

Hadoop MapReduce

OVERVIEW

01 INTRODUCTION

02 HISTORY

03 FEATURES

04 CODE EXAMPLE

Structured

Traditional datasets that have rows, columns and clearly defined data attributes.

- DBMS tables
- Excel, CSV, SQL databases

Semi-structured

Does not conform to fixed fields, but does contain tags and markers.

- XML, HTML, JSON

Unstructured

Data that doesn't reside in a traditional row-column database.

- Images, Videos, PDF files, email attachments.

Big Data

Heterogeneous mix of structured , semi-structured and unstructured data.

Hadoop MapReduce

Hadoop is the most popular open source implementation of the MapReduce framework.

The framework is designed to handle large datasets (multi-terabytes) and processes them in parallel across multiple computing nodes.

Supported by many programming languages: Java, C++, Python

A cloud computing technology.

SIMPLICITY

SCALABILITY

FAULT-TOLERANCE

HISTORY



Origins

MapReduce and the Google File System was created by Google in 2003 to address the need of processing big data.

They also built a DBMS known as Big Table.



Open-Sourcing

Apache Hadoop was created as an open source implementation of the MapReduce system.

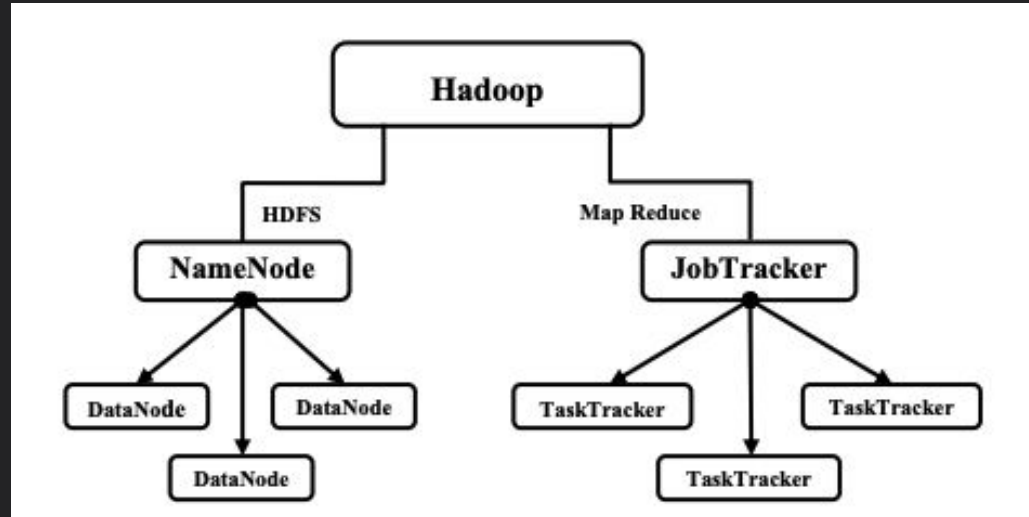
FEATURES

Hadoop Distributed File
System (HDFS)

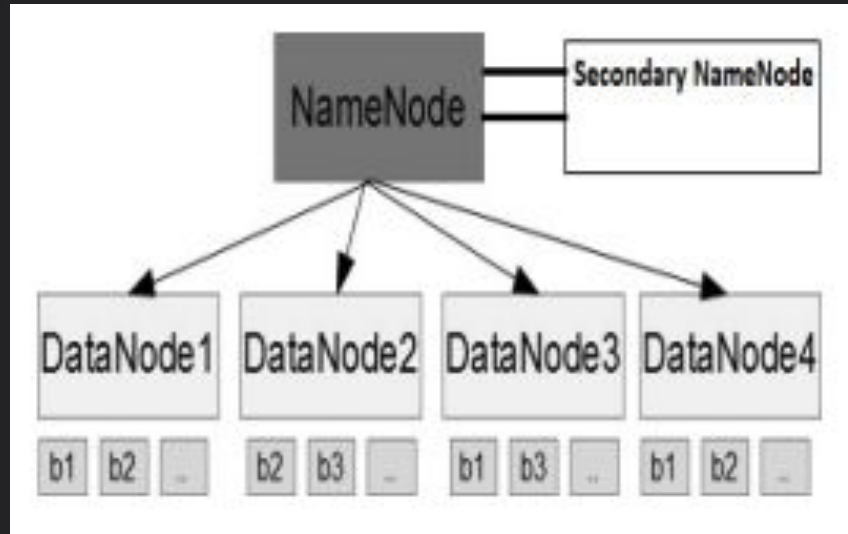
Data Storage

MapReduce Framework

Data Processing



HADOOP



NameNode

Runs on Manager node and dictates what the DataNode computers do with data.

It determines how files are split, where they are stored, and reviews the health of the file system.

DataNode

Controlled by the NameNode. They store data blocks and service read/write requests on files stored on the HDFS.

Secondary NameNode

Regularly reads the file system, logs the changes, and applies them onto a fsimage file.

MapReduce

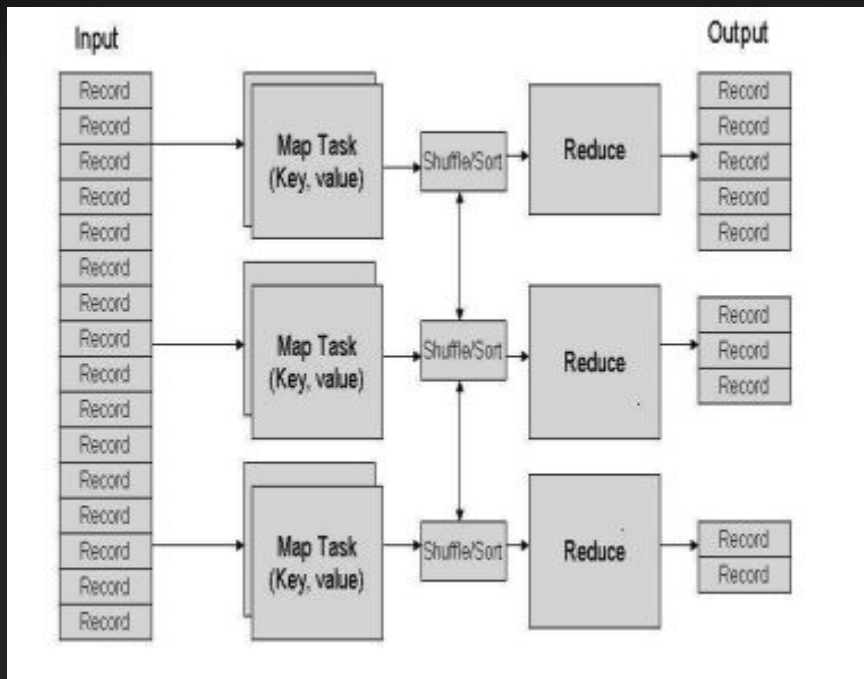
Map Phase

Splits input data into groups of key/value pairs using parallel processing.

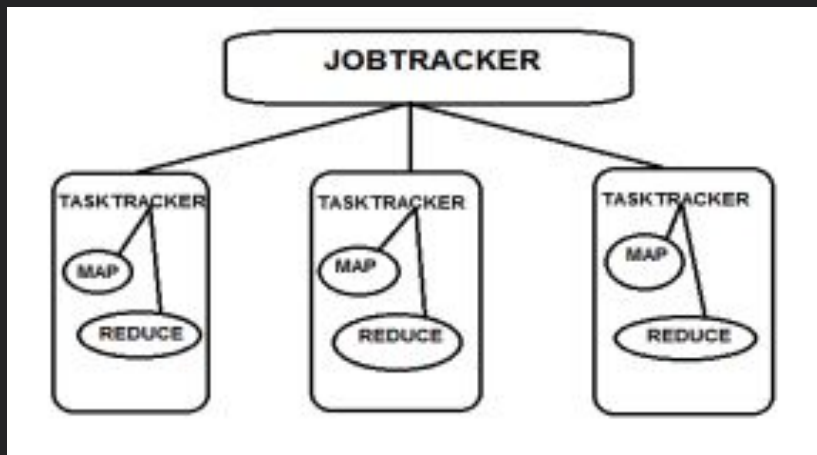
Groups the key/value pairs by their keys and feeds them to the reduce phase.

Reduce Phase

Aggregates and sorts the output from the Map phase to produce the final key/value pair outputs.



MapReduce



JobTracker

Runs on the manager node and manages job assignments, communicates with the NameNode, and monitors the health of the TaskTrackers that it oversees.

TaskTracker

Operates on the worker nodes and accepts and executes from distributed by the JobTracker. The JobTracker manages the task slots on each TaskTracker to ensure there is availability to run each MapReduce operation.

Sample Code

```
1 import java.io.IOException;
2 import java.util.StringTokenizer;
3
4 import org.apache.hadoop.conf.Configuration;
5 import org.apache.hadoop.fs.Path;
6 import org.apache.hadoop.io.IntWritable;
7 import org.apache.hadoop.io.Text;
8 import org.apache.hadoop.mapreduce.Job;
9 import org.apache.hadoop.mapreduce.Mapper;
10 import org.apache.hadoop.mapreduce.Reducer;
11 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
12 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
13
14 public class WordCount {
15
16     public static class TokenizerMapper
17         extends Mapper<Object, Text, Text, IntWritable>{
18
19         private final static IntWritable one = new IntWritable(1);
20         private Text word = new Text();
21
22         public void map(Object key, Text value, Context context
23             ) throws IOException, InterruptedException {
24             StringTokenizer itr = new StringTokenizer(value.toString());
25             while (itr.hasMoreTokens()) {
26                 word.set(itr.nextToken());
27                 context.write(word, one);
28             }
29         }
30     }
31
32     public static class IntSumReducer
33         extends Reducer<Text, IntWritable, Text, IntWritable> {
34         private IntWritable result = new IntWritable();
35
36         public void reduce(Text key, Iterable<IntWritable> values,
37             Context context
38             ) throws IOException, InterruptedException {
39             int sum = 0;
40             for (IntWritable val : values) {
41                 sum += val.get();
42             }
43             result.set(sum);
44             context.write(key, result);
45         }
46     }
47
48     public static void main(String[] args) throws Exception {
49         Configuration conf = new Configuration();
50         Job job = Job.getInstance(conf, "word count");
51         job.setJarByClass(WordCount.class);
52         job.setMapperClass(TokenizerMapper.class);
53         job.setCombinerClass(IntSumReducer.class);
54         job.setReducerClass(IntSumReducer.class);
55         job.setOutputKeyClass(Text.class);
56         job.setOutputValueClass(IntWritable.class);
57         FileInputFormat.addInputPath(job, new Path(args[0]));
58         FileOutputFormat.setOutputPath(job, new Path(args[1]));
59         System.exit(job.waitForCompletion(true) ? 0 : 1);
60     }
61 }
62
```

Mapper Class

Reducer Class

Rising Data Complexity

The demand to process vast amounts of data is increasing. These datasets hold valuable insights, driving the need for effective computational tools that can process large scale data.

Applications

The Hadoop MapReduce Framework allows users a solution for cloud-base large scale data processing and analytics.

Features

Hadoop MapReduce enables distributed processing of large datasets across clusters, offering simple, scalable and fault tolerant solutions for efficient data analysis and computation.

The Future

As big data continues to grow, the Hadoop MapReduce framework will need to continue to progress. New technologies, such as Apache Spark, have gained popularity due to their performance advantages.

CONCLUSION

REFERENCES

Dhavapriya, M., & Yasodha, N. (2016). Big data analytics: challenges and solutions using Hadoop, map reduce and big table. *International Journal of Computer Science Trends and Technology (IJCST)*, 4(1), 5-14.

Dittrich, J., & Quiané-Ruiz, J. A. (2012). Efficient big data processing in Hadoop MapReduce. *Proceedings of the VLDB Endowment*, 5(12), 2014-2015.

Ghazi, M. R., & Gangodkar, D. (2015). Hadoop, MapReduce and HDFS: a developers perspective. *Procedia Computer Science*, 48, 45-50.

Maitrey, S., & Jha, C. K. (2015). MapReduce: simplified data analysis of big data. *Procedia Computer Science*, 57, 563-571.

Mapreduce Tutorial. Apache Hadoop 3.4.0 – MapReduce Tutorial. (2024, March 4).

<https://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>

Pothuganti, A. (2015). Big data analytics: Hadoop-Map reduce & NoSQL databases. *International Journal of Computer Science and Information Technologies*, 6(1), 522-527.

Sagiroglu, S., & Sinanc, D. (2013, May). Big data: A review. In *2013 international conference on collaboration technologies and systems (CTS)* (pp. 42-47). IEEE.

Salo, J. (2012). Pattern definition of MapReduce. In *Proceedings of VikingPLoP 2012 Conference* (p. 97).