# Lab 01: A Gentle Introduction to Hadoop

## How much work, in percent (%),have you finished in each section?
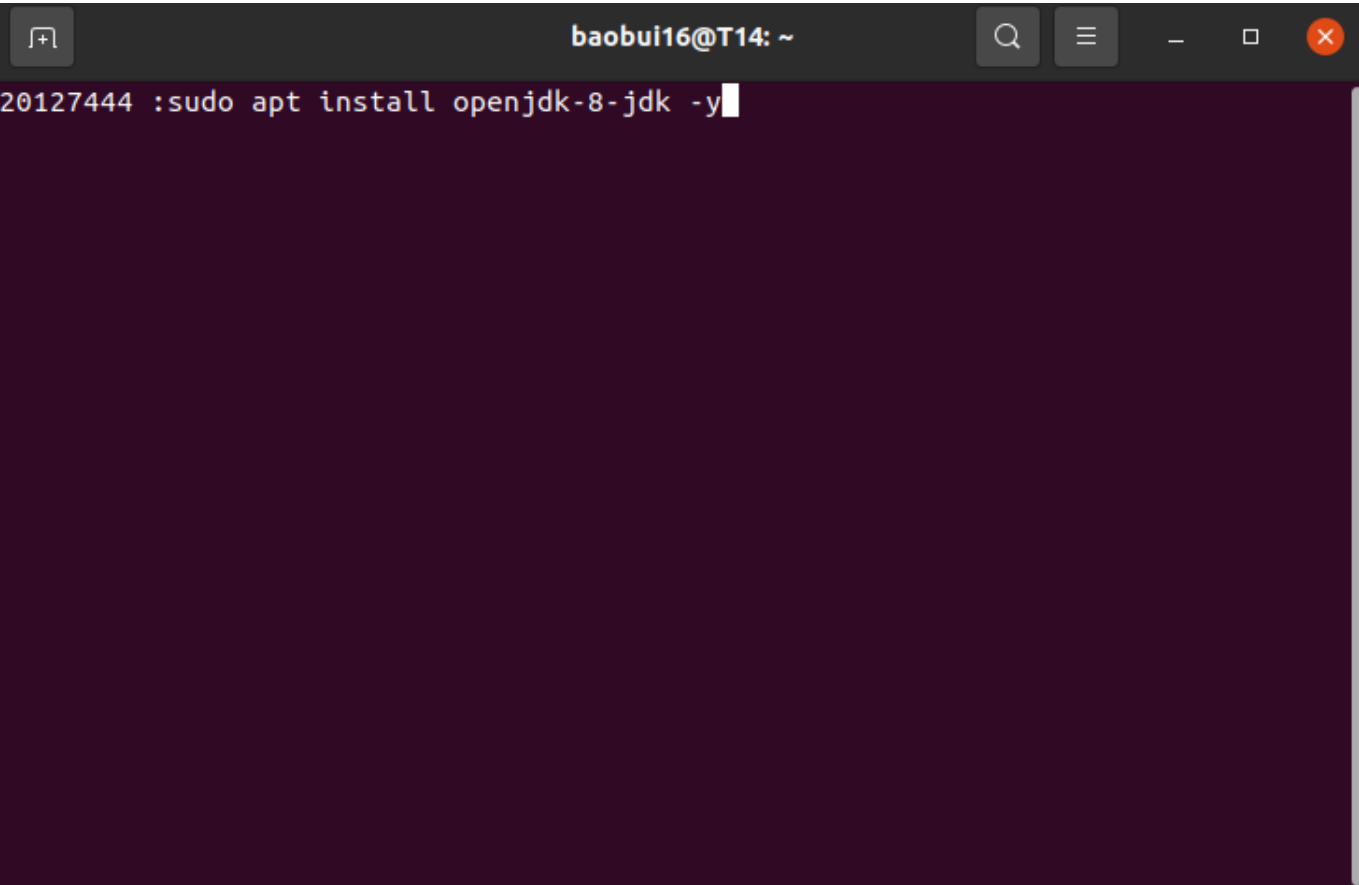
**My team works in section**

```
| Section 1 |  Section 2 |  Section 3  |  Section 4 |
|-----------|------------|-------------|------------|
|    100%   |    100%    |    100%     |     25%    |


(In section 4 we just done section 4.1 and we can't do section 4.2)
```

## Setting up Single-node Hadoop Cluster

- Install java



- Check java settings and path

- Install openSSH



- Create and Install SSH Certificates

```
20127444 :ssh-keygen -t rsa -P '' -f ~/.ssh/id_rsa
Generating public/private rsa key pair.
/home/baobui16/.ssh/id_rsa already exists.
Overwrite (y/n)? y
Your identification has been saved in /home/baobui16/.ssh/id_rsa
Your public key has been saved in /home/baobui16/.ssh/id_rsa.pub
The key fingerprint is:
SHA256:p4EOenn7ngGlIa0SLDOqNteNKkupf8jTvu1/hQz9E4c baobui16@T14
The key's randomart image is:
+---[RSA 3072]----+
|                 |
|  .   .          |
|+ o . o ..    .  |
|.+ . o =. . E .  |
|. . o + So.o o   |
|. .o.+o. +o +    |
|.*.+ooo.o  . .   |
|+.*.+o . o.      |
|.o+=oo++=.       |
+----[SHA256]-----+
20127444 :
```
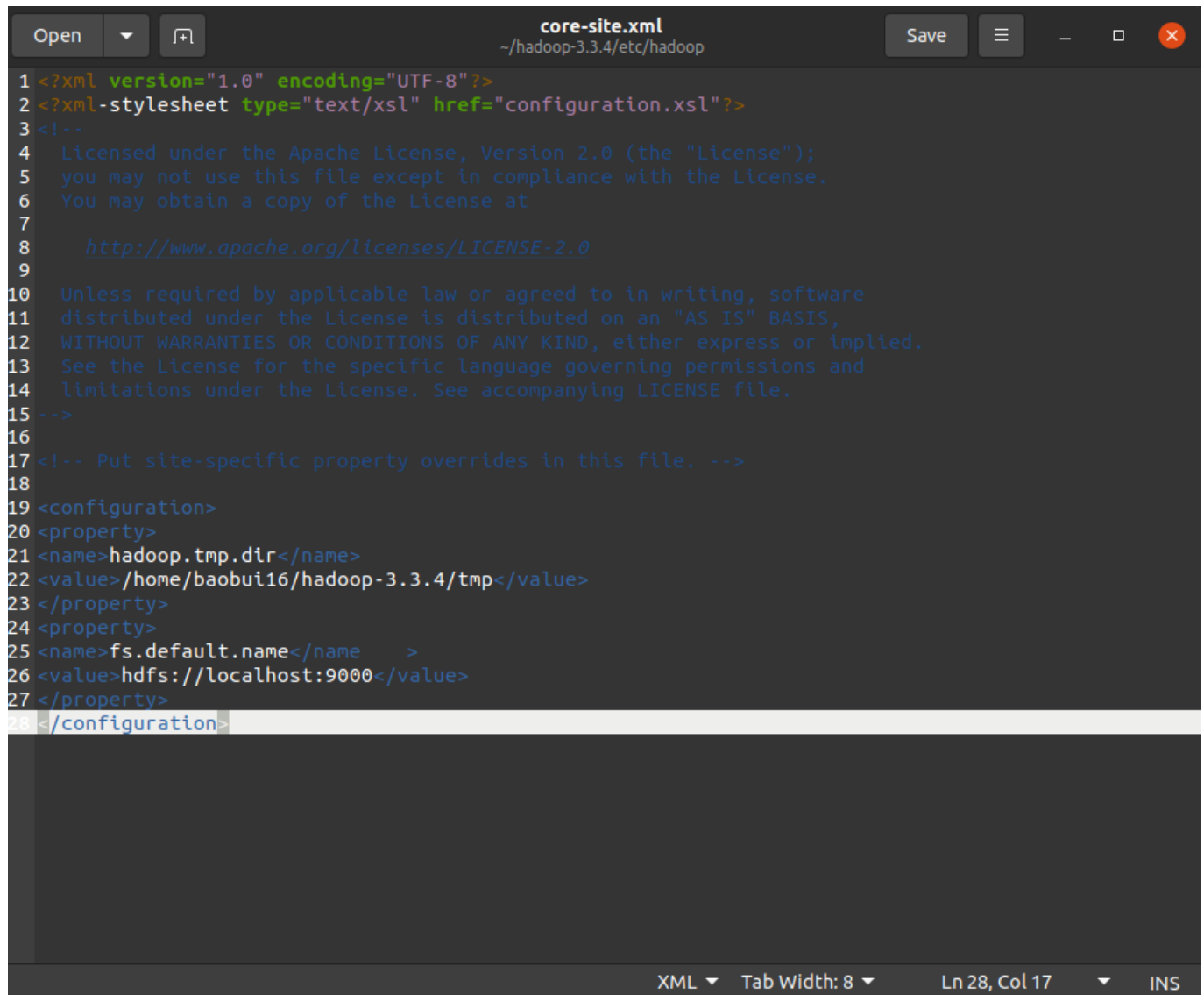
```
20127444 :cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
20127444 :ssh localhost
Welcome to Ubuntu 20.04.3 LTS (GNU/Linux 5.14.0-1057-oem x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:     https://landscape.canonical.com
 * Support:        https://ubuntu.com/advantage

2 devices have a firmware upgrade available.
Run `fwupdmgr get-upgrades` for more information.


416 updates can be applied immediately.
307 of these updates are standard security updates.
To see these additional updates run: apt list --upgradable

Last login: Wed Mar  1 08:20:58 2023 from 127.0.0.1
(base) baobui16@T14:~$
```
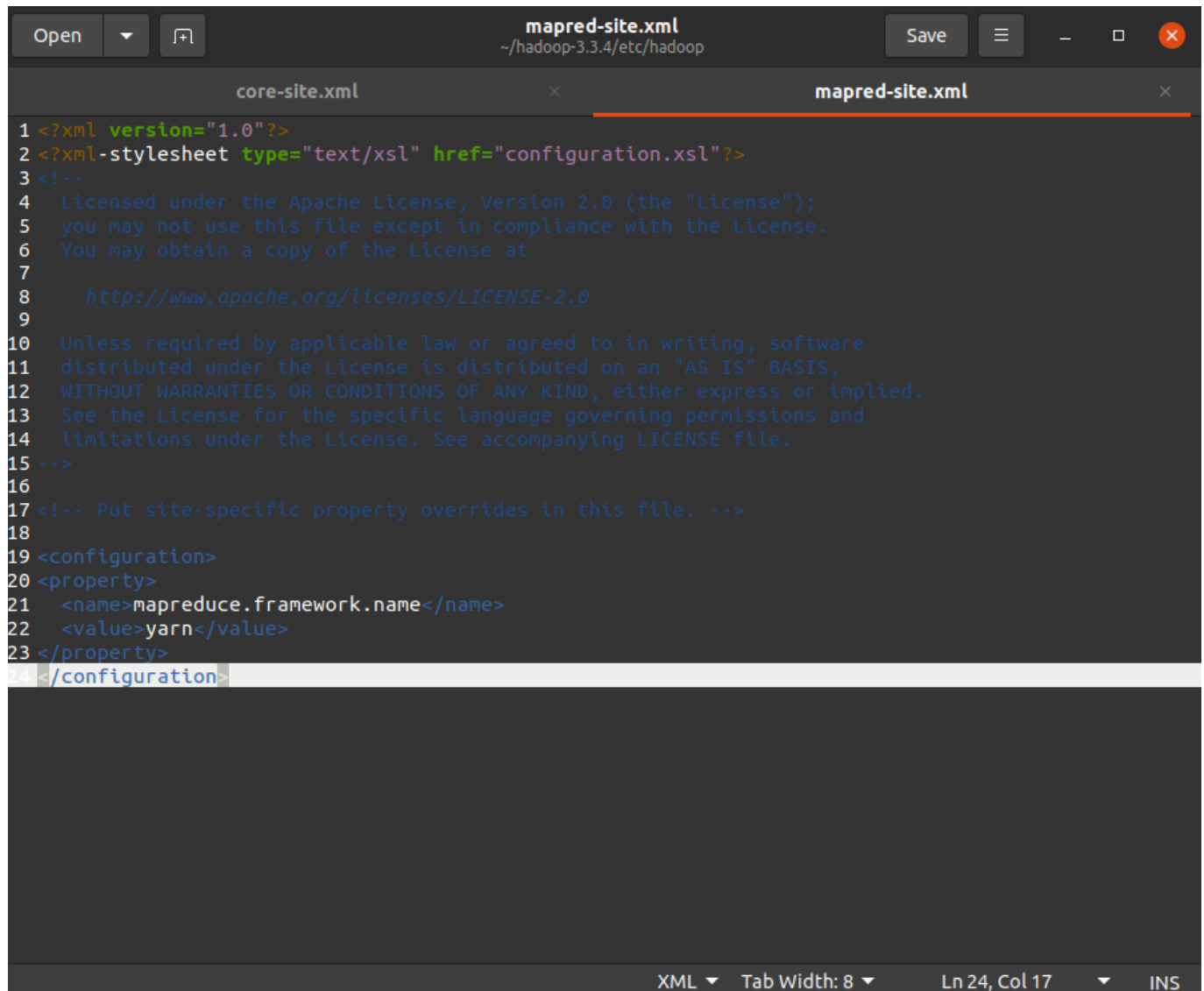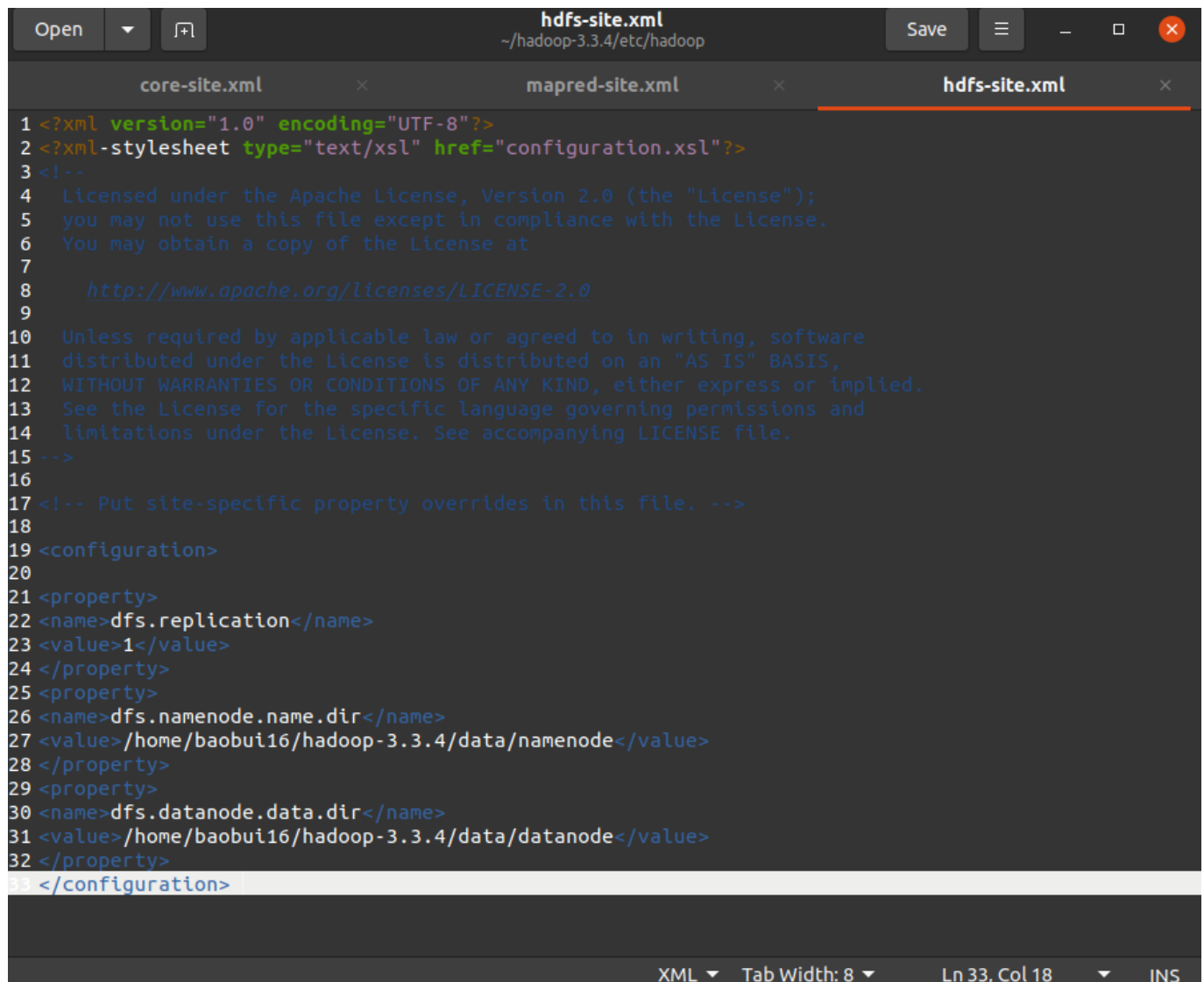
- Configure file core-site.xml

- Configure the file mapred-site.xml

- Configure file hdfs-site.xml

- Configure file yarn-site.xml

- Configure file hadoop-env.sh

```
Open    ▼    ⊞                       hadoop-env.sh                      Save    ☰    _   □   ✕
                                   ~/hadoop-3.3.4/etc/hadoop
15 # See the License for the specific language governing permissions and
16 # limitations under the License.
17
18 # Set Hadoop-specific environment variables here.
19
20 ##
21 ## THIS FILE ACTS AS THE MASTER FILE FOR ALL HADOOP PROJECTS.
22 ## SETTINGS HERE WILL BE READ BY ALL HADOOP COMMANDS.  THEREFORE,
23 ## ONE CAN USE THIS FILE TO SET YARN, HDFS, AND MAPREDUCE
24 ## CONFIGURATION OPTIONS INSTEAD OF xxx-env.sh.
25 ##
26 ## Precedence rules:
27 ##
28 ## {yarn-env.sh|hdfs-env.sh} > hadoop-env.sh > hard-coded defaults
29 ##
30 ## {YARN_xyz|HDFS_xyz} > HADOOP_xyz > hard-coded defaults
31 ##
32
33 # Many of the options here are built from the perspective that users
34 # may want to provide OVERWRITING values on the command line.
35 # For example:
36 #
37 #   JAVA_HOME=/usr/java/testing hdfs dfs -ls
38 export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
39 # Therefore, the vast majority (BUT NOT ALL!) of these defaults
40 # are configured for substitution and not append.  If append
41 # is preferable, modify this file accordingly.
42
43 ###
44 # Generic settings for HADOOP
45 ###
46
47 # Technically, the only required environment variable is JAVA_HOME.
48 # All others are optional.  However, the defaults are probably not
49 # preferred.  Many sites configure these options outside of Hadoop,
50 # such as in /etc/profile.d
51
52 # The java implementation to use. By default, this environment
                                      sh ▼   Tab Width: 8 ▼        Ln 45, Col 4       ▼    INS
```

- Configure file bash

- Successful installation

## SET SUCCESS OF TEAM MEMBERS:

- 20127444 - Bùi Duy Bảo



- 20127049 - Nguyễn Đức Minh

- 20127092 - Nguyễn Minh Tuấn



- 20127448 - Nguyễn Thái Bảo

```
                              bao@baobao: ~
bao@baobao:~$ cat /etc/machine-id
8123fc675b2d48dba9979a6f5804f4cb
bao@baobao:~$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as bao in 10
WARNING: This is not a recommended production deployment configurati
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [baobao]
Starting resourcemanager
Starting nodemanagers
bao@baobao:~$ jps
3044 NodeManager
2612 SecondaryNameNode
2038 NameNode
3419 Jps
2284 DataNode
2879 ResourceManager
bao@baobao:~$ 
```

# Introduction to MapReduce

1. **How do the input keys-values, the intermediate keys-values, and the output keys-values relate?**

The input keys-values represent the input data that is read from HDFS or other data sources

The intermediate keys-values are generated by the map tasks during the processing of the input data, and are in a different format or schema than the input data. These intermediate pairs are then sorted and shuffled across the cluster, and sent to the reduce tasks for further processing.

The output keys-values are the final key-value pairs generated by the reduce tasks, and are the result of aggregating, summarizing, or transforming the intermediate pairs

Overall, the input keys-values, intermediate keys-values, and output keys-values are all important components of the MapReduce data processing model, and they are all related to each other in the sense that they represent different stages in the processing of data on a distributed cluster.

2. **How does MapReduce deal with node failures?**

MapReduce deals with node failure by being designed to be fault-tolerant and by using some techniques, for example data replication and job restart even node failure appears.

About by being designed to be fault-tolerant: it means that MapReduce can handle node failures and continue processing the job. If a node fails, the tasks running on that node are automatically reassigned to other available nodes. The framework also periodically pings the nodes to check if they are still alive. If a node does not respond, it is assumed to have failed, and its tasks are reassigned to other nodes.

About by using data replication: MapReduce uses this technique to ensure that data is not lost when node failure happens, each block of data is replicated across multiple nodes in the cluster. If one node fails, the

data is still available on the other nodes, and the job can continue processing.

About by using job restart: in this situation, the entire job may need to be restarted. This is because if a node fails while processing a task, the output of that task may be lost. If the output of a task is lost, any subsequent tasks that depend on that output will need to be rerun.

3. **What is the meaning and implication of locality? What does it use?**

The locality is input data is stored on the local disks of the machines that make up cluster. It's manager by GFS which divides each file into 64 MB blocks, and stores several copies of each block (typically 3 copies) on different machines,if it had any node fail ,GFS would be specify another node that contains the copy in order to comply with the minimum requirement, 3 copies must be obtained

It's used when running large MapReduce operations on a significant fraction of the workers in a cluster, most input data is read locally and consumes so it not need to spend network bandwidth

4. **Which problem is addressed by introducing a combiner function to the MapReduce model?**

Combiner Function is similar to small Reduce phases of each Mapper on local disk which helps the process decrease number of pairs (key-value) before Reduce phases. As we see, reducing pairs of key-value on local disk decreases the workload that enhance bandwidth quality to run the process faster.

## Running a warm-up problem: Word Count

- Check Hadoop version

```
                                   baobui16@T14: ~

(base) baobui16@T14:~$ cat /etc/machine-id
c8d75585b7cd40b0802a09676ad54d72
(base) baobui16@T14:~$ hadoop version
Hadoop 3.3.4
Source code repository https://github.com/apache/hadoop.git -r a585a73c3e02ac623
50c136643a5e7f6095a3dbb
Compiled by stevel on 2022-07-29T12:32Z
Compiled with protoc 3.7.1
From source with checksum fb9dd8918a7b8a5b430d61af858f6ec
This command was run using /home/baobui16/hadoop-3.3.4/share/hadoop/common/hadoo
p-common-3.3.4.jar
(base) baobui16@T14:~$ javac -version
javac 1.8.0_362
(base) baobui16@T14:~$
```

- The input directory contains the file input.txt

| Name | Size | Modified | Star |
|------|------|----------|------|
| Input | 1 item | 21:16 | ☆ |
| tutorial_class | 0 items | 21:17 | ☆ |
| WordCount.java | 2,1 kB | 19 Thg 2 | ☆ |

15 / 26

```
input.txt
~/Desktop/Lab1/Input

 1 Cristiano
 2 Ronaldo
 3 David
 4 Beckham
 5 Neymar
 6 Messi
 7 Neymar
 8 Ronaldo
 9 Messi
10 Beckham
```

Plain Text    Tab Width: 8    Ln 1, Col 1    INS

- Enter the following path to export the Hadoop classpath to bash
- Make sure it's now exported

```
(base) baobui16@T14:~$ export HADOOP_CLASSPATH=$(hadoop classpath)
(base) baobui16@T14:~$ echo $HADOOP_CLASSPATH
/home/baobui16/hadoop-3.3.4/etc/hadoop:/home/baobui16/hadoop-3.3.4/share/hadoop/
common/lib/*:/home/baobui16/hadoop-3.3.4/share/hadoop/common/*:/home/baobui16/ha
doop-3.3.4/share/hadoop/hdfs:/home/baobui16/hadoop-3.3.4/share/hadoop/hdfs/lib/*
:/home/baobui16/hadoop-3.3.4/share/hadoop/hdfs/*:/home/baobui16/hadoop-3.3.4/sha
re/hadoop/mapreduce/*:/home/baobui16/hadoop-3.3.4/share/hadoop/yarn:/home/baobui
16/hadoop-3.3.4/share/hadoop/yarn/lib/*:/home/baobui16/hadoop-3.3.4/share/hadoop
/yarn/*
(base) baobui16@T14:~$
```

- Create this directory on HDFS and put input.txt . file

## Browse Directory

/WordCount/Input

```
(base) baobui16@T14:~$ hadoop fs -mkdir /WordCount
(base) baobui16@T14:~$ hadoop fs -mkdir /WordCount/Input
(base) baobui16@T14:~$ hadoop fs -put '/home/baobui16/Desktop/Lab1/Input/input.txt' /WordCount/Input
(base) baobui16@T14:~$
```

Show 25 entries

| ☐ | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | -rw-r--r-- | baobui16 | supergroup | 75 B | Mar 07 22:02 | 1 | 128 MB | input.txt | 🗑 |

Showing 1 to 1 of 1 entries                                    Previous  1  Next

## Browse Directory

/

```
(base) baobui16@T14:~$ hadoop fs -mkdir /WordCount
(base) baobui16@T14:~$ hadoop fs -mkdir /WordCount/Input
(base) baobui16@T14:~$ hadoop fs -put '/home/baobui16/Desktop/Lab1/Input/input.txt' /WordCount/Input
(base) baobui16@T14:~$
```

Show 25 entries

| ☐ | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | drwxr-xr-x | baobui16 | supergroup | 0 B | Mar 07 22:02 | 0 | 0 B | WordCount | 🗑 |
| ☐ | drwxr-xr-x | baobui16 | supergroup | 0 B | Mar 04 22:10 | 0 | 0 B | WordCountTutorial | 🗑 |
| ☐ | drwx------ | baobui16 | supergroup | 0 B | Mar 04 22:10 | 0 | 0 B | tmp | 🗑 |

Showing 1 to 3 of 3 entries                                    Previous  1  Next

Hadoop, 2022.

## Browse Directory

/WordCount

```
(base) baobui16@T14:~$ hadoop fs -mkdir /WordCount
(base) baobui16@T14:~$ hadoop fs -mkdir /WordCount/Input
(base) baobui16@T14:~$ hadoop fs -put '/home/baobui16/Desktop/Lab1/Input/input.txt' /WordCount/Input
(base) baobui16@T14:~$
```

Show 25 entries

| ☐ | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | drwxr-xr-x | baobui16 | supergroup | 0 B | Mar 07 22:02 | 0 | 0 B | Input | 🗑 |

Showing 1 to 1 of 1 entries                                    Previous  1  Next

Hadoop, 2022.

- Run file jar on Hadoop

| Name | | Size | Modified | Star |
|---|---|---|---|---|
| 📁 Input | | 1 item | 21:16 | ☆ |
| 📁 tutorial_class | | 3 items | 22:18 | ☆ |
| 📦 firstTutorial.jar | | 3,1 kB | 22:17 | ☆ |
| ☕ WordCount.java | | 2,1 kB | 19 Thg 2 | ☆ |

```
baobui16@T14: ~/Desktop/Lab1                    🔍  ≡  —  ▢  ✕

(base) baobui16@T14:~/Desktop/Lab1/tutorial_class$ jar -cvf firstTutorial.jar -C
 tutorial_class/ .
tutorial_class/. : no such file or directory
added manifest
(base) baobui16@T14:~/Desktop/Lab1/tutorial_class$ cd ..
(base) baobui16@T14:~/Desktop/Lab1$ jar -cvf firstTutorial.jar -C tutorial_class
/ .
added manifest
adding: WordCount$IntSumReducer.class(in = 1739) (out= 739)(deflated 57%)
adding: WordCount$TokenizerMapper.class(in = 1736) (out= 754)(deflated 56%)
adding: WordCount.class(in = 1491) (out= 814)(deflated 45%)
(base) baobui16@T14:~/Desktop/Lab1$ █
```

| Name | | Size | Modified | Star |
|---|---|---|---|---|
| ☕ WordCount.class | | 1,5 kB | 22:06 | ☆ |
| ☕ WordCount$IntSumReducer.class | | 1,7 kB | 22:06 | ☆ |
| ☕ WordCount$TokenizerMapper.class | | 1,7 kB | 22:06 | ☆ |

```
baobui16@T14: ~                    🔍  ≡  —  ▢  ✕

(base) baobui16@T14:~$ hadoop fs -mkdir /WordCount
(base) baobui16@T14:~$ hadoop fs -mkdir /WordCount/Input
(base) baobui16@T14:~$ hadoop fs -put '/home/baobui16/Desktop/Lab1/Input/input.txt'
/WordCount/Input
(base) baobui16@T14:~$ javac -classpath $HADOOP_CLASSPATH -d '/home/baobui16/Desktop
/Lab1/tutorial_class' '/home/baobui16/Desktop/Lab1/WordCount.java'
(base) baobui16@T14:~$ █
```

```
2023-03-07 22:21:47,654 INFO mapreduce.Job: Job job_1678198135509_0001 running i
n uber mode : false
2023-03-07 22:21:47,655 INFO mapreduce.Job:  map 0% reduce 0%
2023-03-07 22:21:52,741 INFO mapreduce.Job:  map 100% reduce 0%
2023-03-07 22:21:57,781 INFO mapreduce.Job:  map 100% reduce 100%
2023-03-07 22:21:57,797 INFO mapreduce.Job: Job job_1678198135509_0001 completed
 successfully
2023-03-07 22:21:57,913 INFO mapreduce.Job: Counters: 54
        File System Counters
                FILE: Number of bytes read=87
                FILE: Number of bytes written=551077
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=187
                HDFS: Number of bytes written=57
                HDFS: Number of read operations=8
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
                HDFS: Number of bytes read erasure-coded=0
        Job Counters
                Launched map tasks=1
                Launched reduce tasks=1
                Data-local map tasks=1
```

- Output

```
                     Peak Reduce Virtual memory (bytes)=2570833920
        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=75
        File Output Format Counters
                Bytes Written=57
(base) baobui16@T14:~/Desktop/Lab1$ hadoop dfs -cat /WordCount/Output/*
WARNING: Use of this script to execute dfs is deprecated.
WARNING: Attempting to execute replacement "hdfs dfs" instead.

Beckham 2
Cristiano       1
David   1
Messi   2
Neymar  2
Ronaldo 2
(base) baobui16@T14:~/Desktop/Lab1$
```

| Hadoop | Overview | Datanodes | Datanode Volume Failures | Snapshot | Startup Progress | Utilities ▾ |

# Browse Directory

| /WordCount | | | | | | Go! |

Show [ 25 ▾ ] entries                                                                 Search: [              ]

| ☐ | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|---|
| ☐ | drwxr-xr-x | baobui16 | supergroup | 0 B | Mar 07 22:02 | 0 | 0 B | Input 🗑 |
| ☐ | drwxr-xr-x | baobui16 | supergroup | 0 B | Mar 07 22:21 | 0 | 0 B | Output 🗑 |

Showing 1 to 2 of 2 entries                                              Previous [ 1 ] Next

Hadoop, 2022.

# Bonus

4.1 Extended Word Count: Unhealthy relationships

- Input

```
Open    ▼   ⊞          input.txt          Save   ☰   —  ▢  ✕
                    ~/Desktop/Bonus/Input
1 A  B
2 A  C
3 B  D
4 C  B
5 D  E
6 E  C                           22 / 26
7 C  F
8 G  F
9
```

```
Plain Text ▼   Tab Width: 8 ▼        Ln 9, Col 1      ▼    INS
```

- Output

```
                    Peak Map Physical memory (bytes)=357171200
                    Peak Map Virtual memory (bytes)=2563768320
                    Peak Reduce Physical memory (bytes)=206716928
                    Peak Reduce Virtual memory (bytes)=2570543104
            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=33
            File Output Format Counters
                    Bytes Written=39
(base) baobui16@T14:~/Desktop/Bonus$ hdfs dfs -cat /src/output/*
A       pos
B       neg
C       eq
D       eq
E       eq
F       neg
G       pos
(base) baobui16@T14:~/Desktop/Bonus$
```

- Directory containing output on hadoop

## Browse Directory

| /src | | | | | | Go! |

Show [ 25 ] entries                                                     Search: [        ]

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|---|
| ☐ | -rw-r--r-- | baobui16 | supergroup | 33 B | Mar 10 22:02 | 1 | 128 MB | input.txt 🗑 |
| ☐ | drwxr-xr-x | baobui16 | supergroup | 0 B | Mar 10 22:07 | 0 | 0 B | output 🗑 |

Showing 1 to 2 of 2 entries                                    Previous [ 1 ] Next

Hadoop, 2022.

## Browse Directory

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | -rw-r--r-- | baobui16 | supergroup | 0 B | Mar 10 22:07 | 1 | 128 MB | _SUCCESS | 🗑 |
| ☐ | -rw-r--r-- | baobui16 | supergroup | 38 B | Mar 10 22:07 | 1 | 128 MB | part-r-00000 | 🗑 |

Showing 1 to 2 of 2 entries

Hadoop, 2022.

- Visualize example



Reflection of your team. (Does your journey to the deadline have any bugs? How have you overcome it? What have you learned after this process? If you cannot overcome the bugs, describe where the bottlenecks are in your work.)

- In process , we have some problem about setting up Hadoop in make enviremnt so we need to uninstall ubuntu and resetup after that
- when we run Wordcount and expand wordcount problems , we had error in output of map and input of reduce not correct so we search in google and see solusion by some people
- After this lab , we can operate MapReduce and custom MapReduce in high level.
- We try to install full hadoop distribution mode but can't.

## References

- Two Cloudera version of WordCount problem:

  - https://docs.cloudera.com/documentation/other/tutorial/CDH5/topics/ht_wordcount2.html
  - https://docs.cloudera.com/documentation/other/tutorial/CDH5/topics/ht_wordcount3.html

- Apache Hadoop

  - [1] Apache Hadoop 3.3.4 – Hadoop Cluster Setup.
    https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/ClusterSetup.html.

  - [2] Apache Hadoop 3.3.4 – Hadoop in Secure Mode.
    https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/SecureMode.html.

  - [3] Apache Hadoop 3.3.4 – Hadoop: Setting up a Single Node Cluster.
    https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/SingleCluster.html.

  - [4] Apache Hadoop 3.3.4 – MapReduce Tutorial
    https://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html#Example:_WordCount_v1.0.

- Book:

  - MapReduce Design Patterns [Donald Miner, Adam Shook, 2012]
  - [5] Jeffrey Dean and Sanjay Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. In OSDI'04: Sixth Symposium on Operating System Design and Implementation, pages 137–150, San Francisco, CA, 2004.

- All of StackOverflow link related.

  - https://stackoverflow.com/questions/56153007/java-lang-exception-java-io-ioexception-wrong-value-class-while-setting-hadoop?fbclid=IwAR1UNjMqPVvXy9nMberqIHAYFNZ02q0we1BAoqZOGUhmEFArVmbnBfyoCiQ

  - https://stackoverflow.com/questions/58272650/what-exactly-does-data-locality-mean-in-hadoop?fbclid=IwAR0jSqFAXazrxBUA7uO8sZsGWbhkry6SYJAf4PtBEbNfilKoZcl6HTBbdLk

- Another link

- https://community.cloudera.com/t5/Support-Questions/class-org-apache-hadoop-io-MapWritable-is-not-class-org/m-p/136602?fbclid=IwAR3VazsXyHo688u4UyNvsBdj_a_2hOwuJ-KpGZD6fg4aMwW8trQlpz-KArA
- https://www.thoughtworks.com/insights/decoder/d/data-locality?fbclid=IwAR1BwoHh0Q8-hylq51ylms9NIWQI300GMKdZm4yydbQkEDBWhrbcqQoM554
- https://www.tutorialspoint.com/map_reduce/map_reduce_quick_guide.html
- https://www.youtube.com/watch?v=6sK3LDY7Pp4&t=483s
- https://www.youtube.com/watch?v=MZ_FUEnbrR4