

Uma introdução ao processamento de dados com Hadoop MapReduce

Alan Cortes



Objetivo:

- Uma introdução

Objetivo:

- Uma introdução
- Uma visão simplificada

Objetivo:

- Uma introdução
- Uma visão simplificada
- Omissão de vários detalhes

Objetivo:

- Uma introdução
- Uma visão simplificada
- Omissão de vários detalhes
- Ao final ter uma idéia de como fazer BY
- Não tratarei de como instalar...

What Is Apache Hadoop?

- The Apache Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models.



Versões:

Version	Release Date
3.0.0-beta1	03 October, 2017
2.8.2	24 Oct, 2017
2.7.4	04 August, 2017
2.6.5	08 October, 2016

Stable: 2.8.2

Módulos do projeto:

- **Hadoop Common:** The common utilities that support the other Hadoop modules.

Módulos do projeto:

- **Hadoop Common:** The common utilities that support the other Hadoop modules.
- **Hadoop Distributed File System (HDFS™):** A distributed file system that provides high-throughput access to application data.

Módulos do projeto:

- **Hadoop Common:** The common utilities that support the other Hadoop modules.
- **Hadoop Distributed File System (HDFS™):** A distributed file system that provides high-throughput access to application data.
- **Hadoop YARN:** A framework for job scheduling and cluster resource management.

Módulos do projeto:

- **Hadoop Common:** The common utilities that support the other Hadoop modules.
- **Hadoop Distributed File System (HDFS™):** A distributed file system that provides high-throughput access to application data.
- **Hadoop YARN (Yet Another Resource Negotiator):** A framework for job scheduling and cluster resource management.
- **Hadoop MapReduce:** A YARN-based system for parallel processing of large data sets.

Módulos do projeto:

ARCHITECTURE COMPARISON Hadoop 1.0 vs. Hadoop 2.0.

Single Use System

Batch Apps

HADOOP 1.0

MapReduce

(cluster resource management
& data processing)

HDFS

(redundant, reliable storage)

Multi Use Data Platform

Batch, Interactive, Online, Streaming, ...

HADOOP 2.0

MapReduce

(batch)

Tez

(interactive)

Others

(varied)

YARN

(operating system: cluster resource management)

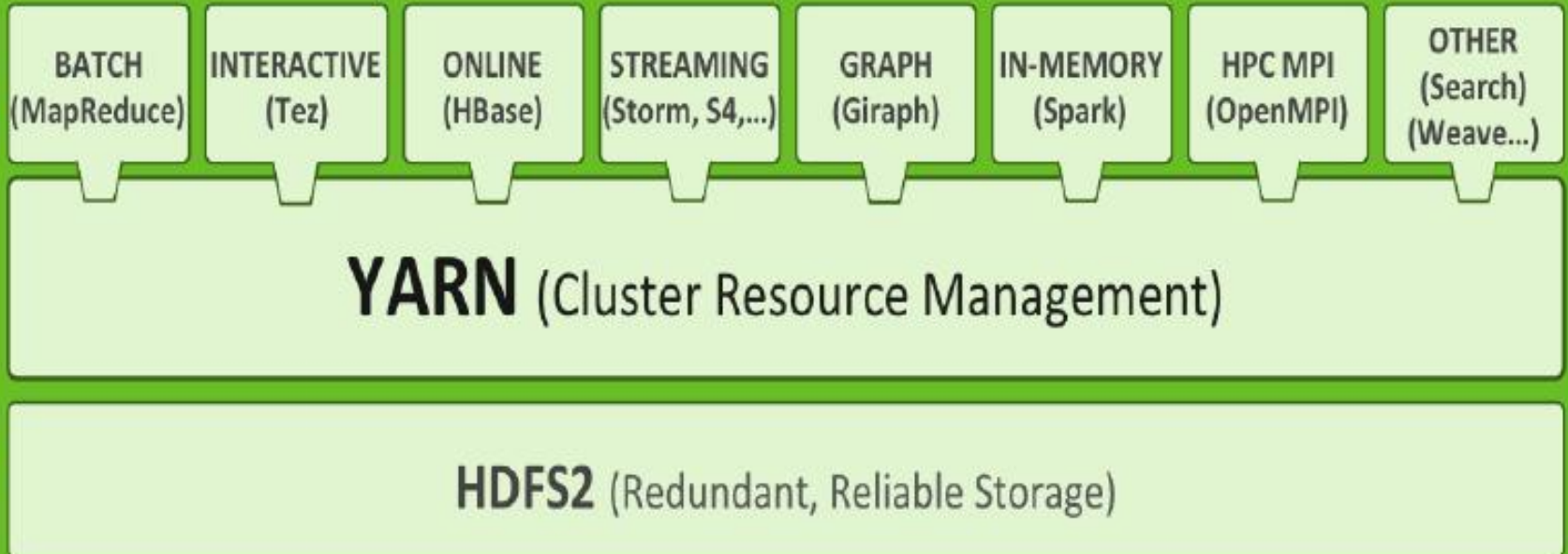
HDFS2

(redundant, reliable storage)

SOURCE: HORTONWORKS

Módulos do projeto:

Applications Run Natively **IN** Hadoop



HDFS:

- Sistema De Arquivos Distribuído
 - Hardware Failure
 - Large Data Sets
 - Designed to be deployed on low-cost hardware

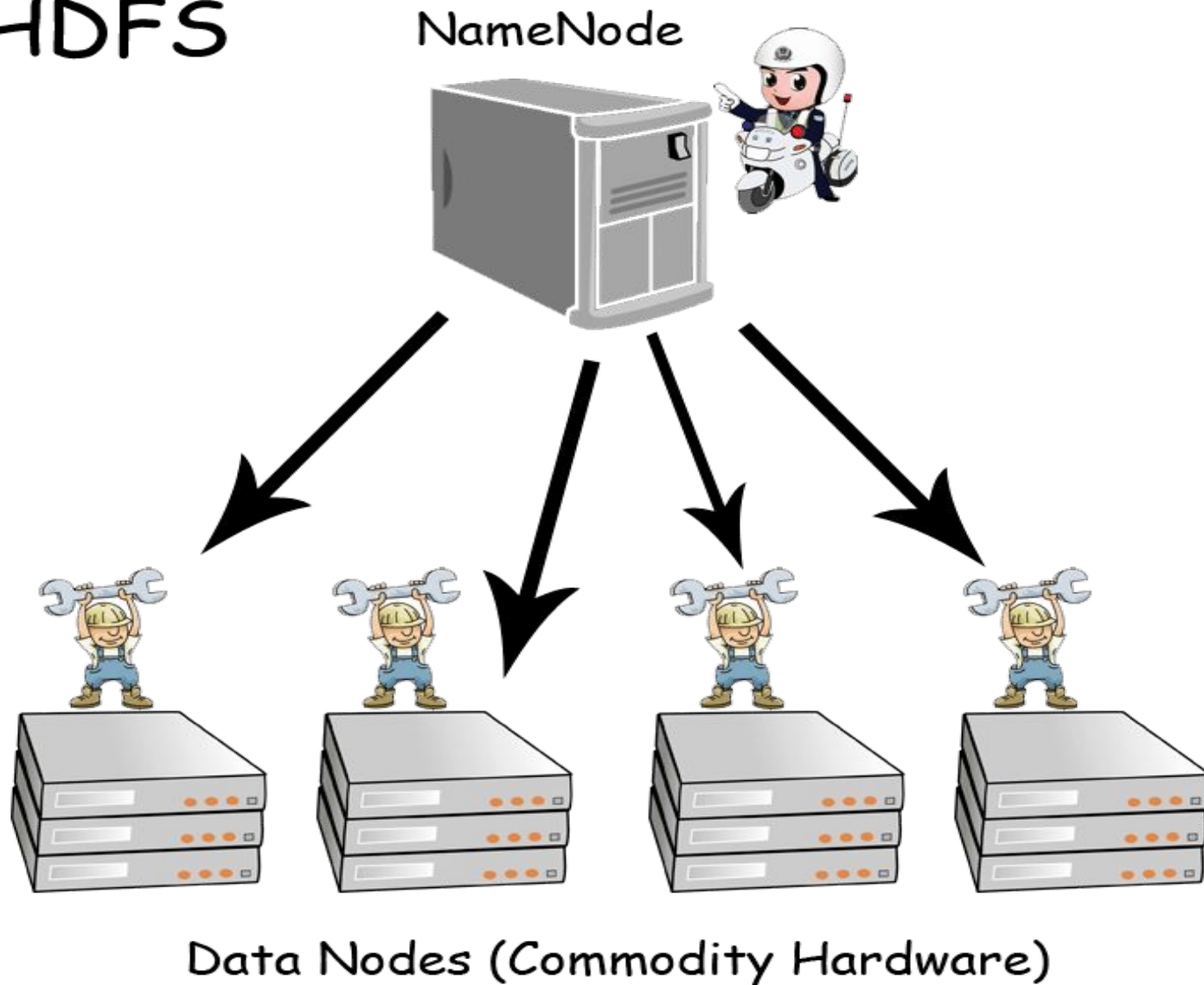
HDFS:

- Sistema De Arquivos Distribuído
 - Hardware Failure
 - Large Data Sets
 - Designed to be deployed on low-cost hardware
- Arquitetura Master/Slave (NameNode/DataNodes)
 - “um” mestre
 - Muitos trabalhadores

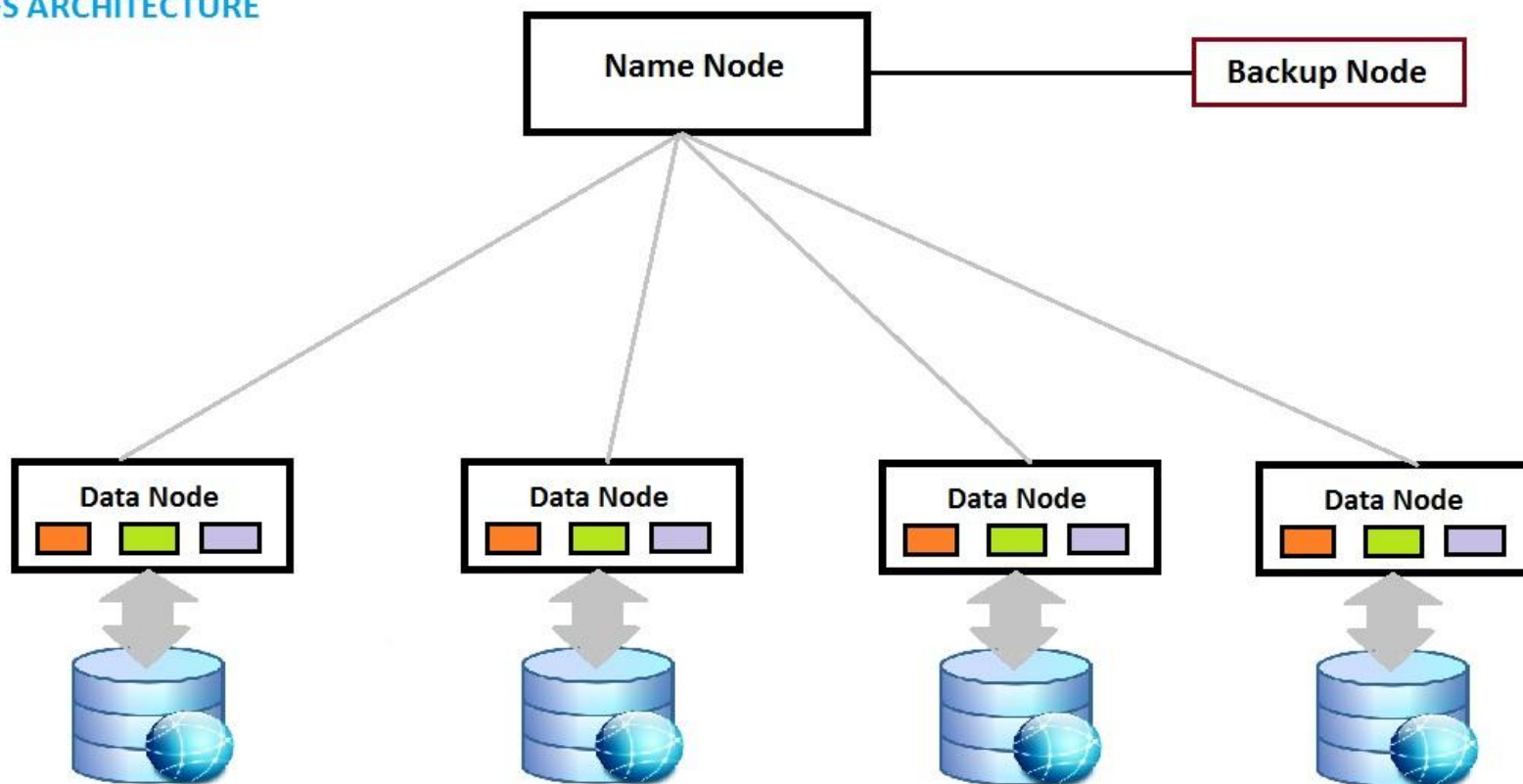
HDFS:

- Sistema De Arquivos Distribuído
 - Hardware Failure
 - Large Data Sets
 - Designed to be deployed on low-cost hardware
- Arquitetura Master/Slave (NameNode/DataNodes)
 - “um” mestre
 - Muitos trabalhadores
- Particionamento dos dados em blocos (256MB)
 - **Replicação (Redundância)**

HDFS



HDFS ARCHITECTURE



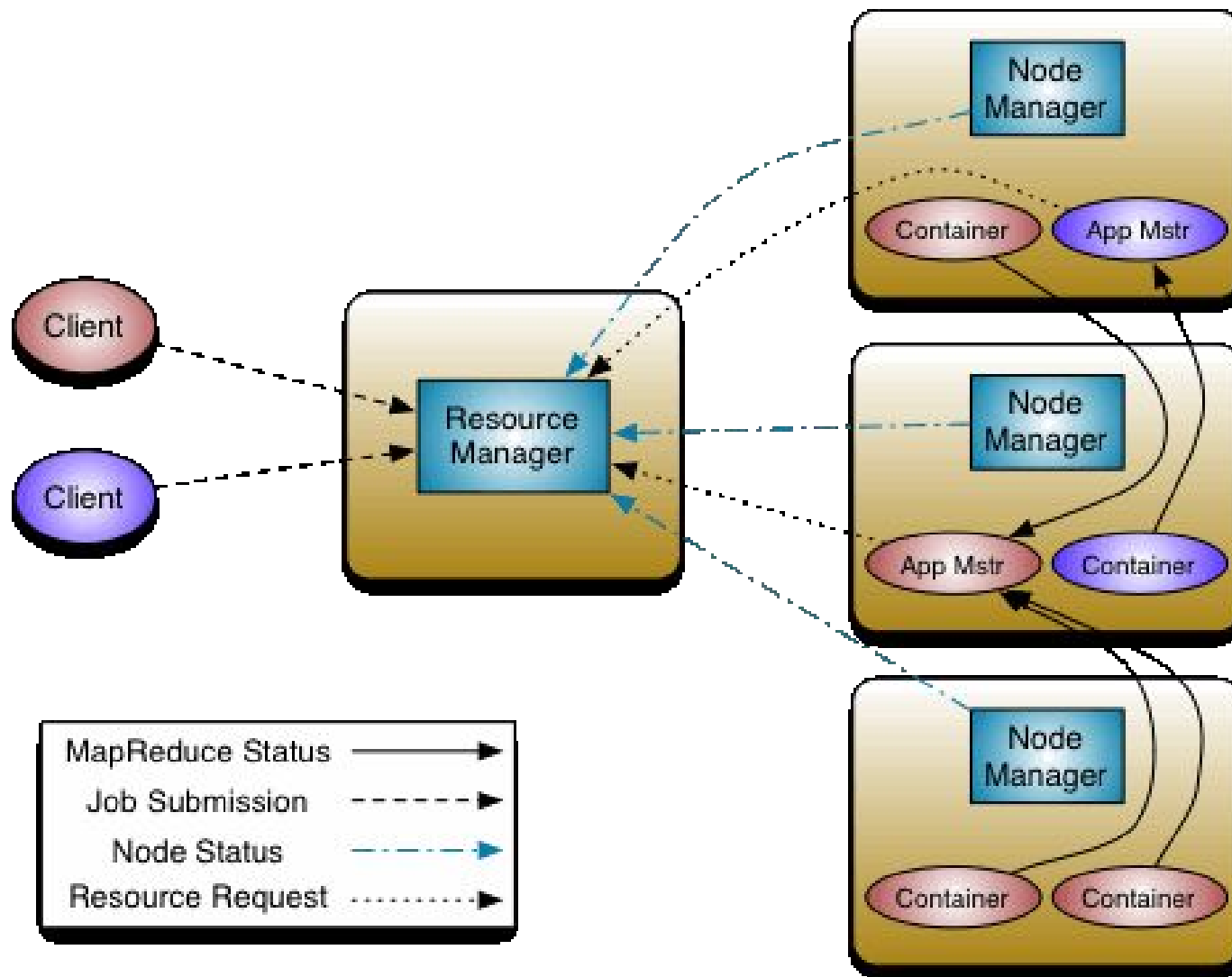
YARN:

- Dividir Gerenciamento de Recursos do gerenciamento dos Jobs (tarefas).
 - Lidar com as Falhas

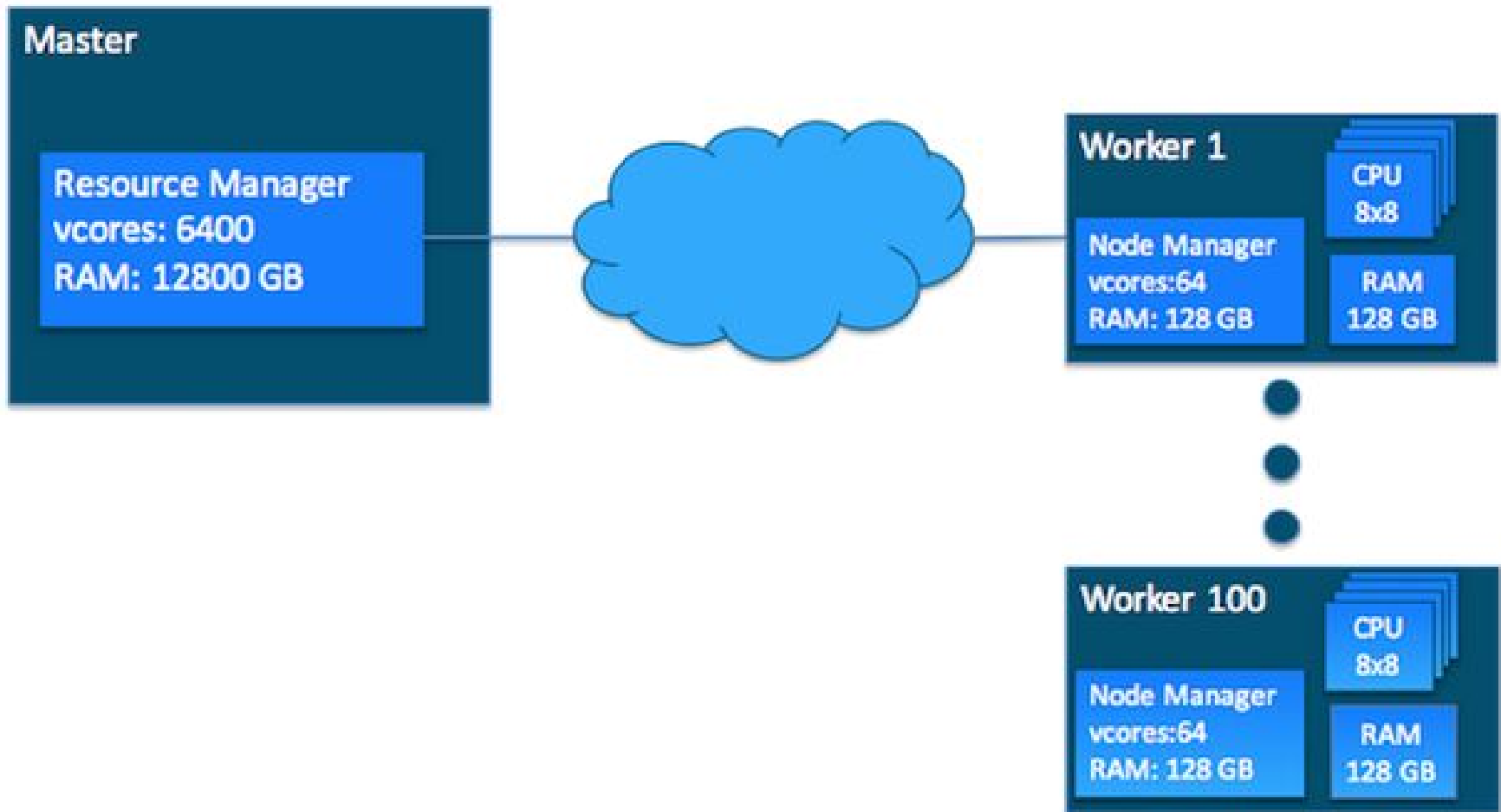
YARN:

- Dividir Gerenciamento de Recursos do gerenciamento dos Jobs (tarefas).
 - Lidar com as Falhas
- **Arquitetura Master/Slave**
 - ResourceManager (Master)
 - NodeManager (Cada nó)
 - ApplicationMaster (Um por aplicação)

YARN:

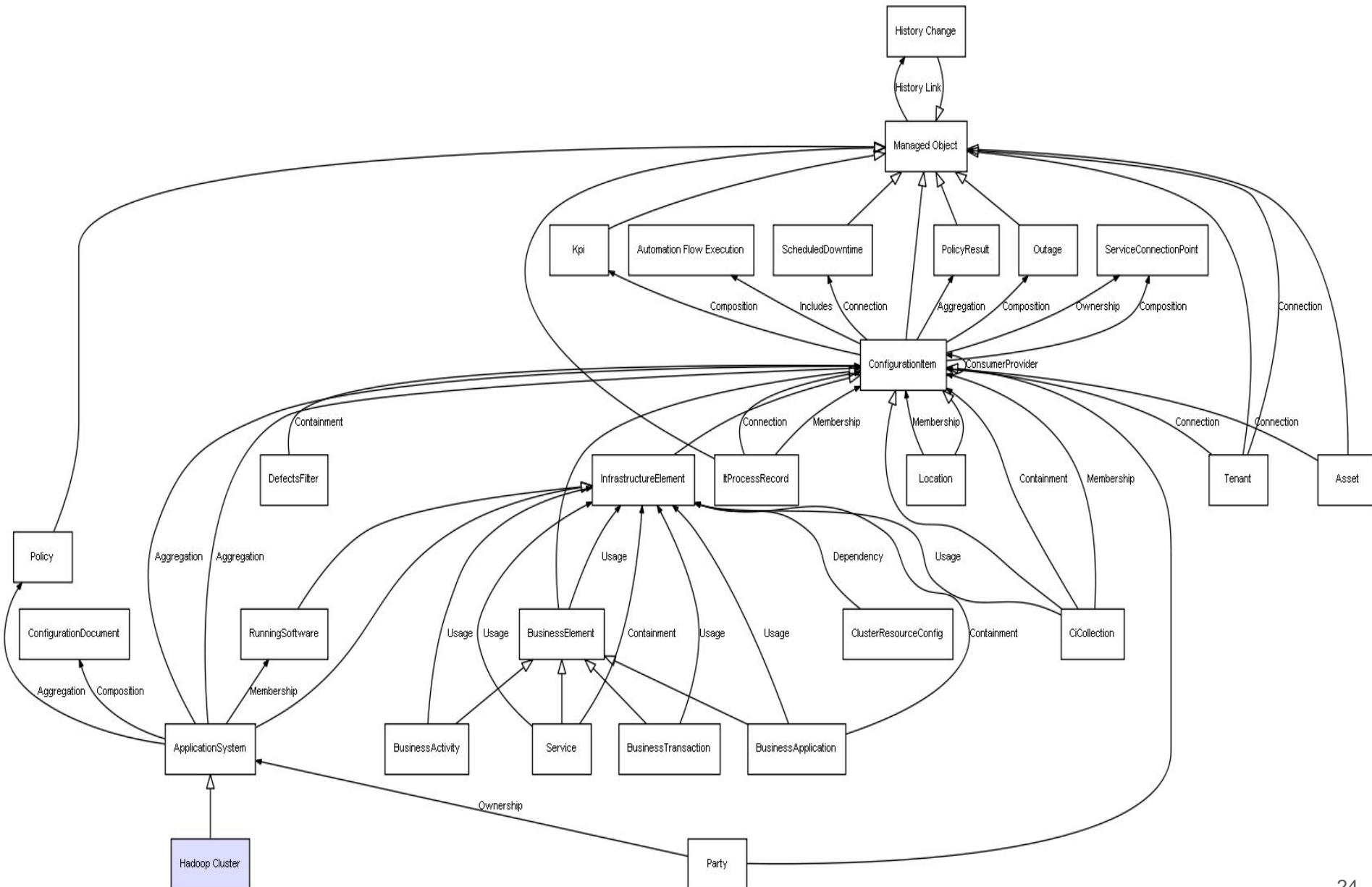


YARN:



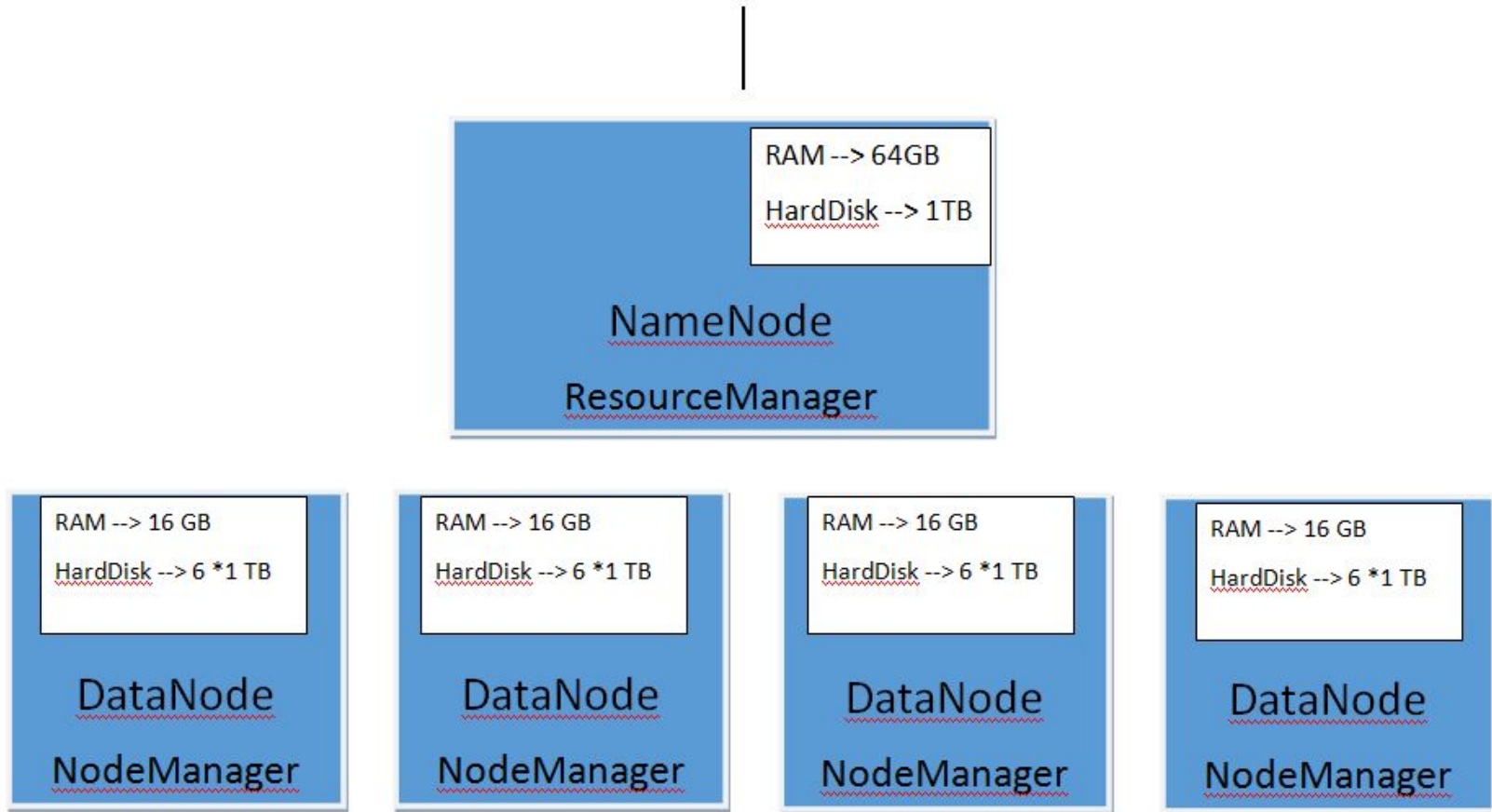
Os dois juntos?

Os dois juntos?



Os dois juntos?

Hadoop Architecture



TIP: processa localmente os blocos disponíveis naquele nó. (mover a computação em vez dos dados)

Map Reduce:

- Paradigma de programação...

Map Reduce:

- Paradigma de programação...
- Não é nada novo...
 - Já existe a muito tempo (linguagens funcionais [Lisp])

Map Reduce:

- Paradigma de programação...
- Não é nada novo...
 - Já existe a muito tempo (linguagens funcionais [Lisp])
- Se baseia em duas coisas.
 - Funções **Map**(<Chave , Valor>) e **Reduce**(<Chave , Valor>)
 - **Map** processa dados, e o **Reduce** processa o resultado do **Map**

$\text{Map}(k1, v1) \rightarrow \text{list}(k2, v2)$

$\text{Reduce}(k2, \text{list}(v2)) \rightarrow \text{list}(k3, v3)$

WordCount...

```
function map(String name, String document):  
    // name: document name  
    // document: document contents  
    for each word w in document:  
        emit (w, 1)  
  
function reduce(String word, Iterator partialCounts):  
    // word: a word  
    // partialCounts: a list of aggregated partial counts  
    sum = 0  
    for each pc in partialCounts:  
        sum += pc  
    emit (word, sum)
```

WordCount...

Texto: José, João, Maria, João, Maria, Felipe, Alan

Map $\rightarrow \{ \langle \text{José}, 1 \rangle, \langle \text{João}, 1 \rangle, \langle \text{Maria}, 1 \rangle, \langle \text{José}, 1 \rangle, \langle \text{João}, 1 \rangle, \langle \text{Maria}, 1 \rangle, \langle \text{Felipe}, 1 \rangle, \langle \text{Alan}, 1 \rangle \}$

Mágica $\rightarrow \langle \text{José}, (1,1) \rangle, \langle \text{João}, (1,1) \rangle, \langle \text{Maria}, (1,1) \rangle, \langle \text{Felipe}, (1) \rangle, \langle \text{Alan}, (1) \rangle$

Reduce($\langle \text{José}, (1,1) \rangle$) $\rightarrow \langle \text{José}, 2 \rangle$

Reduce($\langle \text{João}, (1,1) \rangle$) $\rightarrow \langle \text{João}, 2 \rangle$

Reduce($\langle \text{Maria}, (1,1) \rangle$) $\rightarrow \langle \text{Maria}, 2 \rangle$

Reduce($\langle \text{Felipe}, (1) \rangle$) $\rightarrow \langle \text{Felipe}, 1 \rangle$

Reduce($\langle \text{Alan}, (1) \rangle$) $\rightarrow \langle \text{Alan}, 1 \rangle$

E o Hadoop MapReduce?

E o Hadoop MapReduce?

- Framework para escrever facilmente aplicações MapReduce que processam grandes quantidades de dados, de forma paralela, distribuída em grandes clusters de hardware commodity com tolerância a falhas.

E o Hadoop MapReduce?

- Uma aplicação (*Job*) MapReduce geralmente divide os dados de entrada em *Blocos* independentes que são processados pelos Maps paralelamente.
- O framework organiza (sort/shuffle) as saídas dos Maps que serão as entradas dos Reducers (Mágica).
- Geralmente a entrada e a saída dos *Jobs* estão/são armazenadas no HDFS.
- O framework se encarrega de organizar e monitorar as tarefas, e reexecutar as que falharam.

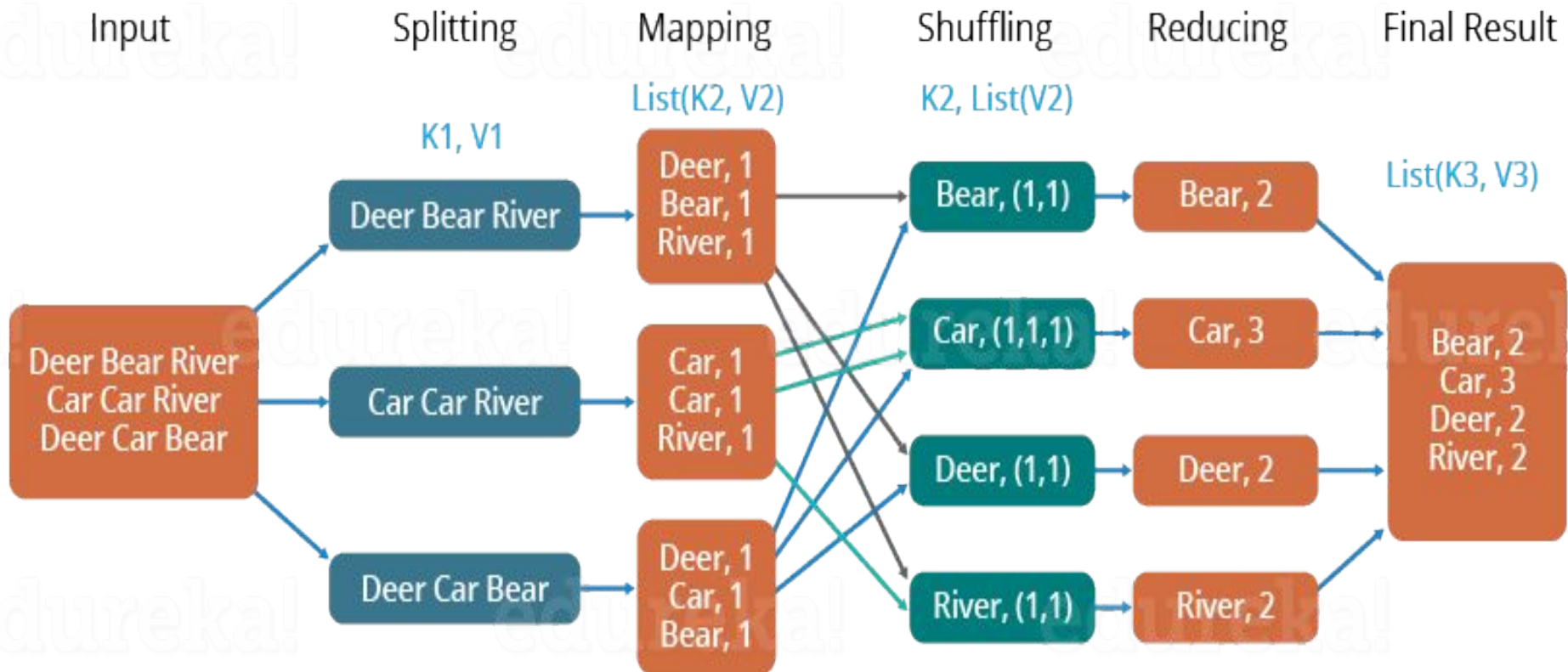
E o Hadoop MapReduce?

- Uma aplicação (*Job*) MapReduce geralmente divide os dados de entrada em *Blocos* independentes que são processados pelos Maps paralelamente.
 - O framework organiza (sort/shuffle) as saídas dos Maps que serão as entradas dos Reducers (Mágica).
 - Geralmente a entrada e a saída dos *Jobs* estão/são armazenadas no HDFS.
 - O framework se encarrega de organizar e monitorar as tarefas, e reexecutar as que falharam.
-
- ❖ Minimamente um Job MapReduce precisa:
 - Especificar a localização da entrada e da saída;
 - Prover as funções Map e Reduce;
 - Implementação das interfaces apropriadas/classes abstratas
 - Esses e outros parâmetros compõem a configuração do *Job*;
 - ❖ O cliente então submete o Jar/Executável para o ResourceManager
 - O ResourceManager assume daí;
 - Ele se vira com o resto;

WordCount...

WordCount...

The Overall MapReduce Word Count Process



Agora um pouco de Java...

origem:

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

Agora um pouco de 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;

public class WordCount {
    .
    .
    .
}
```

Agora um pouco de Java...

```
public class WordCount {  
  
    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);  
            }  
        }  
    }  
}
```

.
.
.

Agora um pouco de Java...

```
public class WordCount {  
    .  
    .  
    .  
    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);  
        }  
    }  
    .  
    .  
    .  
}
```


Agora um pouco de Java...

```
public class WordCount {  
    .  
    .  
    .  
  
    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);  
    }  
  
} //fim WordCount
```

Vamos rodar de verdade...

DESAFIO:

desafio.txt

Cidade	Dia da Semana	Temp. Mínima	Temp. Máxima	Elevação
Goiânia	1	22	33	223
Goiânia	2	23	25	223
-	-	-	-	-
-	-	-	-	-
Goiânia	7	29	31	223