# Principles of HDFS & MapReduce2:

Name: อวยชัย    ภิรมย์รื่น

Tel.   : 086-813-5354
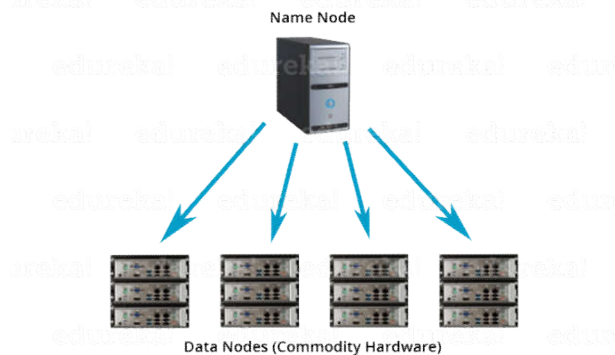
e-mail : p.Auoychai@gmail.com

# Wrap up

# # **HDFS Architecture :**

## **- What is HDFS ,** Java based distributed file system and allow to store large

data across multiple nodes in a Hadoop Cluster.



**Distributed Storage:**

**Can't use Linux command to see it**



*MasterNode:*

*NameNode , ResourceManager*

*SlaveNode:*

*DataNode , NodeManager*

# # HDFS Write-Read Process :

What the HDFS difference than traditional File System

# HDFS Write Process :

1. Set up of Pipeline

2. Data streaming and replication

3. Shutdown of Pipeline (Acknowledgement stage)

# #  HDFS Write Process :

## 1). Set up of Pipeline



Setting up HDFS - Write Pipeline

https://www.edureka.co

# #  HDFS Write Process :

## 2). Data streaming and replication



HDFS - Write Pipeline

https://www.edureka.co

# #  HDFS Write Process :

## 3). Shutdown of Pipeline



Acknowledgement in HDFS - Write

https://www.edureka.co

# # HDFS Read Process :



## HDFS - Read Architecture

BLK A

BLK B

**HDFS Client**

Client JVM

**Client Node**

① Read Request - Block A & B →

② ← IP Addresses: DN1 & DN3

**NameNode**

Read From DN1 & DN3

**Core Switch**

④

**Switch**

BLK A

DataNode 1  BLK A

3B

DataNode 2

3A

BLK B

DataNode 3  BLK B

Rack 1

**Switch**

DataNode 4  BLK A

DataNode 5

DataNode 6  BLK A

Rack 5

**Switch**

DataNode 7  BLK B

DataNode 8

DataNode 9  BLK B

Rack 7

# **Hadoop Command** :

- File System

# # **Interact to HDFS** :

- Interacting with HDFS is primarily performed from the command line using the script named **hdfs**. The **hdfs** script has the following usage:

```
$ hadoop fs [-option <arg>]
```

- The **-option** argument is the name of a specific option for the specified command, and **<arg>** is one or more arguments that that are specified for this option.

- For example, show help

```
$ hadoop fs -help
```

# Interact to HDFS :

- List directory contents

  - use -ls command:     **hdfs dfs -ls**

    ```
    $ hadoop fs -ls
    ```

    - Running the -ls command on a new cluster will not return any results. This is because the -ls command, without any arguments, will attempt to display the contents of the user's home directory on HDFS.

  - Providing -ls with the forward slash (/) as an argument displays the contents of the root of HDFS:

    ```
    $ hadoop fs -ls /
    ```
    **hdfs dfs -ls**

# Interact to HDFS :

- Creating a directory

  - To create the books directory within HDFS, use the -mkdir command: **hdfs dfs -mkdir**

    ```
    $ hadoop fs -mkdir [directory name]
    ```

  - For example, create books directory in home directory

    ```
    $ hadoop fs -mkdir books
    ```

  - Use the -ls command to verify that the previous directories were created:

    ```
    $ hadoop fs -ls
    ```

# Interact to HDFS :

- Copy Data onto HDFS

  - After a directory has been created for the current user, data can be uploaded to the user's HDFS home directory with the **-put** command:

    ```
    $ hadoop fs -put [source file] [destination file]
    ```

  - For example, copy book file from local to HDFS

    ```
    $ hadoop fs -put pg20417.txt books/pg20417.txt
    ```

  - Use the **-ls** command to verify that pg20417.txt was moved to HDFS:

    ```
    $ hadoop fs -ls books
    ```

# **Interact to HDFS** :

- Retrieve (view) Data from HDFS

  - Multiple commands allow data to be retrieved from HDFS.

  - To simply view the contents of a file, use the **-cat** command. **-cat** reads a file on HDFS and displays its contents to stdout.

  - The following command uses -cat to display the contents of pg20417.txt

    ```
    $ hadoop fs -cat books/pg20417.txt
    ```

# # **Interact to HDFS** :

- Retrieve (view) Data from HDFS

  - Data can also be copied from HDFS to the local filesystem using the **-get** command. The **-get** command is the opposite of the **-put** command:

    ```
    $ hadoop fs -get [source file] [destination file]
    ```

  - For example, This command copies pg20417.txt from HDFS to the local filesystem.

    ```
    $ hadoop fs -get pg20417.txt .
    ```

# **Wrap up HDFS** :



https://www.slideshare.net/eakasit_dpu/introduction-to-hadoop-and-mapreduce-75323926?qid=6dc387d1-f614-4646-a2cb-b4805805235a&v=&b=&from_search=3

# **#  What is MapReduce :**

Google released a paper on MapReduce technology in December, 2004. This became the genesis of the Hadoop Processing Model. So, MapReduce is a programming model that allows us to perform parallel and distributed processing on huge data sets.

# # MapReduce best fit for ? :
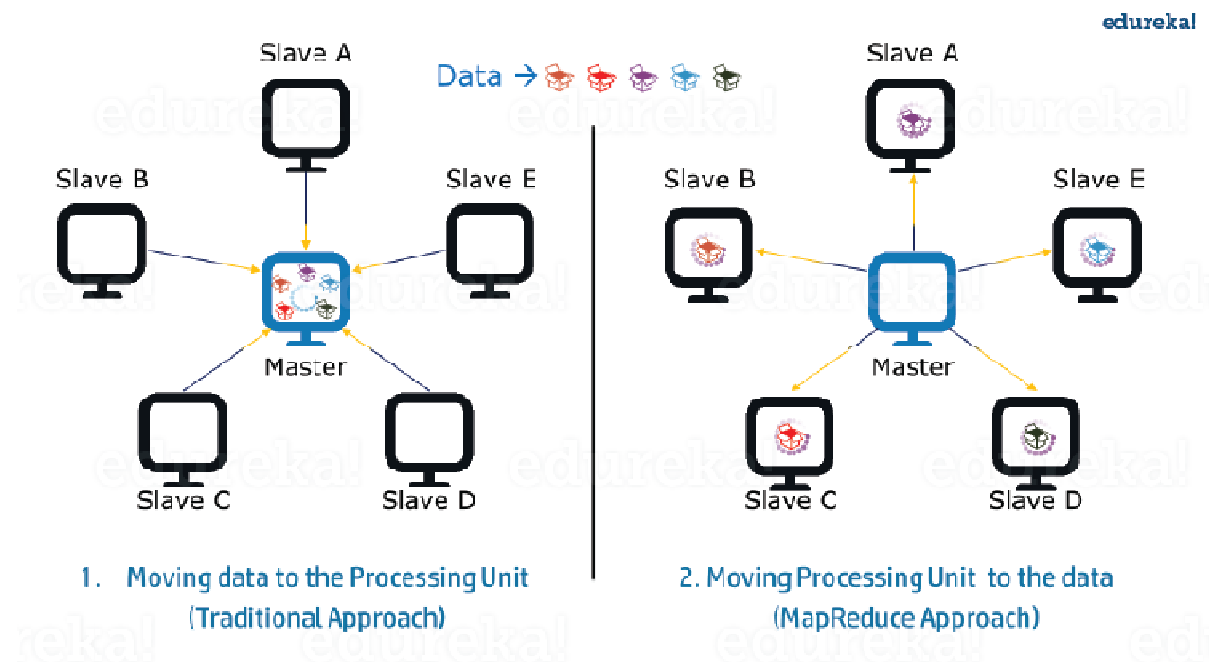
**Principle of MapReduce :**

-

# Data Processing Model : Traditional vs. Distributed

The biggest advantage of MapReduce :

1. Parallel Processing:

2. Data Locality:

# MapReduce Concept :

What is MapReduce?

# MapReduce Process Flow :



The Overall MapReduce Word Count Process

edureka!

| Input | Splitting | Mapping | Shuffling | Reducing | Final Result |
|---|---|---|---|---|---|

List(K2, V2)    K2, List(V2)

K1, V1    List(K3, V3)

Deer Bear River → Deer, 1 / Bear, 1 / River, 1 → Bear, (1,1) → Bear, 2

Deer Bear River / Car Car River / Deer Car Bear → Car Car River → Car, 1 / Car, 1 / River, 1 → Car, (1,1,1) → Car, 3 → Bear, 2 / Car, 3 / Deer, 2 / River, 2

Deer Car Bear → Deer, 1 / Car, 1 / Bear, 1 → Deer, (1,1) → Deer, 2

River, (1,1) → River, 2

# MapReduce Programming Model :

```java
public class WordCount {

  public static class TokenizerMapper
       extends Mapper<Object, Text, Text, IntWritable>{
    ...
  }

  public static class IntSumReducer
       extends Reducer<Text,IntWritable,Text,IntWritable> {
    ....
  }

  public static void main(String[] args) throws Exception {

    Configuration conf = new Configuration();

    Job job = Job.getInstance(conf, "word count");

    job.setJarByClass(WordCount.class);

    job.setMapperClass(TokenizerMapper.class);

    job.setCombinerClass(IntSumReducer.class);

    job.setReducerClass(IntSumReducer.class);

    job.setOutputKeyClass(Text.class);

    job.setOutputValueClass(IntWritable.class);

    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));

    System.exit(job.waitForCompletion(true) ? 0 : 1);
  }
}
```

# MapReduce Programming Model :

```java
public static class TokenizerMapper
     extends Mapper<Object, Text, Text, IntWritable>{

  private final static IntWritable one = new IntWritable(1);
  private Text word = new Text();

  public void map(Object key, Text value, Context context
                  ) throws IOException, InterruptedException {

    StringTokenizer itr = new StringTokenizer(value.toString());
    while (itr.hasMoreTokens()) {
      word.set(itr.nextToken());
      context.write(word, one);

    }
  }
}


public static class IntSumReducer
     extends Reducer<Text,IntWritable,Text,IntWritable> {
  private IntWritable result = new IntWritable();

  public void reduce(Text key, Iterable<IntWritable> values,
                     Context context
                     ) throws IOException, InterruptedException {

    int sum = 0;
    for (IntWritable val : values) {
      sum += val.get();
    }
    result.set(sum);
    context.write(key, result);
  }

}
```

# MapReduce Programming Model :

```java
import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

# MapReduce Programming Model :

```java
public static void main(String[] args) throws Exception {

    Configuration conf = new Configuration();

    Job job = Job.getInstance(conf, "word count");

    job.setJarByClass(WordCount.class);

    job.setMapperClass(TokenizerMapper.class);

    job.setCombinerClass(IntSumReducer.class);

    job.setReducerClass(IntSumReducer.class);

    job.setOutputKeyClass(Text.class);

    job.setOutputValueClass(IntWritable.class);

    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));

    System.exit(job.waitForCompletion(true) ? 0 : 1);
}
```

# MapReduce Programming :

1). Write Code Program

    1-1:  Create Directory /home/auoychai/units

    /home/auoychai/&lt;xxxx&gt;.java

2). Compile & Package

    $javac –classpath hadoop-core-1.2.1.jar -d units  &lt;xxxx&gt;.java

    $jar –cvf  &lt;xxxx&gt;.jar –C units/ .

3). Run Job-Map/Reduce

# Hand-On : File System

---

1). Create Directory ใน Hadoop

      /user/ [ lab  | dataset ]

      /lab/[ input  | output ]

      /dataset/[csv | txt | json | db

2). Copy file เข้า Hadoop

      /home/auoychai/dataset

      ** copy file ของแต่ละหมวดหมู่ไปไว้ใน Hadoop Directory /user/dataset/ [ … ]
          ตามกลุ่มของประเภทไฟล์

3).  List Directory ใน Hadoop

4).  ลบ File / Directory ใน Hadoop

      ** Delete file ของแต่ละหมวดหมู่ไปไว้ใน Hadoop Directory /user/dataset/ [ … ]

# Hand-On : File System : MapReduce Programming :

1). Write Code Program

    1-1:  Create Directory /home/auoychai/units

    /home/auoychai/Wordcount.java

2). Compile & Package

    $javac –classpath hadoop-core-1.2.1.jar -d units WordCount.java

    $jar –cvf  WordCount.jar –C units/ .

3). Run Job-Map/Reduce

    $hadoop jar WordCount.jar WordCount /user/lab/input/sample.txt
/user/lab/output/o3

# Principle HDFS & MapReduce :

# Wrap up