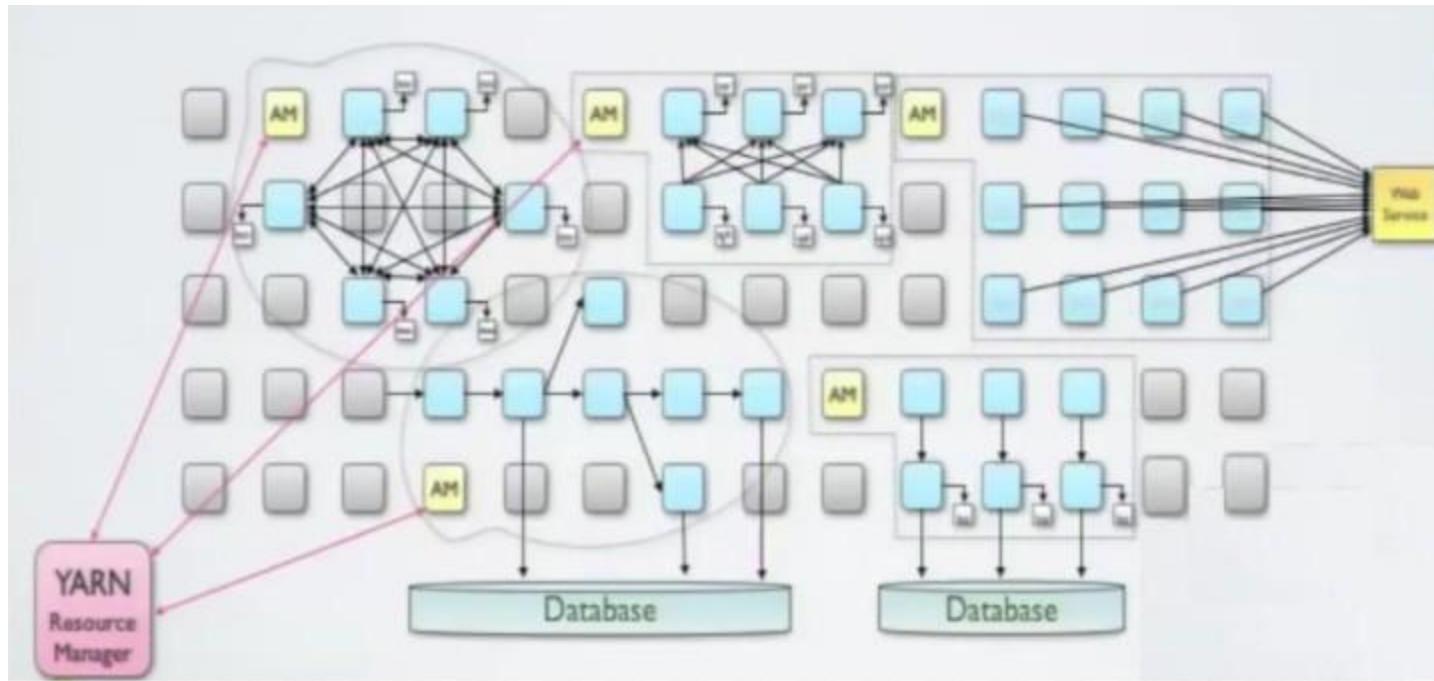


SE256 :Scalable Systems for Data Science

Lab Session : 1

YARN Cluster



YARN Config.

1. The hostname of the RM
2. Memory to allocate to each container request at the Resource Manager
3. CPU cores to allocate to each container request at the Resource Manager
4. container related settings

HDFS Config.

1. DataNode : Paths on the local filesystem where it should store its blocks
2. NameNode : Path on the local filesystem where the NameNode stores the namespace and transaction logs persistently.
3. NameNode URI

JDK Installation

Hadoop requires a working Java 5+ (Java 1.5+) installation. We now describe the installation of OpenJDK 7.

sudo apt-get install openjdk-7-jdk [Ubuntu]

su -c "yum install java-1.7.0-openjdk-devel" [Fedora]

```
dream@dream:~$ sudo apt-get install openjdk-7-jdk
Reading package lists... Done
Building dependency tree
Reading state information... Done
The following extra packages will be installed:
  acl at-spi2-core ca-certificates-java colord cpp cpp-4.8 dbus-x11
  dconf-gsettings-backend dconf-service desktop-file-utils fontconfig
  fontconfig-config fonts-dejavu-core fonts-dejavu-extra gconf-service
  gconf-service-backend gconf2 gconf2-common gvfs gvfs-common gvfs-daemons
  gvfs-libs hicolor-icon-theme java-common libasound2 libasound2-data
  libasyncns0 libatasmart4 libatk-bridge2.0-0 libatk-wrapper-java
  libatk-wrapper-java-jni libatk1.0-0 libatk1.0-data libatspi2.0-0
  libavahi-glib1 libbonobo2-0 libbonobo2-common libcairo-gobject2 libcairo2
  libcanberra0 libcloog-isl4 libcolorl1 libcolorhug1 libdatrie1 libdconf1
  libdrm-intel1 libdrm-nouveau2 libdrm-radeon1 libexif12 libflac8
  libfontconfig1 libfontenc1 libgconf-2-4 libgconf2-4 libgd3
  libgdk-pixbuf2.0-0 libgdk-pixbuf2.0-common libgif4 libgl1-mesa-dri
  libgl1-mesa-glx libglapi-mesa libgnome2-0 libgnome2-bin libgnome2-common
  libgnomevfs2-0 libgnomevfs2-common libgphoto2-6 libgphoto2-l10n
  libgphoto2-port10 libgraphite2-3 libgtk-3-0 libgtk-3-bin libgtk-3-common
  libgtk2.0-0 libgtk2.0-bin libgtk2.0-common libgudev-1.0-0 libgusb2
  libharfbuzz0b libice-dev libice6 libidl-common libidl0 libieee1284-3
  libisl10 libjasper1 libjbig0 libjpeg-turbo8 libjpeg8 liblcms2-2 libllvm3.4
  libltdl7 libmpc3 libmpfr4 libnssr4 libnss3 libnss3-nssdb libogg0
  liborbit-2-0 liborbit2 libpango-1.0-0 libpangocairo-1.0-0 libpangoft2-1.0-0
  libpciaaccess0 libpcsc-lite1 libpixman-1-0 libpthread-stubs0-dev libpulse0
  libsane libsane-common libscrypt1 libsecret-1-0 libsecret-common libsm-dev
  libsm6 libsndfile1 libthai-data libthai0 libtiff5 libtxc-dxtn-s2tc0
```

JDK Installation

It is also necessary to add the following lines to the end of the file **\$HOME/.bashrc**.

[For **64-bit** Linux]

```
export JAVA_HOME=/usr/lib/jvm/java-7-openjdk-amd64  
export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:$JAVA_HOME/jre/lib/amd64/server
```

[For **32-bit** Linux]

```
export JAVA_HOME=/usr/lib/jvm/java-7-openjdk-i386  
export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:$JAVA_HOME/jre/lib/i386/server
```

Compile the file with the command **source \$HOME/.bashrc**.

SSH Configuration

```
sudo apt-get install openssh-server [Ubuntu]
```

```
su -c "yum install openssh-server" [Fedora]
```

```
dream@dream:~$ sudo apt-get install openssh-server
Reading package lists... Done
Building dependency tree
Reading state information... Done
openssh-server is already the newest version.
0 upgraded, 0 newly installed, 0 to remove and 0 not upgraded.
dream@dream:~$ █
```

configure the password-less SSH connection. First, create an RSA key pair

```
$ ssh-keygen -t rsa
```

This creates files **id_rsa** and **id_rsa.pub** under the default directory **\$HOME/.ssh/**. Then, we copy the public key **id_rsa.pub** to another file **authorized_keys**:

```
$ cat $HOME/.ssh/id_rsa.pub >> $HOME/.ssh/authorized_keys
```

Downloading hadoop

Download hadoop using wget command:

```
wget http://www.eu.apache.org/dist/hadoop/common/hadoop-2.6.0/hadoop-2.6.0.tar.gz
```

```
dream@dream ~$ wget http://www.eu.apache.org/dist/hadoop/common/hadoop-2.6.0/hadoop-2.6.0.tar.gz
--2016-01-04 16:32:16--  http://www.eu.apache.org/dist/hadoop/common/hadoop-2.6.0/hadoop-2.6.0.tar.gz
Resolving www.eu.apache.org (www.eu.apache.org)... 88.198.26.2, 2a01:4f8:130:2192::2
Connecting to www.eu.apache.org (www.eu.apache.org) | 88.198.26.2|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 195257604 (186M) [application/x-gzip]
Saving to: 'hadoop-2.6.0.tar.gz'

4% [=====] 8,128,024      1.62MB/s eta 2m 36s
```

Extracting files

```
tar xvf hadoop-2.6.0.tar.gz --gzip
```

```
hadoop-2.6.0/bin/test-container-executor  
hadoop-2.6.0/bin/container-executor  
hadoop-2.6.0/bin/hadoop.cmd  
hadoop-2.6.0/bin/rcc  
hadoop-2.6.0/bin/hdfs  
hadoop-2.6.0/bin/namenode  
hadoop-2.6.0/bin/hadoop  
hadoop-2.6.0/bin/yarn.cmd  
hadoop-2.6.0/bin/nodemanager.cmd  
hadoop-2.6.0/bin/yarn  
hadoop-2.6.0/include/  
hadoop-2.6.0/include/TemplatenFactor.y.hh  
hadoop-2.6.0/include/StringUtils.hh  
hadoop-2.6.0/include/hdfs.h  
hadoop-2.6.0/include/Pipes.hh  
hadoop-2.6.0/include/SerialUtils.hh  
drream@drream ~$ ls  
hadoop-2.6.0 hadoop-2.6.0.tar.gz  
drream@drream ~$ rm hadoop-2.6.0.tar.gz  
drream@drream ~$ tar xvf hadoop-2.6.0.tar.gz --gzip
```

Setting hadoop env variables

```
export HADOOP_PREFIX="/home/dream/hadoop-2.6.0" # Change this to where you unpacked hadoop to.

export HADOOP_HOME=$HADOOP_PREFIX
export HADOOP_COMMON_HOME=$HADOOP_PREFIX
export HADOOP_CONF_DIR=$HADOOP_PREFIX/etc/hadoop
export HADOOP_HDFS_HOME=$HADOOP_PREFIX
export HADOOP_MAPRED_HOME=$HADOOP_PREFIX
export HADOOP_YARN_HOME=$HADOOP_PREFIX
```

```
export HADOOP_PREFIX="/home/dream/hadoop-2.6.0" # Change this to where you unpacked hadoop to.
export HADOOP_HOME=$HADOOP_PREFIX
export HADOOP_COMMON_HOME=$HADOOP_PREFIX
export HADOOP_CONF_DIR=$HADOOP_PREFIX/etc/hadoop
export HADOOP_HDFS_HOME=$HADOOP_PREFIX
export HADOOP_MAPRED_HOME=$HADOOP_PREFIX
export HADOOP_YARN_HOME=$HADOOP_PREFIX
```

Making directories required for hdfs

```
dream@dream:~$ mkdir -p hadoop-2.6.0/hdfs/datanode/ hadoop-2.6.0/hdfs/namenode  
dream@dream:~$ ls hadoop-2.6.0/hdfs/  
datanode namenode
```

Configuring hdfs-site.xml

```
File Edit View Search Terminal Help
GNU nano 2.2.6          File: /home/dream/hadoop-2.6.0/etc/hadoop/hdfs-site.xml          Modified

<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!--
Licensed under the Apache License, Version 2.0 (the "License");
you may not use this file except in compliance with the License.
You may obtain a copy of the License at

http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License. See accompanying LICENSE file.
-->

<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
  <name>dfs.datanode.data.dir</name>
  <value>file:///home/dream/hadoop-2.6.0/hdfs/datanode</value>
  <description>Comma separated list of paths on the local filesystem of a DataNode where it should store its blocks.</description>
</property>

<property>
  <name>dfs.namenode.name.dir</name>
  <value>file:///home/dream/hadoop-2.6.0/hdfs/namenode</value>
  <description>Path on the local filesystem where the NameNode stores the namespace and transaction logs persistently.</description>
</property>
</configuration>
```

Configuring core-site.xml

File Edit View Search Terminal Help

GNU nano 2.2.6

File: /home/dream/hadoop-2.6.0/etc/hadoop/core-site.xml

```
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!--
 Licensed under the Apache License, Version 2.0 (the "License");
 you may not use this file except in compliance with the License.
 You may obtain a copy of the License at

 http://www.apache.org/licenses/LICENSE-2.0

 Unless required by applicable law or agreed to in writing, software
 distributed under the License is distributed on an "AS IS" BASIS,
 WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
 See the License for the specific language governing permissions and
 limitations under the License. See accompanying LICENSE file.
-->

<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
    <name>fs.defaultFS</name>
    <value>hdfs://localhost/</value>
    <description>NameNode URI</description>
</property>
</configuration>
```

Configuring yarn-site.xml

```
GNU nano 2.2.6          File: /home/dream/hadoop-2.6.0/etc/hadoop/yarn-site.xml          Modified

:xml version="1.0"?&gt;
&lt;configuration&gt;

<!-- Site specific YARN configuration properties --&gt;
&lt;property&gt;
    &lt;name&gt;yarn.scheduler.minimum-allocation-mb&lt;/name&gt;
    &lt;value&gt;128&lt;/value&gt;
    &lt;description&gt;Minimum limit of memory to allocate to each container request at the Resource Manager.&lt;/description&gt;
&lt;/property&gt;
&lt;property&gt;
    &lt;name&gt;yarn.scheduler.maximum-allocation-mb&lt;/name&gt;
    &lt;value&gt;2048&lt;/value&gt;
    &lt;description&gt;Maximum limit of memory to allocate to each container request at the Resource Manager.&lt;/description&gt;
&lt;/property&gt;
&lt;property&gt;
    &lt;name&gt;yarn.scheduler.minimum-allocation-vcores&lt;/name&gt;
    &lt;value&gt;1&lt;/value&gt;
    &lt;description&gt;The minimum allocation for every container request at the RM, in terms of virtual CPU cores. Requests lower than t$&lt;/description&gt;
&lt;/property&gt;
&lt;property&gt;
    &lt;name&gt;yarn.scheduler.maximum-allocation-vcores&lt;/name&gt;
    &lt;value&gt;2&lt;/value&gt;
    &lt;description&gt;The maximum allocation for every container request at the RM, in terms of virtual CPU cores. Requests higher than $&lt;/description&gt;
&lt;/property&gt;
&lt;property&gt;
    &lt;name&gt;yarn.nodemanager.resource.memory-mb&lt;/name&gt;
    &lt;value&gt;4096&lt;/value&gt;
    &lt;description&gt;Physical memory, in MB, to be made available to running containers&lt;/description&gt;
&lt;/property&gt;
&lt;property&gt;
    &lt;name&gt;yarn.nodemanager.resource.cpu-vcores&lt;/name&gt;</pre
```

Starting

Now that we've finished configuring everything, it's time to setup the folders and start the daemons:

```
# Format the namenode directory (DO THIS ONLY ONCE, THE FIRST TIME)
$HADOOP_PREFIX/bin/hdfs namenode -format

# Start the namenode daemon
$HADOOP_PREFIX/sbin/hadoop-daemon.sh start namenode

# Start the datanode daemon
$HADOOP_PREFIX/sbin/hadoop-daemon.sh start datanode

## Start YARN daemons

# Start the resourcemanager daemon
$HADOOP_PREFIX/sbin/yarn-daemon.sh start resourcemanager

# Start the nodemanager daemon
$HADOOP_PREFIX/sbin/yarn-daemon.sh start nodemanager
```

```
dream@dream:~$ $HADOOP_PREFIX/bin/hdfs namenode -format
16/01/04 16:54:10 INFO namenode.NameNode: STARTUP_MSG:
*****STARTUP_MSG: Starting NameNode
STARTUP_MSG: host = dream/10.12.0.4
STARTUP_MSG: args = [-format]
STARTUP_MSG: version = 2.6.0
STARTUP_MSG: classpath = /home/dream/hadoop-2.6.0/etc/hadoop:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/commons-collections-3.2
l.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/commons-beanutils-core-1.8.0.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/
'jetty-6.1.26.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/servlet-api-2.5.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/
avro-1.7.4.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/netty-3.6.2.Final.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/
ava-xmlobjbuilder-0.4.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/apacheds-kerberos-codec-2.0.0-M15.jar:/home/dream/hadoop-2.6.0/
share/hadoop/common/lib/jersey-server-1.9.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/commons-beanutils-1.7.0.jar:/home/dream/hadoop-2.6.0/
share/hadoop/common/lib/commons-clients-3.1.0.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/commons-io-2.4.jar:/home/dream/hadoop-2.6.0/
share/hadoop/common/lib/junit-4.11.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/mockito-all-1.8.5.jar:/home/dream/hadoop-2.6.0/
share/hadoop/common/lib/api-util-1.0.0-M20.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/commons-net-3.1.jar:/home/dream/hadoop-2.6.0/
share/hadoop/common/lib/jsch-0.1.42.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/commons-configuration-1.6.jar,
/home/dream/hadoop-2.6.0/share/hadoop/common/lib/guava-11.0.2.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/hadoop-annotations-2.0.
0.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/jasper-compiler-5.5.23.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/com-
mons-el-1.0.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/xmlenc-0.52.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/common-
.lang-2.6.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/activation-1.1.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/jsr31-
3.1.3.9.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/xz-1.0.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/hadoop-rest-core-
3.3.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/slf4j-api-1.7.5.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/api-asn1-ap-
i-1.0.0-M20.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/asm-3.2.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/jackson-m-
oper-asl-1.9.13.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/commons-digester-1.8.jar:/home/dream/hadoop-2.6.0/share/hadoop/com-
mon/lib/jaxb-impl-2.2.3-1.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/apacheds-118n-2.0.0-M15.jar:/home/dream/hadoop-2.6.0/share/
'hadoop/common/lib/httpcore-4.2.5.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/slf4j-log4j12-1.7.5.jar:/home/dream/hadoop-2.6.0/
share/hadoop/common/lib/Log4j-1.2.17.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/curator-recipes-2.6.0.jar:/home/dream/hadoop-2.
6.0/share/hadoop/common/lib/jasper-runtime-5.5.23.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/snappy-java-1.0.4.1.jar:/home/dre-
am/hadoop-2.6.0/share/hadoop/common/lib/protobuf-java-2.5.0.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/jackson-xc-1.9.13.jar
/home/dream/hadoop-2.6.0/share/hadoop/common/lib/httpclient-4.2.5.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/ibus-2.2.2.ja-
r:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/commons-math3-3.1.1.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/zookeeper-3
1.6.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/stax-api-1.0-2.jar:/home/dream/hadoop-2.6.0/share/hadoop/common/lib/commons-1.0-
```

Starting Hadoop daemon process

File Edit View Search Terminal Help

```
dream@dream:~$ $HADOOP_PREFIX/sbin/hadoop-daemon.sh start namenode
starting namenode, logging to /home/dream/hadoop-2.6.0/logs/hadoop-dream-namenode-dream.out
dream@dream:~$ jps
11042 Jps
10967 NameNode
dream@dream:~$ $HADOOP_PREFIX/sbin/hadoop-daemon.sh start datanode
starting datanode, logging to /home/dream/hadoop-2.6.0/logs/hadoop-dream-datanode-dream.out
dream@dream:~$ jps
11144 Jps
10967 NameNode
11065 DataNode
dream@dream:~$ $HADOOP_PREFIX/sbin/yarn-daemon.sh start resourcemanager
starting resourcemanager, logging to /home/dream/hadoop-2.6.0/logs/yarn-dream-resourcemanager-dream.out
dream@dream:~$ jps
11403 Jps
11173 ResourceManager
10967 NameNode
11065 DataNode
dream@dream:~$ $HADOOP_PREFIX/sbin/yarn-daemon.sh start nodemanager
starting nodemanager, logging to /home/dream/hadoop-2.6.0/logs/yarn-dream-nodemanager-dream.out
dream@dream:~$ jps
11431 NodeManager
11173 ResourceManager
10967 NameNode
11525 Jps
11065 DataNode
dream@dream:~$ █
```

File Edit View Search Terminal Help

```
dream@dream:~/hadoop-2.6.0$ $HADOOP_PREFIX/bin/hadoop jar ./share/hadoop/mapreduce/hadoop-mapreduce-examples-2.6.0.jar pi 10 100
Number of Maps = 10
Samples per Map = 100
Wrote input for Map #0
Wrote input for Map #1
Wrote input for Map #2
Wrote input for Map #3
Wrote input for Map #4
Wrote input for Map #5
Wrote input for Map #6
Wrote input for Map #7
Wrote input for Map #8
Wrote input for Map #9
Starting Job
16/01/04 16:59:45 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
16/01/04 16:59:45 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName=JobTracker, sessionId=
16/01/04 16:59:45 INFO input.FileInputFormat: Total input paths to process : 10
16/01/04 16:59:45 INFO mapreduce.JobSubmitter: number of splits:10
16/01/04 16:59:46 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local183546265_0001
16/01/04 16:59:46 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
16/01/04 16:59:46 INFO mapreduce.Job: Running job: job_local183546265_0001
16/01/04 16:59:46 INFO mapred.LocalJobRunner: OutputCommitter set in config null
16/01/04 16:59:46 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
16/01/04 16:59:46 INFO mapred.LocalJobRunner: Waiting for map tasks
16/01/04 16:59:46 INFO mapred.LocalJobRunner: Starting task: attempt_local183546265_0001_m_000000_0
16/01/04 16:59:46 INFO mapred.Task: Using ResourceCalculatorProcessTree : [ ]
16/01/04 16:59:46 INFO mapred.MapTask: Processing split: hdfs://localhost/user/dream/QuasiMonteCarlo_1451926782358_1678543576/in/part0
+118
16/01/04 16:59:46 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
16/01/04 16:59:46 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
16/01/04 16:59:46 INFO mapred.MapTask: soft limit at 83886080
16/01/04 16:59:46 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
16/01/04 16:59:46 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
16/01/04 16:59:46 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
16/01/04 16:59:46 INFO mapred.LocalJobRunner:
16/01/04 16:59:46 INFO mapred.MapTask: Starting flush of map output
```

HDFS: Number of write operations=135

Map-Reduce Framework

Map input records=10
Map output records=20
Map output bytes=180
Map output materialized bytes=280
Input split bytes=1420
Combine input records=0
Combine output records=0
Reduce input groups=2
Reduce shuffle bytes=280
Reduce input records=20
Reduce output records=0
Spilled Records=40
Shuffled Maps =10
Failed Shuffles=0
Merged Map outputs=10
GC time elapsed (ms)=77
CPU time spent (ms)=0
Physical memory (bytes) snapshot=0
Virtual memory (bytes) snapshot=0
Total committed heap usage (bytes)=5345640448

Shuffle Errors

BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0

File Input Format Counters

Bytes Read=1180

File Output Format Counters

Bytes Written=97

Job Finished in 2.913 seconds

Estimated value of Pi is 3.14800000000000000000000000

██████████ /███████ ██████████ ███

Demo

WordCount Example

-log

-output