create** `/export/rocks/install/site-profiles/5.1/graphs/default/{gridftp,hadoop}.xml`
- **rocks create distro**
- **Kickstart the namenode**
- **Kickstart data/gridftp nodes**



TODO



Use generic values as defaults in rpm to fish out caltech site-specific settings

- ★ Find single place in rocks to define site-specific settings

Package BestMan and gridftp as rpms

Build hadoop rpm from source

Finish hadoop ganglia+MonALISA integration

Make a hadoop roll