# Hadoop setup / Spark with HDFS

Lab 2
October 19th, 2017

Jonghyun Bae(jonghbae@snu.ac.kr)

Computer Science and Engineering

Seoul National University

# Index

- **Hadoop distributed file system (HDFS)**

- **Installation**

- **Configuration**

- **Starting / Stopping cluster**

- **Hadoop Web Interface**

- **Running Spark with HDFS**

- **Exercise**

# Before we start...

■ **Please connect your VM using SSH**

■ **https://docs.google.com/spreadsheets/d/1X9Uavr2PACqgfLC3 rOcQ7Gqo4-NQNKoBhvcaL_86g9E/edit?usp=sharing**

```
1  # Please your public IP address in xxx.xxx.xxx.xxx

2  student@computer:~$ ssh -X -i bde3.pem ubuntu@xxx.xxx.xxx.xxx

3  Welcome to Ubuntu 14.04.5 LTS (GNU/Linux 3.13.0-125-generic x86_64)

4  [...snipp...]

5  ubuntu@ip-x-x-x:~$
```

# Hadoop distributed file system (1)
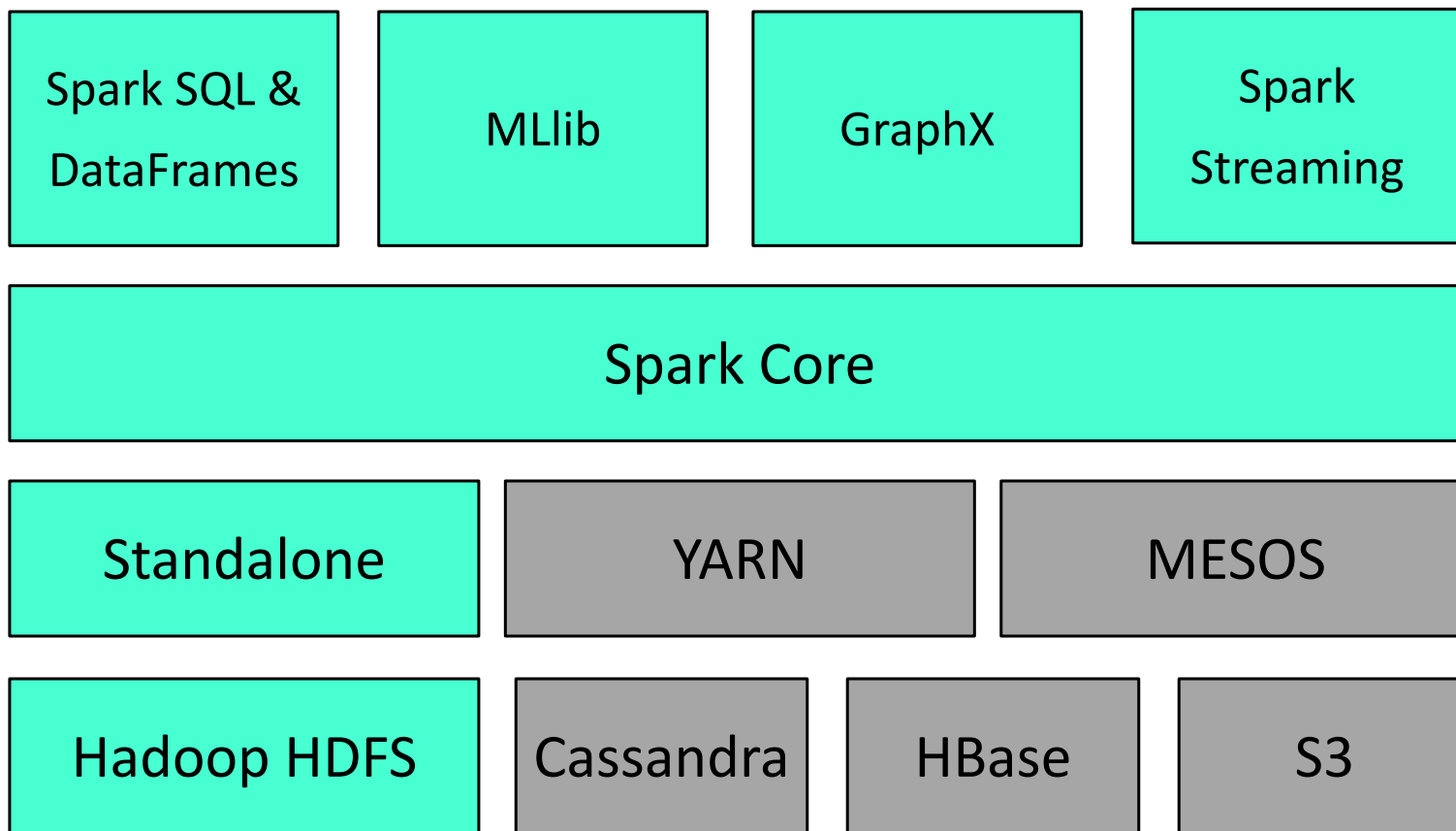
- **What is Hadoop?**



  - Developed by Doug Cutting and Mike Cafarella in 2006
  - Open-source software for reliable, scalable, distributed computing

# Hadoop distributed file system (2)

- **Why do we need Hadoop?**
  - Spark needs file system for saving / loading data

| Spark SQL & DataFrames | MLlib | GraphX | Spark Streaming |
|:---:|:---:|:---:|:---:|

| Spark Core |
|:---:|

| Standalone | YARN | MESOS |
|:---:|:---:|:---:|

| Hadoop HDFS | Cassandra | HBase | S3 |
|:---:|:---:|:---:|:---:|

* Image from https://www.safaribooksonline.com/library/view/data-analytics-with/9781491913734/ch04.html

# Installation (1)

- **Download Hadoop from <u>Apache Download Mirrors</u>**

```
   ubuntu@ip-x-x-x:~$ wget
1  http://mirror.navercorp.com/apache/hadoop/common/hadoop-2.7.4/hadoop-
   2.7.4.tar.gz
```

- **Unzip hadoop-2.7.4.tar.gz**

```
1  ubuntu@ip-x-x-x:~$ tar xzf hadoop-2.7.4.tar.gz

2  ubuntu@ip-x-x-x:~$ cd hadoop-2.7.4

3  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ cd

4  ubuntu@ip-x-x-x:~$
```

# Installation (2)

■ **Update $HOME/.bashrc**

```
1   ubuntu@ip-x-x-x:~$ vi ~/.bashrc
```

| $HOME/.bashrc |
| --- |

```
1   # ~/.bashrc: executed by bash(1) for non-login shells.

(...)                        (...)

2   export JAVA_HOME=/usr/lib/jvm/java-8-oracle

3   export SPARK_HOME=/home/ubuntu/spark-2.1.0

4   # Add new environment variable for Hadoop!

5   export HADOOP_HOME=/home/ubuntu/hadoop-2.7.4
```

# Installation (3)
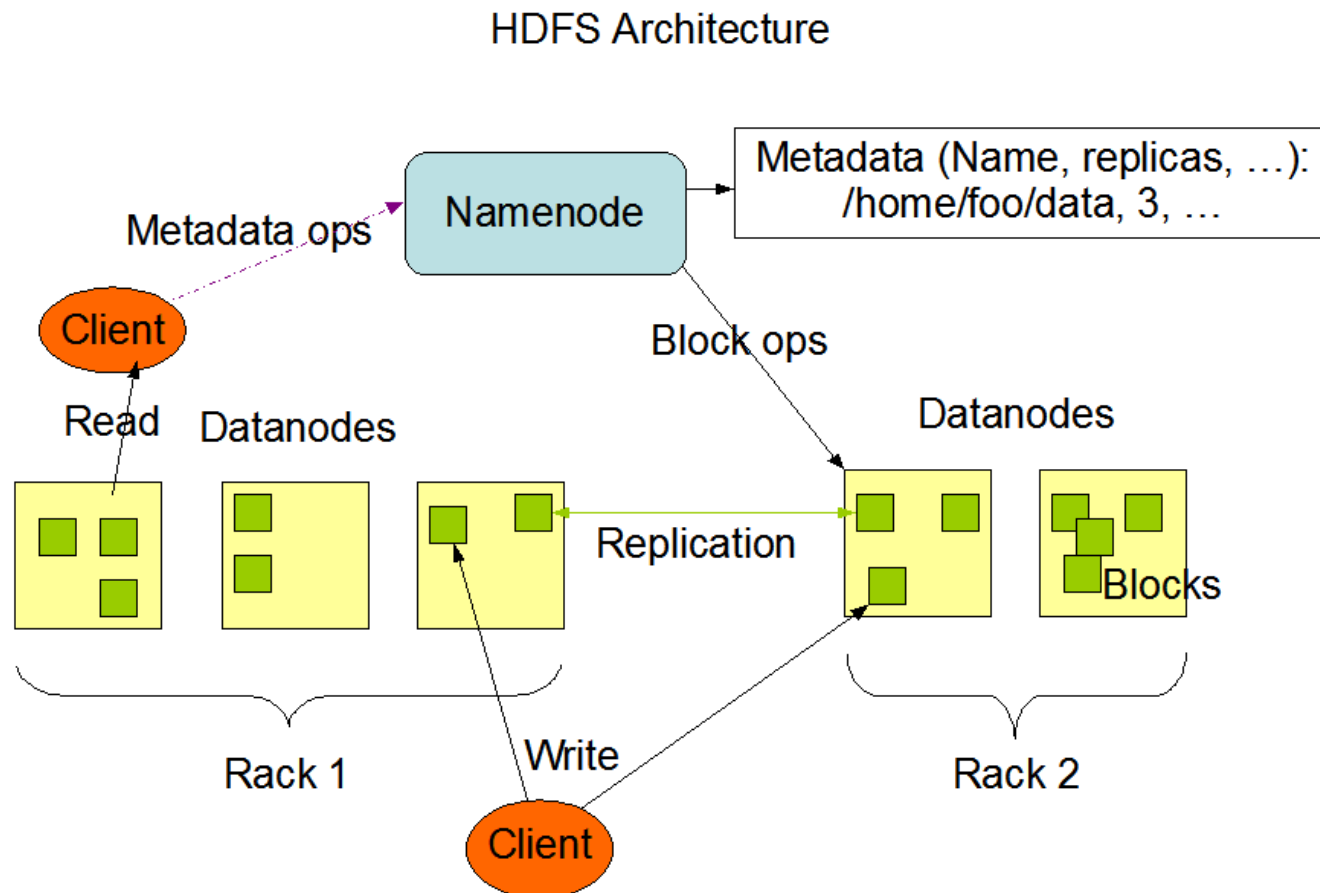
- **Apply changed setup in $HOME/.bashrc**

```
1   ubuntu@ip-x-x-x:~$ source ~/.bashrc

2   ubuntu@ip-x-x-x:~$ cd $HADOOP_HOME

3   ubuntu@ip-x-x-x:~/hadoop-2.7.4$
```

# Configuration (1)

■ **Architecture of namenode and datanode**

HDFS Architecture



Metadata ops

Namenode

Metadata (Name, replicas, …):
/home/foo/data, 3, …

Client

Read

Datanodes

Block ops

Datanodes

Replication

Blocks

Rack 1

Write

Rack 2

Client

* Image from http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HdfsDesign.html

# Configuration (2)

■ **hadoop-env.sh**

　　■ Set your default configuration for Hadoop

```
1  ubuntu@ip-x-x-x:~$ vi $HADOOP_HOME/etc/hadoop/hadoop-env.sh
```

$HADOOP_HOME/etc/hadoop/hadoop-env.sh

```
1  # The java implementation to use.

2  export JAVA_HOME=/usr/lib/jvm/java-8-oracle

3

4  # Set your Hadoop configuration directory

5  export HADOOP_CONF_DIR=/home/ubuntu/hadoop-2.7.4/etc/hadoop
```

# Configuration (3)

■ **core-site.xml**

- ▪ You can set cluster information for master and slave model
- ▪ Write properties between `<configuration>` and `</configuration>`

```
1   ubuntu@ip-x-x-x:~$ vi $HADOOP_HOME/etc/hadoop/core-site.xml
```

$HADOOP_HOME/etc/hadoop/core-site.xml

```
1   <property>
2       <name>fs.defaultFS</name>
3       <value>hdfs://localhost:9000</value>
4   </property>
```

# Configuration (4)

■ **hdfs-site.xml**

- ▪ You can set internal HDFS information for namenode and datanode
- ▪ Write properties between `<configuration>` and `</configuration>`

| $HADOOP_HOME/etc/hadoop/hdfs-site.xml |
|---|

```
1  <property>
2      <name>dfs.replication</name>
3      <value>1</value>
4  </property>
5  <property>
6      <name>dfs.namenode.name.dir</name>
7      <value>file:/home/ubuntu/hadoop-2.7.4/hdfs/namenode</value>
8  </property>
   (continued to next page)
```

# Configuration (5)

- **hdfs-site.xml**
  - You can set internal HDFS information for namenode and datanode
  - Write properties between `<configuration>` and `</configuration>`

| $HADOOP_HOME/etc/hadoop/hdfs-site.xml (continued) |
|---|
| 9  `<property>` |
| 10  `    <name>dfs.datanode.data.dir</name>` |
| 11  `    <value>file:/home/ubuntu/hadoop-2.7.4/hdfs/datanode</value>` |
| 12  `</property>` |

# Configuration (6)

- **Formatting the HDFS file system**

```
1  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs namenode -format
```

- **The output will look like this:**

```
1    ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs namenode -format
2    INFO namenode.NameNode: STARTUP_MSG:
(...)                              (...)
3    INFO Storage directory (...)/hdfs/namenode has been successfully formatted.
(...)                              (...)
4    /*********************************************************
5    SHUTDOWN_MSG: Shutting down NameNode at 'username'/127.0.0.1
6    *********************************************************/
```

# Starting / Stopping cluster (1)

■ **Run the command:**

```
1   ubuntu@ip-x-x-x:~/hadoop-2.7.4$ sbin/start-dfs.sh
```

# Starting / Stopping cluster (2)

■ **The output will look like this:**

```
1   ubuntu@ip-x-x-x:~/hadoop-2.7.4$ sbin/start-dfs.sh

2   Starting namenodes on [localhost]

3   localhost: starting namenode, logging to 'Namenode logging directory'

4   localhost: starting datanode, logging to 'Datanode logging directory'

5   Starting secondary namenodes [0.0.0.0]

6   The authenticity of host '0.0.0.0 (0.0.0.0)' can't be established.

7   ECDSA key fingerprint is 16:20:01:83:ef:85:41:fb:ad:90:19:20:59:e1:7e:65.

8   Are you sure you want to continue connecting (yes/no)? yes          Typing!!!

9   0.0.0.0: starting secondarynamenode, logging to 'logging directory'

10  ubuntu@ip-x-x-x:~/hadoop-2.7.4$
```

# Starting / Stopping cluster (3)

- **Checking whether the Hadoop processes are running**

```
1  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ jps

2  1001 NameNode

3  1002 DataNode

4  1003 SecondaryNameNode

5  1004 Jps

7  ubuntu@ip-x-x-x:~/hadoop-2.7.4$
```

# Starting / Stopping cluster (4)

- **Run the command:**

```
1  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ sbin/stop-dfs.sh
```

- **The output will look like this:**

```
1  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ sbin/stop-dfs.sh
2  Stopping namenodes on [localhost]
3  localhost: stopping namenode
4  localhost: stopping datanode
5  Stopping secondary namenodes on [0.0.0.0]
6  0.0.0.0: stopping secondarynamenode
7  ubuntu@ip-x-x-x:~/hadoop-2.7.4$
```

# Hadoop Web Interface (1)

- **Web UI of the NameNode daemon**
  - http://localhost:50070

```
1  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ firefox
```

# Hadoop Web Interface (2)

# Running Spark with HDFS (1)

■ **WordCount example**

   ▪ Save three ebooks from Project Gutenberg

```
1   ubuntu@ip-x-x-x:~/hadoop-2.7.4$ wget
    http://www.gutenberg.org/cache/epub/20417/pg20417.txt

2   ubuntu@ip-x-x-x:~/hadoop-2.7.4$ wget http://www.gutenberg.org/files/5000/5000-8.txt

3   ubuntu@ip-x-x-x:~/hadoop-2.7.4$ wget http://www.gutenberg.org/files/4300/4300-0.txt
```

# Running Spark with HDFS (2)

■ **WordCount example**

▪ Put .txt files into HDFS

```
1  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –mkdir /input

2  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –put pg20417.txt /input/sample1.txt

3  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –put 5000-8.txt /input/sample2.txt

4  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –put 4300-0.txt /input/sample3.txt

5  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –ls /input

6  drwxr-xr-x - ubuntu supergroup    674570  2017-10-19 14:00 /input/sample1.txt

7  drwxr-xr-x - ubuntu supergroup   1428841  2017-10-19 14:00 /input/sample2.txt

8  drwxr-xr-x - ubuntu supergroup   1580890  2017-10-19 14:00 /input/sample3.txt

9  ubuntu@ip-x-x-x:~/hadoop-2.7.4$
```

# Running Spark with HDFS (3)

■ **WordCount example**

  ▪ Check the text files in HDFS

# Running Spark with HDFS (4)

## Browse Directory

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| / | | | | | | | | Go! |

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| drwxr-xr-x | ubuntu | supergroup | 0 B | 9/26/2017, 2:58:08 AM | 0 | 0 B | input |

## Browse Directory

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| /input | | | | | | | | Go! |

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | ubuntu | supergroup | 658.76 KB | 9/26/2017, 2:57:47 AM | 1 | 128 MB | sample1.txt |
| -rw-r--r-- | ubuntu | supergroup | 1.36 MB | 9/26/2017, 2:57:58 AM | 1 | 128 MB | sample2.txt |
| -rw-r--r-- | ubuntu | supergroup | 1.51 MB | 9/26/2017, 2:58:08 AM | 1 | 128 MB | sample3.txt |

# Running Spark with HDFS (5)

- **Start and open your Spark shell**

```
1   ubuntu@ip-x-x-x:~/hadoop-2.7.4$ cd $SPARK_HOME

2   ubuntu@ip-x-x-x:~/spark-2.1.0$ sbin/start-all.sh

3   ubuntu@ip-x-x-x:~/spark-2.1.0$ bin/pyspark

4   Python 2.7.6 (default, Oct 26 2016, 20:30:19)

5   [GCC 4.8.4] on linux2

6   Type "help", "copyright", "credits" or "license" for more information.

7   [...snipp...]

8   Using Python version 2.7.6 (default, Oct 26 2016 20:30:19)

9   SparkSession available as 'spark'.

10  >>>
```

# Running Spark with HDFS (6)

■ **Example 1: Count the total words in text files**

```
1   >>> text_file = sc.textFile("hdfs://localhost:9000/input")

2   >>> counts = text_file.flatMap(lambda line: line.split(" ")) \

3   ...                    .count()

4   [Stage 0:>>>>>>>>>>>>>                                    (1 + 1) / 3]

5   >>>
```

# Running Spark with HDFS (7)

■ **Example 1: Check the output**

```
1  >>> print counts

2  664559

3  >>>
```

# Running Spark with HDFS (8)

■ **Example 2: Count the occurrence of each word in text files**

```
1   >>> text_file = sc.textFile("hdfs://localhost:9000/input")

2   >>> wordcounts = text_file.flatMap(lambda line: line.split(" ")) \

3   ...                  .map(lambda word: (word, 1)) \

4   ...                  .reduceByKey(lambda a, b: a + b)

5   >>> wordcounts.saveAsTextFile("hdfs://localhost:9000/output")

6   [Stage 0:>>>>>>>>>>>>>>                                    (1 + 1) / 3]

7   [Stage 1:>>>>>>>>>>>>>>                                    (1 + 1) / 3]

8   >>> exit()
```

# Running Spark with HDFS (9)

■ **Example 2: Check the output**

```
1  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –cat /output /part-00000

2  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –getmerge /output result.txt

3  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ vi result.txt
```

# Running Spark with HDFS (10)

- **Example 3: Count the occurrence of each word in text files and sort words by frequency**

```
1  >>> text_file = sc.textFile("hdfs://localhost:9000/input")

2  >>> sortcounts = text_file.flatMap(lambda line: line.split(" ")) \

3  ...                 .map(lambda word: (word, 1)) \

4  ...                 .reduceByKey(lambda a, b: a + b) \

5  ...                 .sortBy(lambda x: -x[1])

5  >>> sortcounts.saveAsTextFile("hdfs://localhost:9000/output2")

6  [Stage 0:>>>>>>>>>>>>>>                                    (1 + 1) / 3]

7  [Stage 1:>>>>>>>>>>>>>>                                    (1 + 1) / 3]

8  >>> exit()
```

# Running Spark with HDFS (11)

■ **Example 3: Check the output**

```
1  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –cat /output2 /part-00000

2  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –getmerge /output2 result2.txt

3  ubuntu@ip-x-x-x:~/hadoop-2.7.4$ vi result2.txt
```

# Exercise 1

■ **Find the top 5 most used words only**

■ **Hints**

- ▪ takeOrdered(*N, function*)
  - ▪ Description: get the *N* elements from an RDD ordered in ascending order or specified by the optional *function*

■ **Please show me your result like this!**

```
1   [(u'the', 42098), (u'', 34667), (u'of', 23947), (u'and', 16921), (u'a', 12060)]
```

# Exercise 2

■ **Make a bigram count program using pyspark**

    ▪ Example

apple  banana  banana  apple  banana  banana

⬇

((apple, banana), 1) , ((banana, banana), 1) , ((banana, apple), 1) ,
((apple, banana), 1) , ((banana, banana), 1)

⬇

((apple, banana), 2) , ((banana, banana), 2) , ((banana, apple), 1)

# Exercise 2

- **Fill in the blank (???)**

```
1  >>> text_file = sc.textFile("hdfs://localhost:9000/input")

2  >>> bicounts = text_file.map(lambda line: line.split(" ")) \

3  ...                .flatMap(lambda x: ???) \

4  ...                .reduceByKey(lambda a, b: a + b)

5  >>> bicounts.saveAsTextFile("hdfs://localhost:9000/output3")

6  [Stage 0:>>>>>>>>>>>>>>                              (1 + 1) / 3]

7  [Stage 1:>>>>>>>>>>>>>>                              (1 + 1) / 3]

8  >>> exit()
```

# Exercise 2

■ **Check the output**

```
1   ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –cat /output3

2   ubuntu@ip-x-x-x:~/hadoop-2.7.4$ bin/hdfs dfs –getmerge /output3 result3.txt

3   ubuntu@ip-x-x-x:~/hadoop-2.7.4$ vi result3.txt
```

# Appendix