# - Caraduação



# DATA SCIENCE BIG DATA ARCHITECTURING & DATA INTEGRATION Prof. Dr. Renê de Ávila Mendes

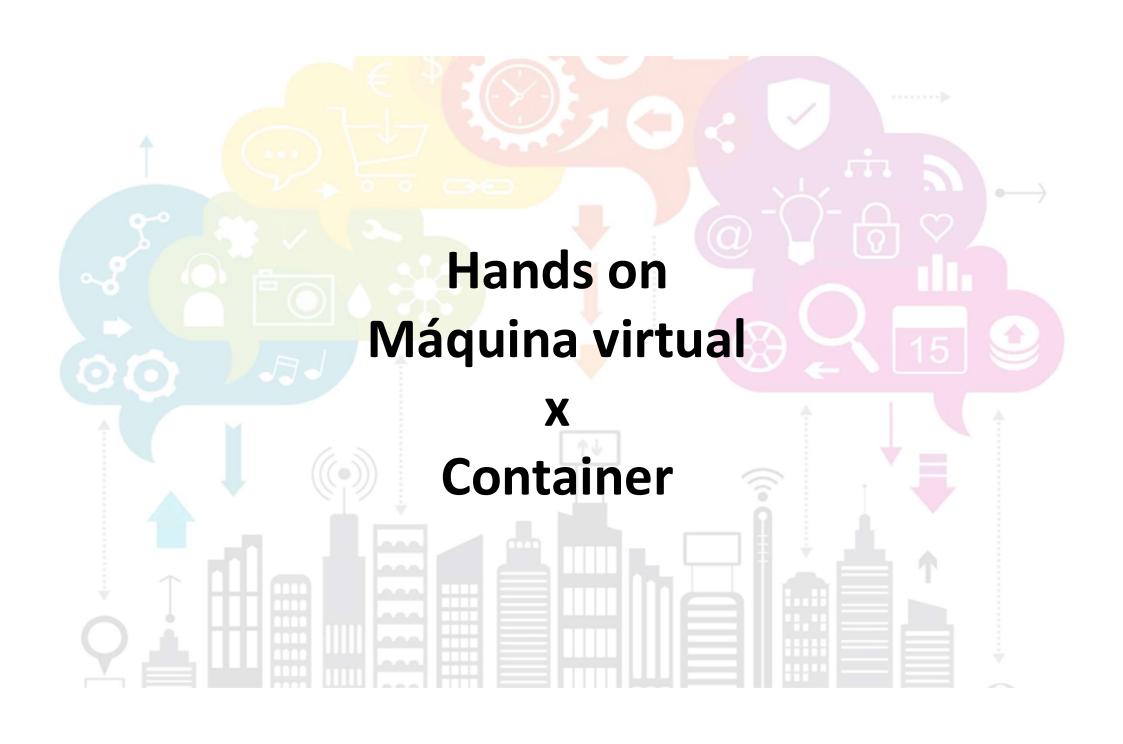
### Objetivos da disciplina

**DISCIPLINA:** Big Data Architecturing & Data Integration

OBJETIVOS: Entenda as principais arquiteturas para ingestão, processamento e análise de grandes volumes de dados. Conheça as principais ferramentas open-source de Big Data como Hadoop, MapReduce, Spark, Sqoop, NiFi, Flume, Kafka, Zookeeper, HBase, Hive e as integre com as ferramentas de extração, transformação e carga de dados em modelos dimensionais. Entenda conceitos sobre computação paralela e distribuída, aplicação do Hadoop e bases Apache e arquiteturas serverless e desacopladas. Veja como visualizar os dados estruturados ou não estruturados com ferramentas de Self-Service Business Intelligence como PowerBI, utilizando as melhores práticas de visualização de dados.

### Assuntos – 1º Semestre

- Introdução a Big Data
- Conceitos Computação Paralela e Distribuída, Lei de Moore, Sistema HDFS
- Aplicação Hadoop. Usos e Administração de Ambientes Hadoop
- Introdução a MapReduce
- Introdução à Integração de Dados
- Integração entre SQL e Hadoop SQOOP
- Bases Apache
- Bases Apache PIG
- Introdução da Data Streaming FLUME
- Introdução a análise de dados com SPARK

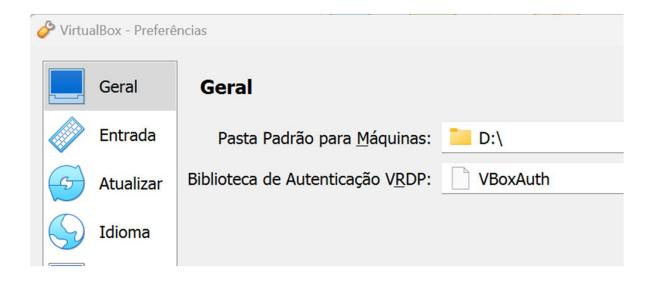


### VM Hadoop

Alterar a configuração do local:

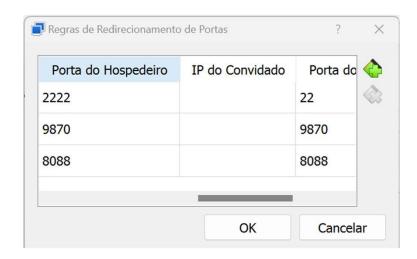
Arquivo > Preferências > Pasta padrão para máquinas: "D:\"

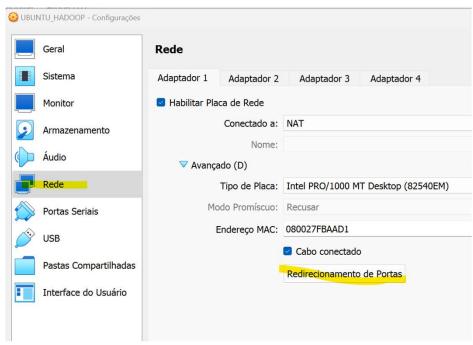
- Importar a imagem
- Usuário e senha
  - hadoop / hadoop

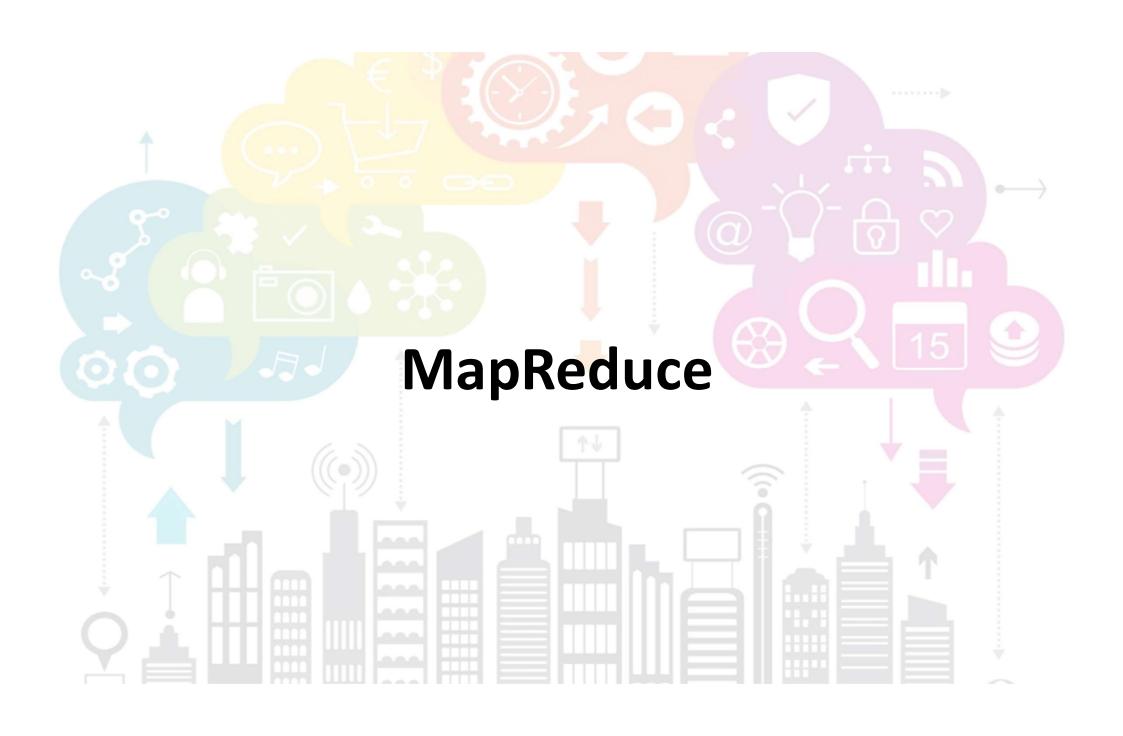


# VM Hadoop

- Mapear portas
  - -SSH









### Map Reduce

- Em 2004 o Google publicou o artigo chamado MapReduce: Simplified Data Processing on Large Clusters.
- MapReduce é um modelo de programação para processamento de grandes volumes de dados.
- Esta abstração foi inspirada nas primitivas map e reduce do Lisp e outras linguagens funcionais.

MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

Google, Inc.

### Abstract

MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a map function that processes a key/value pair to generate a set of intermediate key/value pairs, and a reduce function that merges all intermediate values associated with the same intermediate key. Many real world tasks are expressible in this model, as shown in the paper.

Programs written in this functional style are automatically parallelized and executed on a large cluster of commodity machines. The run-time system takes care of the details of partitioning the input data, scheduling the program's execution across a set of machines, handling machine failures, and managing the required inter-machine communication. This allows programmers without any experience with parallel and distributed systems to easily utilize the resources of a large distributed system.

Our implementation of MapReduce runs on a large cluster of commodity machines and is highly scalable: a typical MapReduce computation processes many terabytes of data on thousands of machines. Programmers find the system easy to use: hundreds of MapReduce programs have been implemented and upwards of one thousand MapReduce jobs are executed on Google's clusters every day.

### 1 Introduction

Over the past five years, the authors and many others at Google have implemented hundreds of special-purpose computations that process large amounts of raw data, such as crawled documents, web request logs, etc., to compute various kinds of derived data, such as inverted indices, various representations of the graph structure of web documents, summaries of the number of pages crawled per host, the set of most frequent queries in a

given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines in order to finish in a reasonable amount of time. The issues of how to parallelize the computation, distribute the data, and handle failures conspire to obscure the original simple computation with large amounts of complex code to deal with these issues.

As a reaction to this complexity, we designed a new abstraction that allows us to express the simple computa tions we were trying to perform but hides the messy details of parallelization, fault-tolerance, data distribution and load balancing in a library. Our abstraction is inspired by the map and reduce primitives present in Lisp and many other functional languages. We realized that most of our computations involved applying a map operation to each logical "record" in our input in order to compute a set of intermediate key/value pairs, and then applying a reduce operation to all the values that shared the same key, in order to combine the derived data anpropriately. Our use of a functional model with user specified map and reduce operations allows us to parallelize large computations easily and to use re-execution as the primary mechanism for fault tolerance.

The major contributions of this work are a simple and powerful interface that enables automatic parallelization and distribution of large-scale computations, combined with an implementation of this interface that achieves high performance on large clusters of commodity PCs.

Section 2 describes the basic programming model and gives several examples. Section 3 describes an implementation of the MapReduce interface tallored towards our cluster-based computing environment. Section 4 describes several refinements of the programming model that we have found useful. Section 5 has performance measurements of our implementation for a variety of tasks. Section 6 explores the use of MapReduce within Google including our experiences in using it as the basis

To appear in OSDI 2004



### Map e Reduce no Lisp

 Lisp é uma linguagem funcional, projetada por John McCarthy que apareceu pela primeira vez em 1958.

 A primitiva map recebe como parâmetros um operador unário e uma lista, onde este operador será chamada para cada elemento da lista

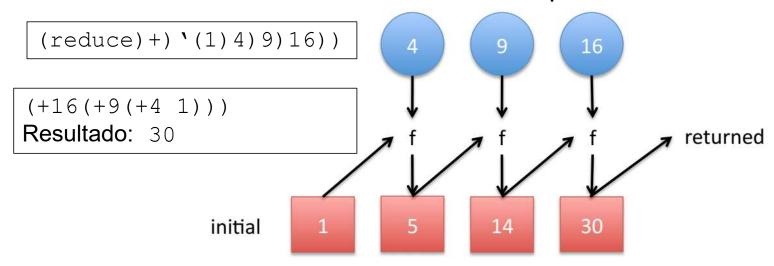
Soma do quadrado de uma lista (map) square) '(1)2)3)4))) **Resultado**: (1 4 9 16)

https://cse.buffalo.edu/~stevko/courses/cse704/fall10/papers/cse704-mapreduce.pdf



### Map e Reduce no Lisp

 A primitiva reduce recebe como parâmetros um operador binário e uma lista, onde os elementos da lista serão combinados utilizando este operador





### Map Reduce em Lisp

- Map
  - processa cada registro individualmente
- Reduce
  - Processa (combina) o conjunto de todos os registros em lote



### O que a Google percebeu

- Pesquisa por palavra-chave
- Encontre uma palavra-chave em cada página da web individualmente e, se ela for encontrada, retorne sua URL
- Reduce
- Combine todos os resultados (URLs) e devolva-os
- Contagem do número de ocorrências de cada palavra
  - Conte o número de ocorrências em cada página da web individualmente e retorne a lista de <palavra, número>
- Reduce

Map

Para cada palavra, some (combine) a contagem



### Visão Geral do Map Reduce



Os dados na fase map são divididos entre os nós task tracker onde os dados estão localizados. Cada nó na fase map emite pares chave-valor com base no registro de entrada, um registro por vez.



A fase shuffle é tratada pela estrutura Hadoop. Ela transfere os resultados dos mappers para os reducers juntando os resultados por meio chave.

Reduce

A saída dos mappers é enviado para os reducers. Os dados na fase reduce são divididos em partições onde cada reducer lê uma chave e uma lista de valores associados a essa chave. Os reducers emitem zero ou mais pares chavevalor com base na lógica utilizada

Visão Geral do Map Reduce FIAP

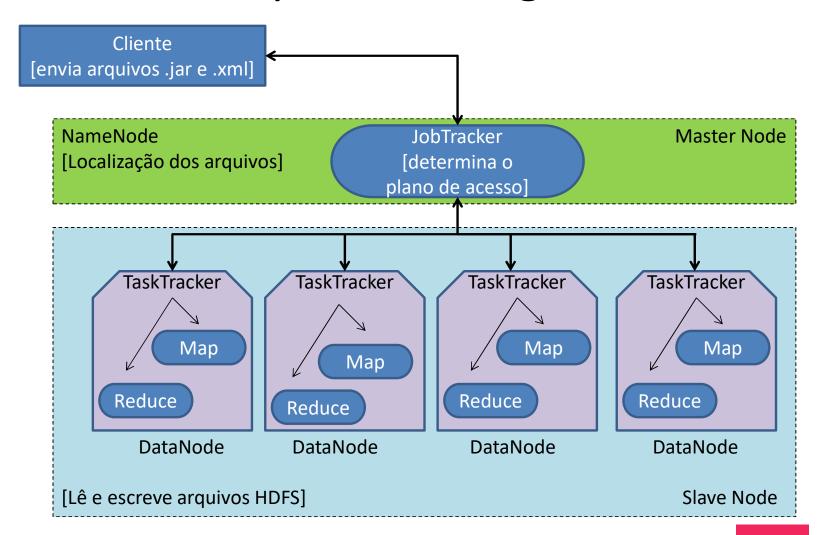


### Visão Geral do Map Reduce FIAP

A 1.1 bela 2 bela bela bela 1.1.1 belo 3 jovem bela belo belo cidade cidade saiu A bela jovem saiu com belo com com um belo desta desta com um belo rapaz belo cidade jovem jovem belo para rapaz com no no desta para para para belo passeio jovem um parque parque parque desta bela belo no passeio passeio para passeio rapaz " rapaz cidade no parque saiu saiu 1.1 belo passeio um parque rapaz desta saiu bela cidade um

# MapReduce Engine







 Incluindo as palavras dog, cat, mouse e hippo em um banco orientado a documentos.

```
db.items.insert({tags: ['dog', 'cat']})
db.items.insert({tags: ['dog']})
db.items.insert({tags: ['dog', 'mouse']})
db.items.insert({tags: ['dog', 'mouse', 'hippo']})
db.items.insert({tags: ['dog', 'mouse', 'hippo']})
db.items.insert({tags: ['dog', 'hippo']})
```



Criando o mapeamento

```
var map = function() {
  this.tags.forEach(function(t) {
    emit(t, {count: 1});
  });
}
```



Criando o reduce

```
var reduce = function(key, values) {
  var count = 0;
  for(var i=0, len=values.length; i<len; i++) {
    count += values[i].count;
  }
  return {count: count};
}</pre>
```



Executando o MapReduce



Executando o MapReduce

> db[result.result].find()

```
{ "_id" : "cat", "value" : { "count" : 1 } } 
{ "_id" : "dog", "value" : { "count" : 6 } } 
{ "_id" : "hippo", "value" : { "count" : 3 } } 
{ "_id" : "mouse", "value" : { "count" : 3 } }
```



# Submetendo um job MapReduce

```
mapear a porta 22 para 2222
copiar os arquivos para uma pasta da VM
formatar o namenode
iniciar o hadoop
copiar o arquivo shakespeare.txt
hdfs dfs -put shakespeare.txt /
```

https://hadoop.apache.org/docs/r3.3.4/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html#Example:\_WordCount\_v1.0

```
export JAVA_HOME=/usr/java/default
export PATH=${JAVA_HOME}/bin:${PATH}
export HADOOP_CLASSPATH=${JAVA_HOME}/lib/tools.jar
hadoop com.sun.tools.javac.Main WordCount.java
jar cf wc.jar WordCount*.class
hadoop jar wc.jar WordCount /shakespeare.txt /saida
```