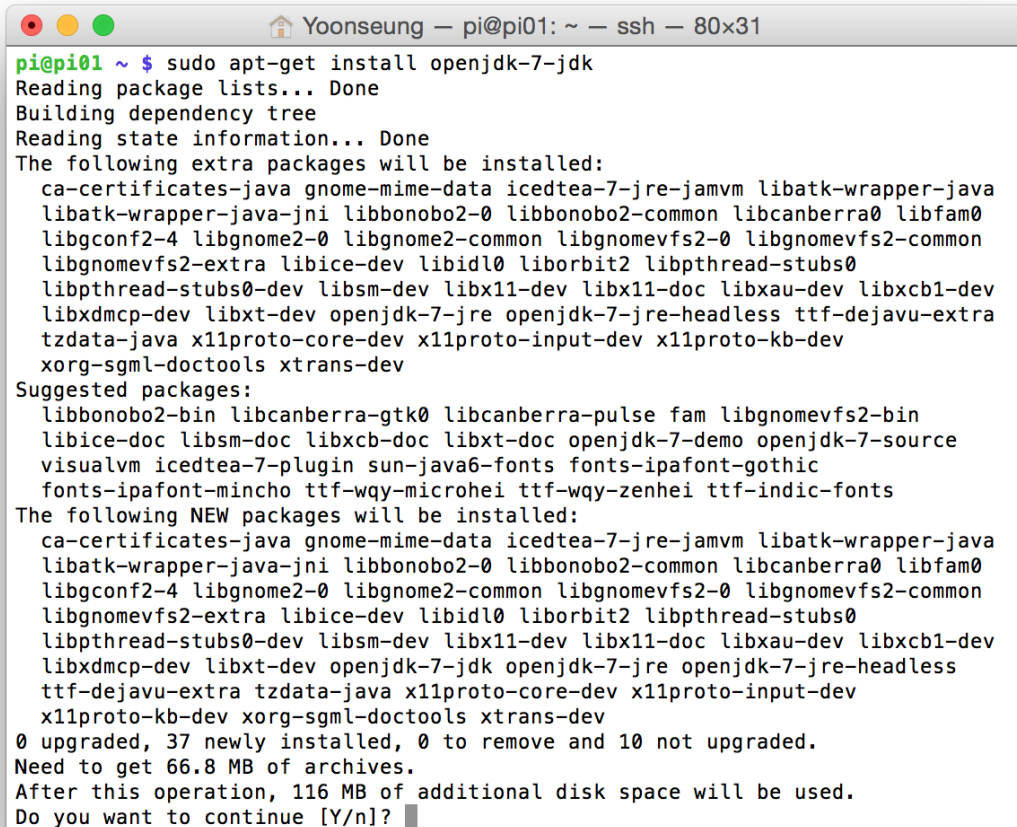# Hadoop Cluster

installation

# Outline

- For Raspberry Pi
  - Install JAVA & HADOOP
  - Set Path
  - Wordcount Example
- For Ubuntu 14.04

# For Raspberry Pi

Hadoop MapReduce Cluster installation

# Install JAVA Developer Kit

• sudo apt-get install openjdk-7-jdk

# Install Hadoop

- wget <link for hadoop>

- sudo tar vxzf <filename> -C /usr/local



ex >
wget https://archive.apache.org/dist/hadoop/common/hadoop-1.2.1/hadoop-1.2.1.tar.gz
sudo tar vxzf hadoop-1.2.1.tar.gz  -C /usr/local

# Set Path

- export JAVA_HOME=/usr/lib/jvm/java-7-openjdk-armhf

- export HADOOP_HOME=/usr/local/<hadoop ver>

```
pi@pi02 ~ $ export JAVA_HOME=/usr/lib/jvm/java-7-openjdk-armhf
pi@pi02 ~ $ export HADOOP_HOME=/usr/local/hadoop-1.2.1/
pi@pi02 ~ $ echo $JAVA_HOME
/usr/lib/jvm/java-7-openjdk-armhf
pi@pi02 ~ $ echo $HADOOP_HOME
/usr/local/hadoop-1.2.1/
pi@pi02 ~ $
```

# Set Path

- export PATH=$PATH:$HADOOP_HOME/bin

# Wordcount example
for single node

- hadoop jar
  $HADOOP_HOME/hadoop-
  examples-1.2.1.jar
  wordcount
  $HADOOP_HOME/README.
  txt ~/wordcount-opt

- cat ~/wordcount-opt/part-r-
  00000

# Distribute System [Master]

- sudo vi /etc/ssh/sshd_config
  *PubkeyAuthentication yes AuthorizedKeysFile
  .ssh/authorized_keys*

- mkdir ~/.ssh

- ssh-keygen -t rsa -P ""

- cp /home/pi/.ssh/id_rsa.pub
  /pi/stat/.ssh/authorized_keys

# For Ubuntu 14.04

Hadoop MapReduce Cluster installation

# configuration

- **Console mode booting**
  *$ sudo vi /etc/default/grub*
  **-> changes these following lines**
  *GRUB_CMDLINE_LINUX_DEFAULT=""*
  *GRUB_CMDLINE_LINUX="text"*
  **-> after that, update conf & reboot**
  *$ sudo update-grub*
  *$ sudo reboot*

- **Root password setting**
  *$ sudo passwd root*

# configuration

- **Install java**
  $ *sudo add-apt-repository ppa:webupd8team/java*
  $ *sudo apt-get update*
  $ *sudo apt-get install oracle-jdk7-installer*

- **Download & Install Hadoop**
  $ *wget <http link>*
  $ *cp hadoop-x.x.x.tar.gz  /usr/local*
  $ *rm hadoop-x.x.x.tar.gz*
  $ *cd /usr/local*
  $ *tar zxvf hadoop-x.x.x.tar.gz*

# configuration

- **PATH setting**
  *$ sudo vi ~/.profile*
  **-> add these following lines**
  *export JAVA_HOME=/usr/lib/jvm/java-7-oracle*
  *export HADOOP_HOME=/usr/local/hadoop-1.2.1*
  *export PATH=$PATH:$HADOOP_HOME/bin*
  **-> add path by following command**
  *$ source ~/.profile*

- **Check the correct PATH by following command**
  *$ echo $HADOOP_HOME*

# configuration

- **Install ssh software**
  *$ sudo apt-get install ssh*
  *$ sudo apt-get install rsync*

- **Setup passphraseless ssh**
  *$ ssh-keygen -t dsa -P '' -f ~/.ssh/id_dsa*
  *$ cat ~/.ssh/id_dsa.pub >> ~/.ssh/authorized_keys*

# MapReduce configuration

# mapred-site.xml

- io.sort.factor (default value= 10)
  : The number of streams to merge at once while sorting files. This determines the number of open file handles.

- io.sort.mb (default value= 100)
  : The total amount of buffer memory to use while sorting files, in megabytes. By default, gives each merge stream 1MB, which should minimize seeks.