

Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3



Bem-vindo(a)



Instrutores



Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

Data Science Academy



A **Data Science Academy (DSA)** é um portal de ensino online especializado em **Big Data, Machine Learning, Inteligência Artificial, Desenvolvimento de Chatbots, Blockchain e tecnologias relacionadas.**

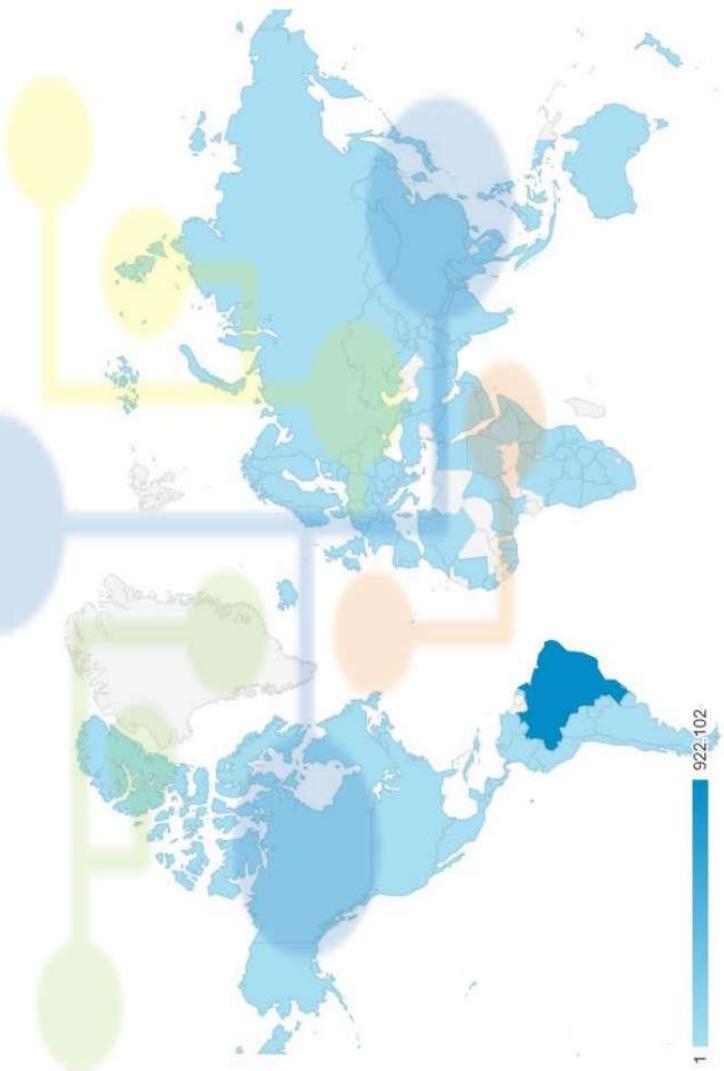
Nosso objetivo é fornecer aos alunos conteúdo de alto nível por meio do uso de computador, tablet ou smartphone, em qualquer lugar, a qualquer hora, 100% online e 100% em português.

Instrutores

Data Science
Academy

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

No Brasil e no Mundo

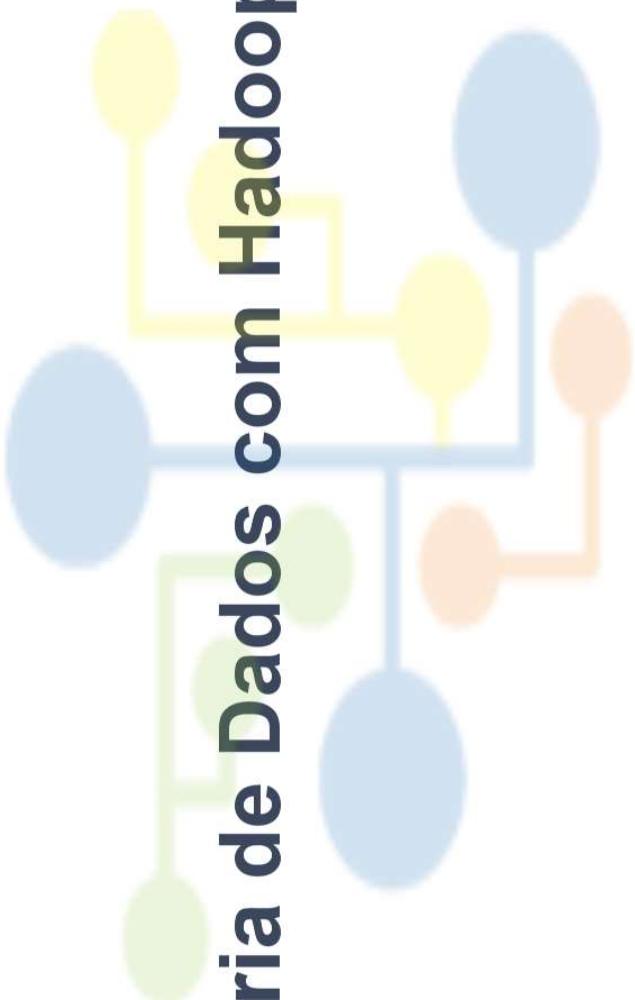




Data Science
Academy

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

Engenharia de Dados com Hadoop e Spark



Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Cluster
Hadoop
Capítulos
2, 3 e 4

Armazenamento
de Dados
Capítulos
5, 6 e 7

Machine
Learning
Capítulo
8

Hadoop e
Spark
Capítulo
9



Engenharia de Dados com Hadoop e Spark

Data Science Academy phelipe.ursempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

O que você vai aprender neste curso?

Conceitos e definições de Big Data, Hadoop,
Ecossistema Hadoop e Spark

Como planejar, instalar e configurar um
cluster Hadoop

Como planejar, instalar e configurar o Ecossistema
Hadoop (Hive, Hbase, Zookeeper, Flume, Oozie,
Ambari, Sqoop, Spark e Storm)

Configuração e utilização do HDFS e
configurações avançadas do cluster Hadoop

Administração e Manutenção do Hadoop e
Spark



Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

O que você vai aprender
neste curso?

Machine Learning com Apache Mahout

Importação/exportação de dados e ETL com Sqoop

Principais distribuições Hadoop do mercado:
Cloudera e Hortonworks

Infraestrutura de Big Data

Análise de Big Data



Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Curso Big Data Fundamentos 2.0

E quais são os pré-requisitos?

- | | | |
|---|--|---|
| Conhecimentos básicos de sistema operacional Linux (desejável) | Conhecimentos básicos de linguagem de programação (desejável) | Muita vontade de aprender e entrar no mundo do Big Data (mandatório) |
|---|--|---|



Engenharia de Dados com Hadoop e Spark

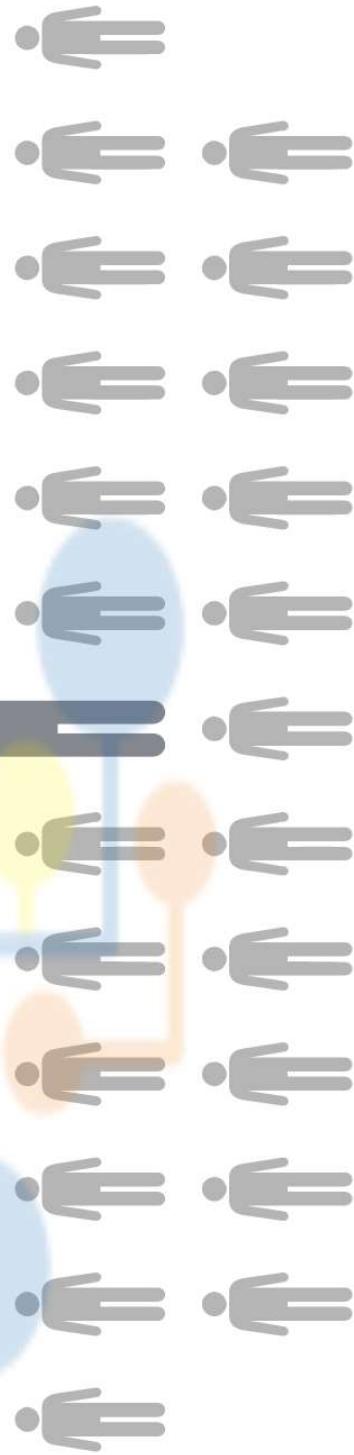


Data Science Academy phelipe.ufssemprboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Benefícios em realizar este curso

São muitos os benefícios em realizar este curso:

Big Data é uma das áreas que mais crescem atualmente. Há um déficit de profissionais no mercado e estima-se que até 2020 o mercado precisará de mais de 200 mil profissionais habilitados em Big Data.



Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Benefícios em realizar este curso

São muitos os benefícios em realizar este curso:

Hadoop é a tecnologia base da infraestrutura de Big Data, que está revolucionando o mundo como o conhecemos. Ele permite a análise de grandes volumes de dados para tomada de decisão.



Engenharia de Dados com Hadoop e Spark

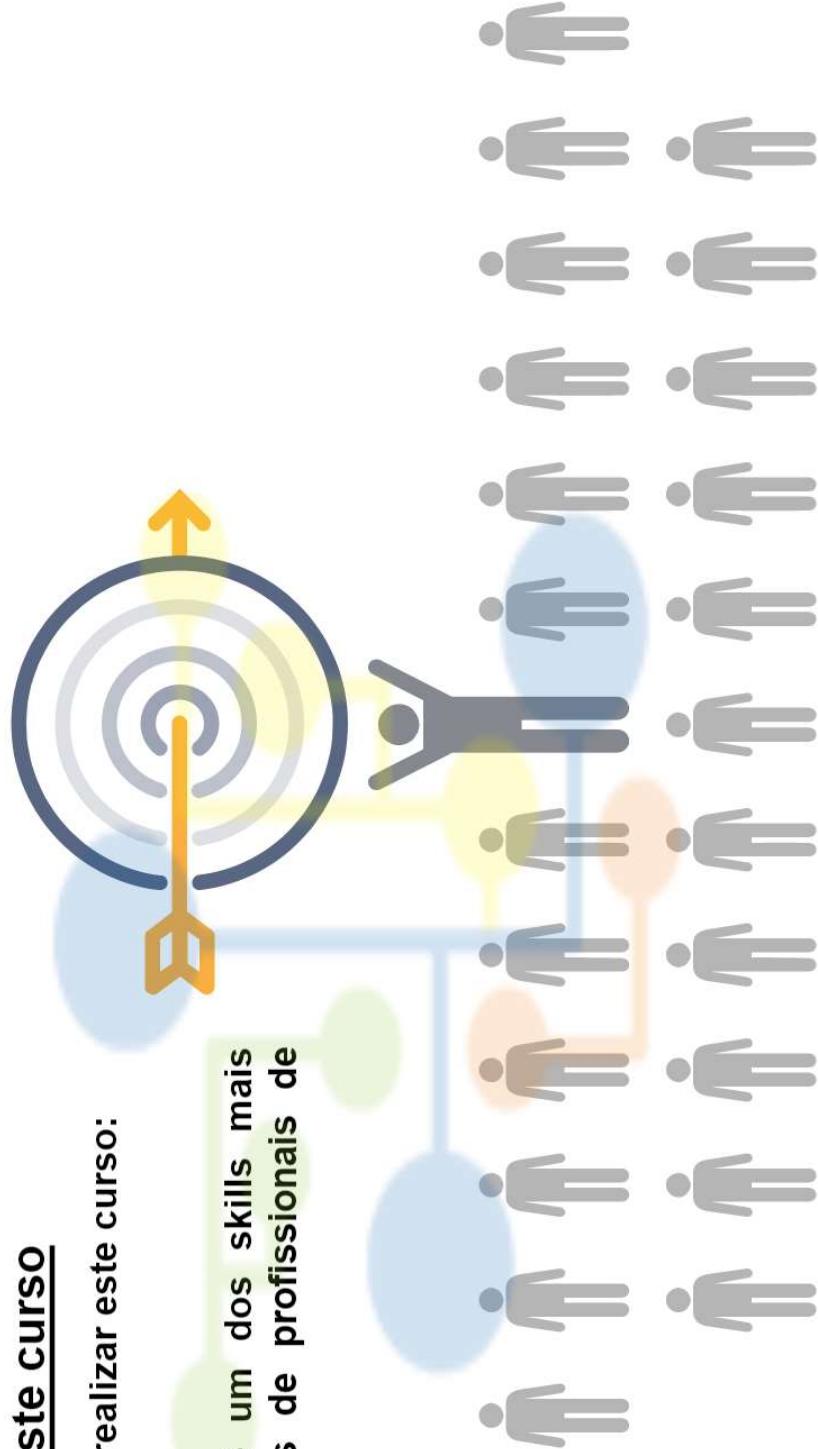


Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Benefícios em realizar este curso

São muitos os benefícios em realizar este curso:

Conhecimento de Hadoop é um dos skills mais procurados por recrutadores de profissionais de Big Data.



Engenharia de Dados com Hadoop e Spark

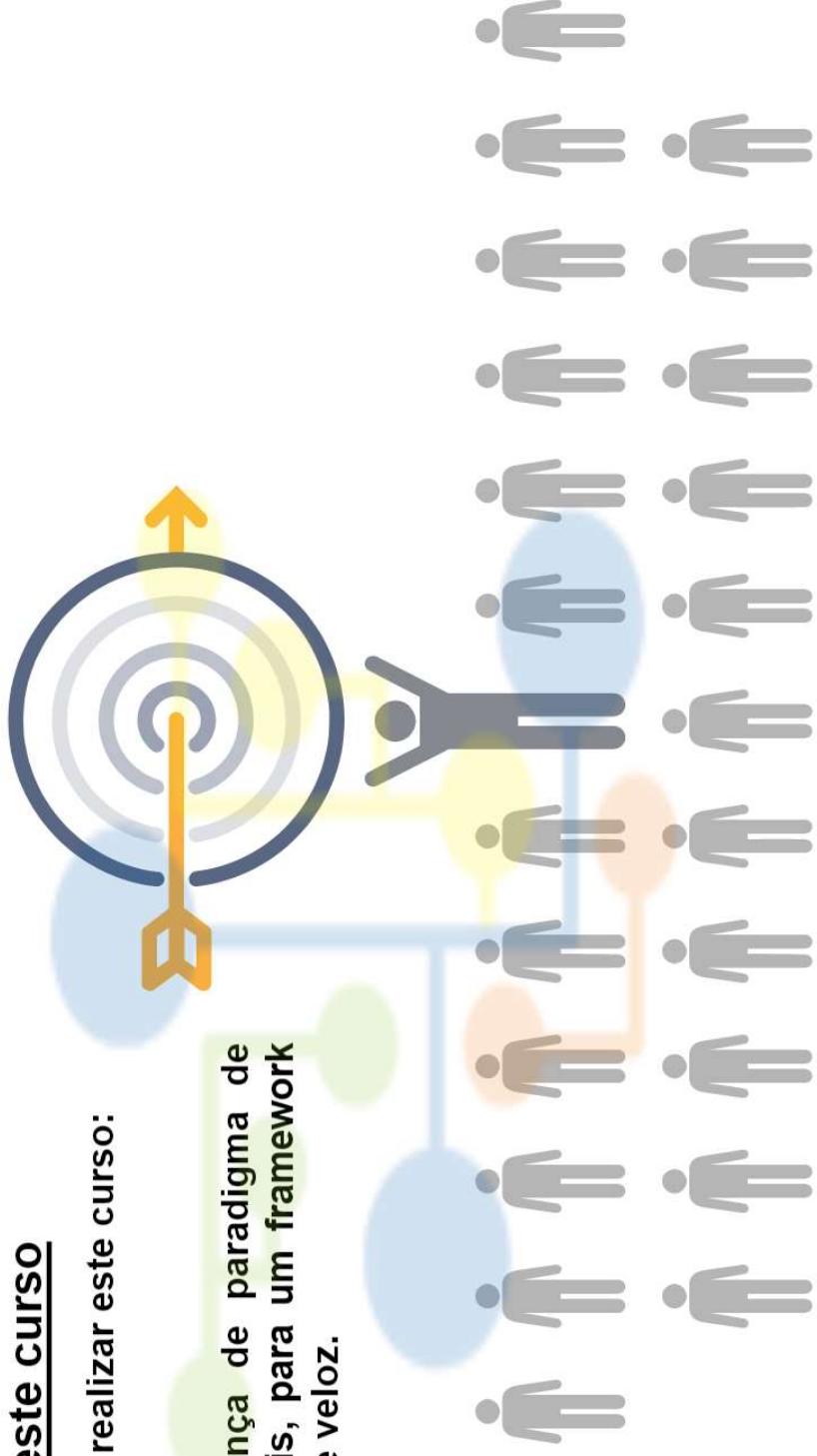


Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Benefícios em realizar este curso

São muitos os benefícios em realizar este curso:

O Hadoop permite a mudança de paradigma de bancos de dados tradicionais, para um framework de dados versátil, adaptável e veloz.



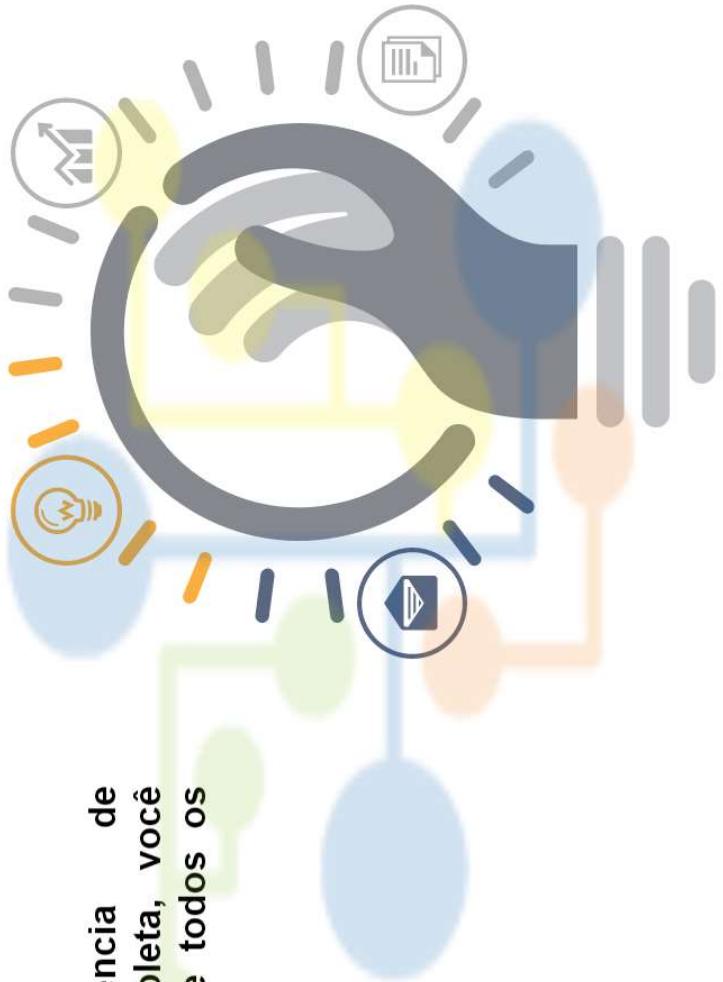
Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Estrutura do Curso

Para tornar sua experiência de aprendizagem ainda mais completa, você terá quizzes e labs ao longo de todos os capítulos.



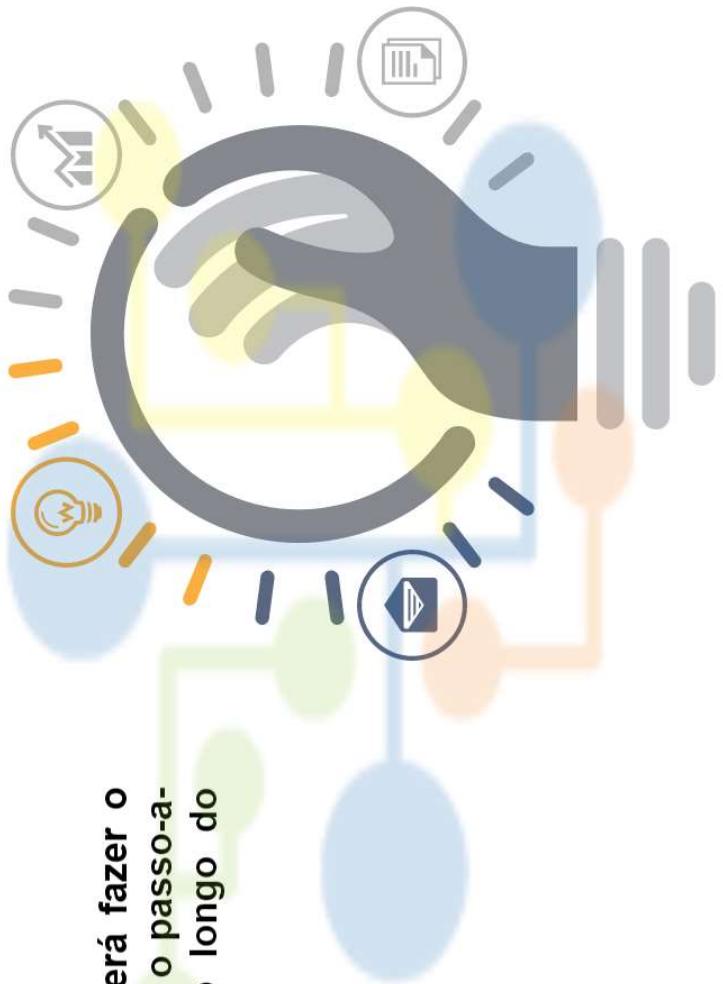
Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Estrutura do Curso

Você também terá acesso e poderá fazer o download dos e-books com todo o passo-a-passo de cada lab realizado ao longo do curso.



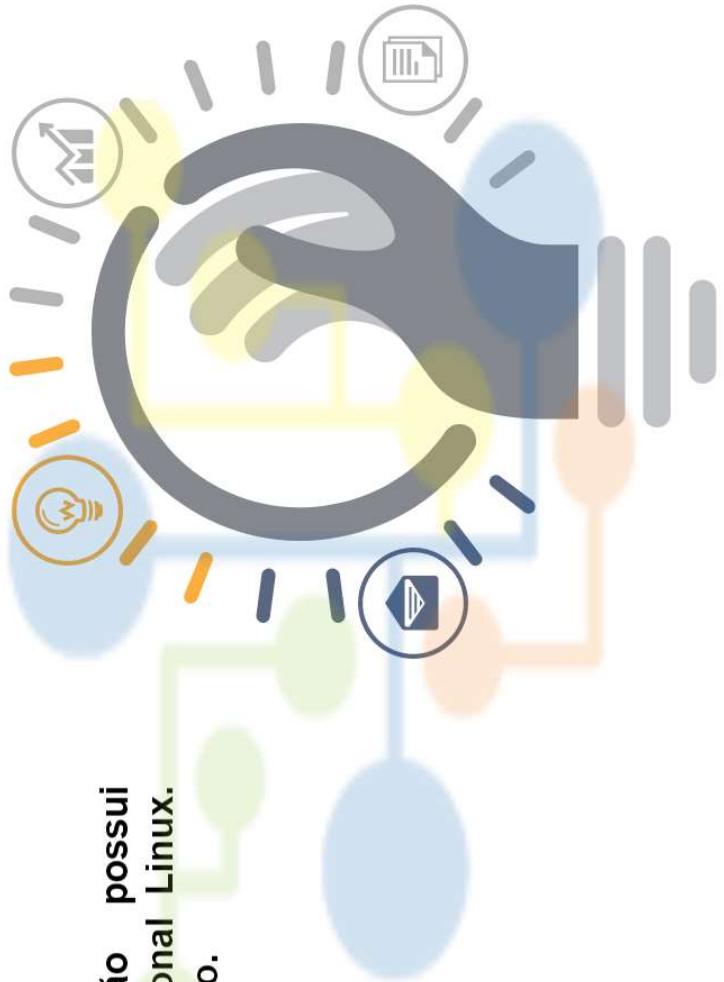
Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Estrutura do Curso

Fique tranquilo se você não possui experiência em sistema operacional Linux. Tudo será explicado passo a passo.



Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Mini-Projetos

Mini-Projeto 1

Importando Dados
do Banco de Dados
Oracle para o
Hadoop com Sqoop



Mini-Projeto 2

Prevendo Casos de
Doenças Cardíacas



Mini-Projeto 3

Design de um Job
MapReduce para
os Gastos Totais
por Cliente



Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Projetos com Feedback

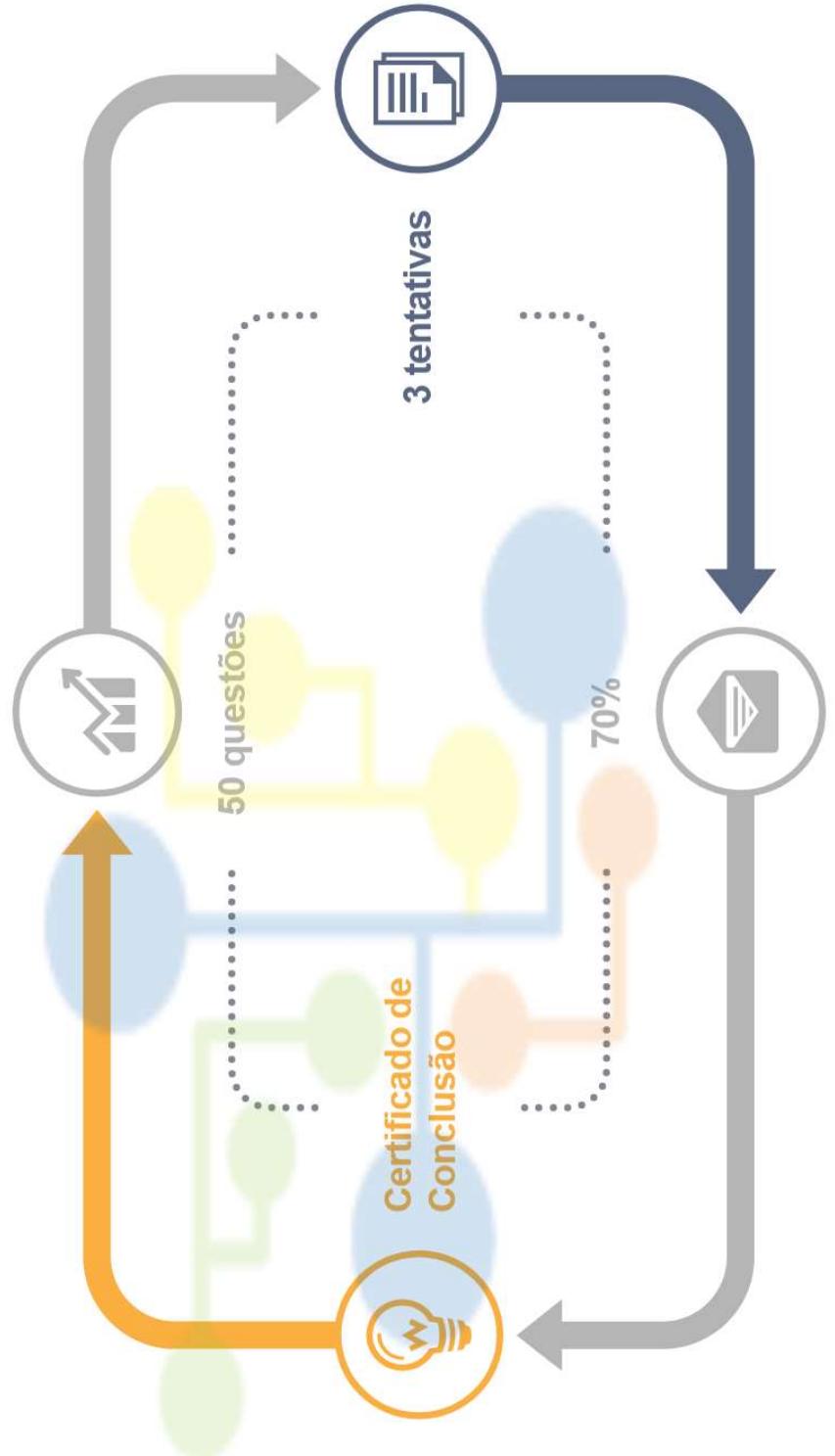


Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufssemprboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Avaliação Final

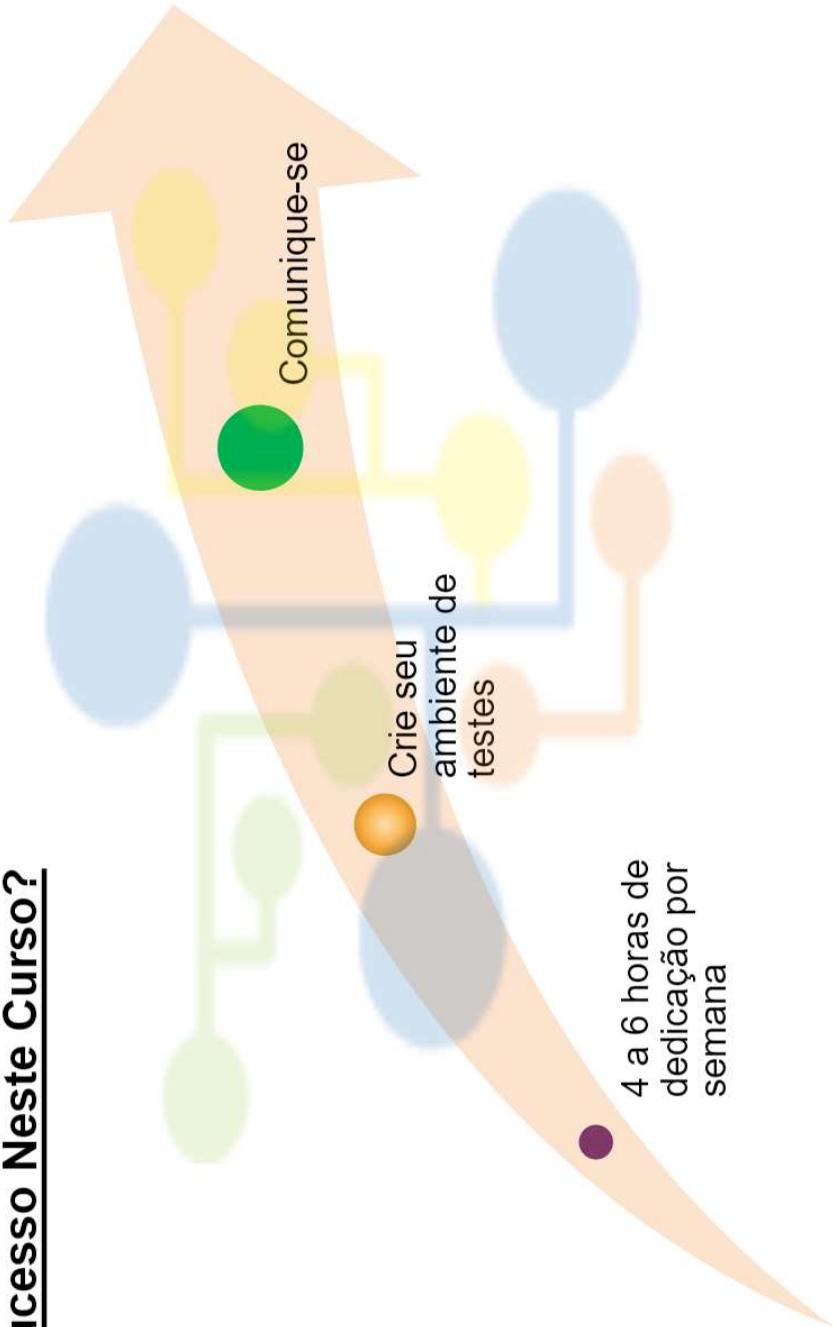


Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

Como Obter Sucesso Neste Curso?



Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3



Engenharia de Dados com Hadoop e Spark



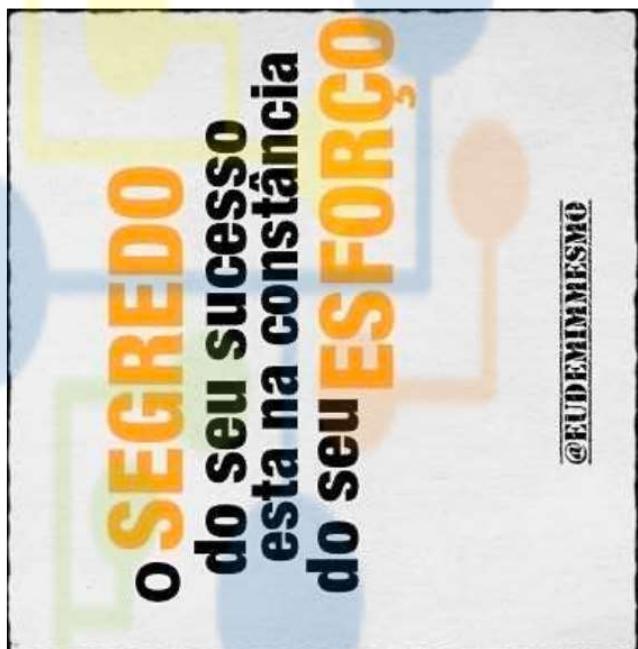
Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3



Engenharia de Dados com Hadoop e Spark



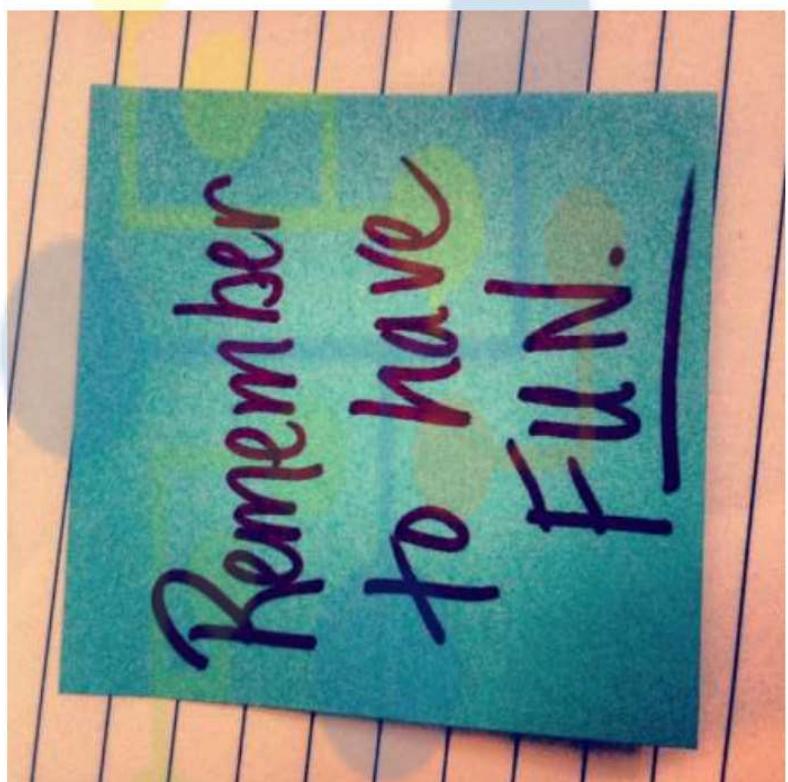
Data Science Academy phelipe.ufssemprboni@outlook.com 5c8a62005e4cded1aeh8b45a3



Engenharia de Dados com Hadoop e Spark



Data Science Academy phelipe.ufsempreboni@outlook.com 5c8a62005e4cded1aeh8b45a3

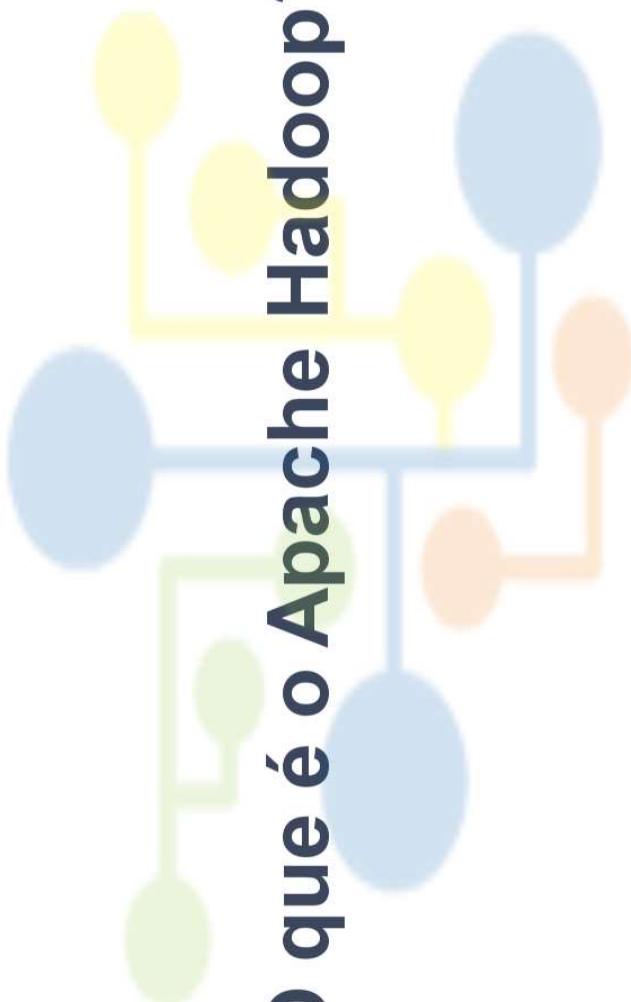




Data Science
Academy

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

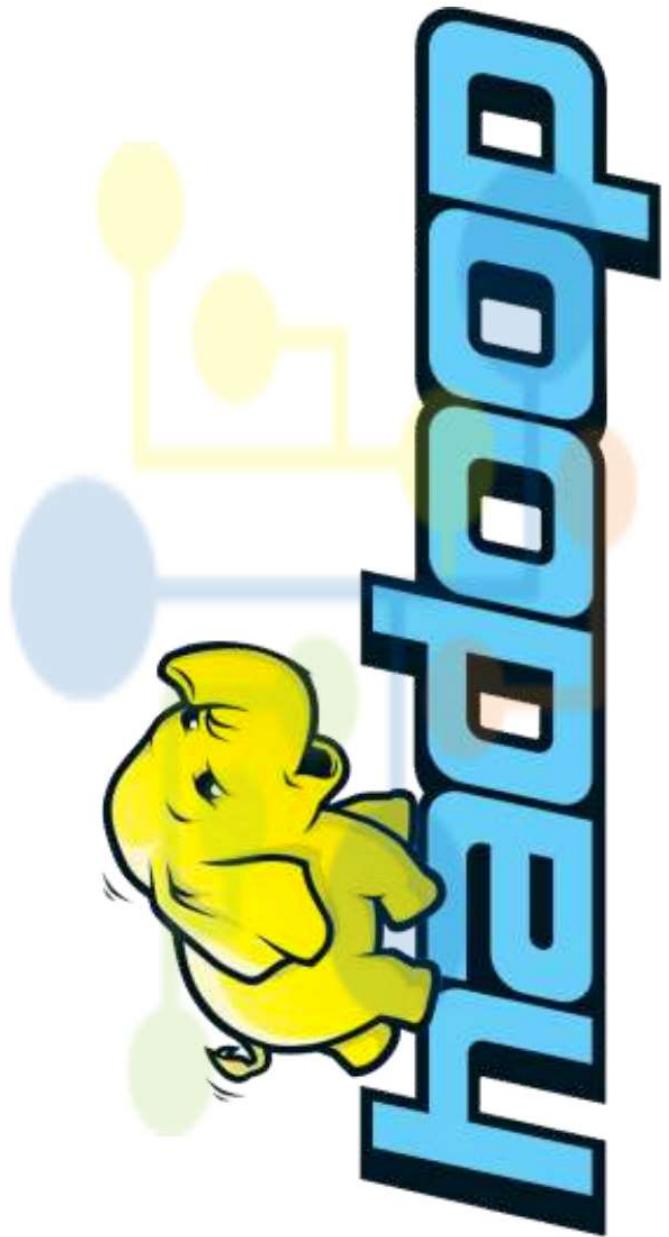
O que é o Apache Hadoop?



O que é o Apache Hadoop?



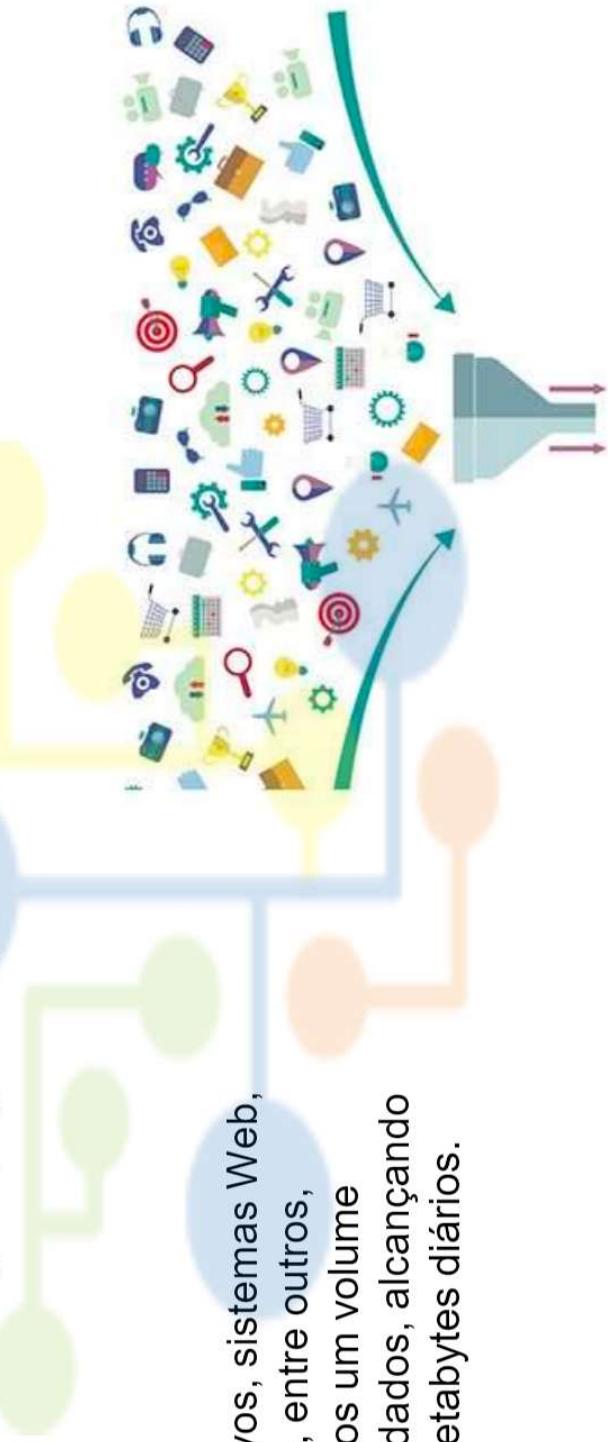
Data Science Academy nheine.ufsemprebon@gmail.com 5c8a62005e4cd1acb8b45a3



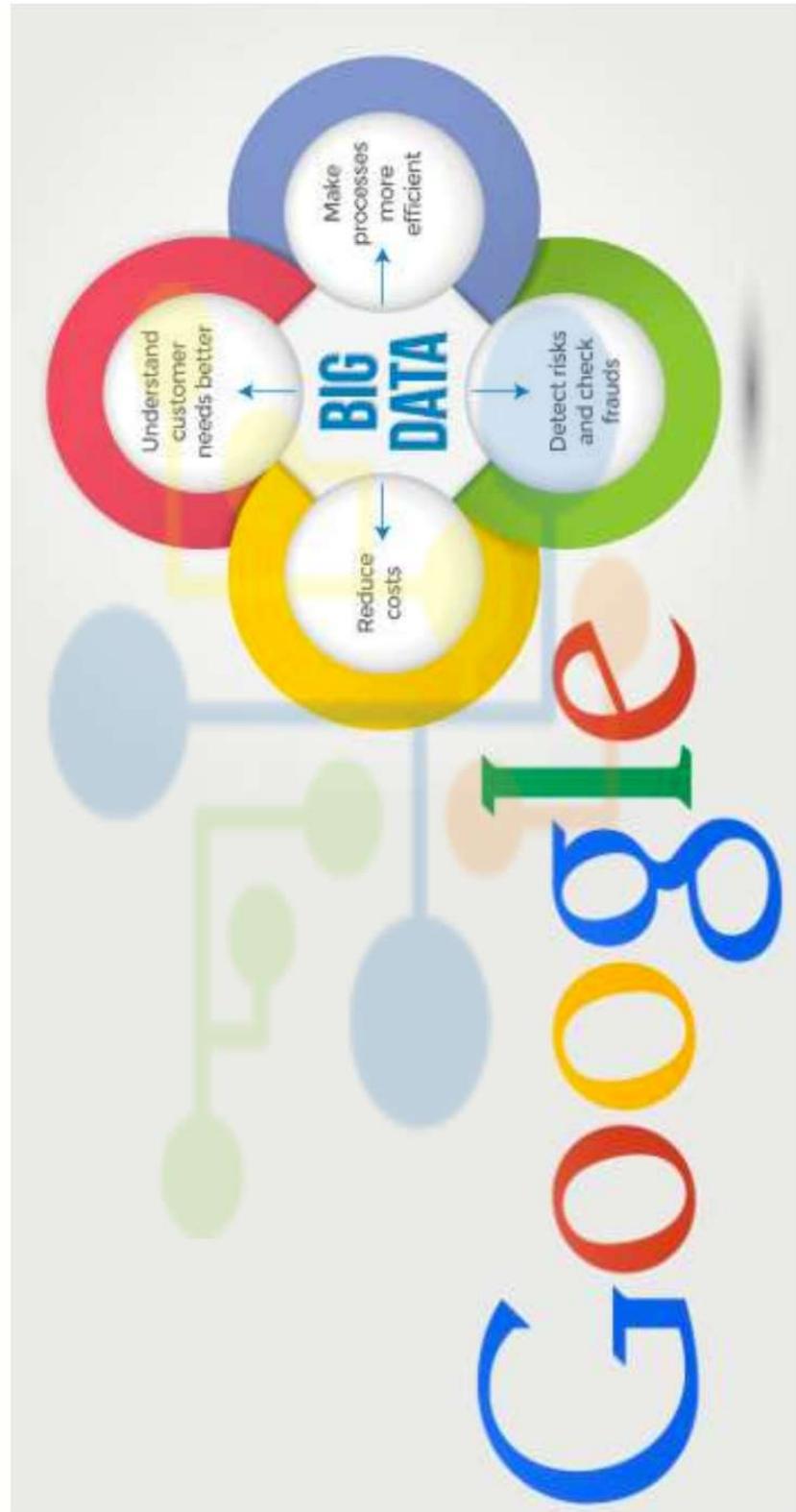
O que é o Apache Hadoop?

Data Science Academy nphelipe.ufsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

Um dos grandes desafios computacionais da atualidade é armazenar, manipular e analisar, de forma inteligente, a grande quantidade de dados existente.



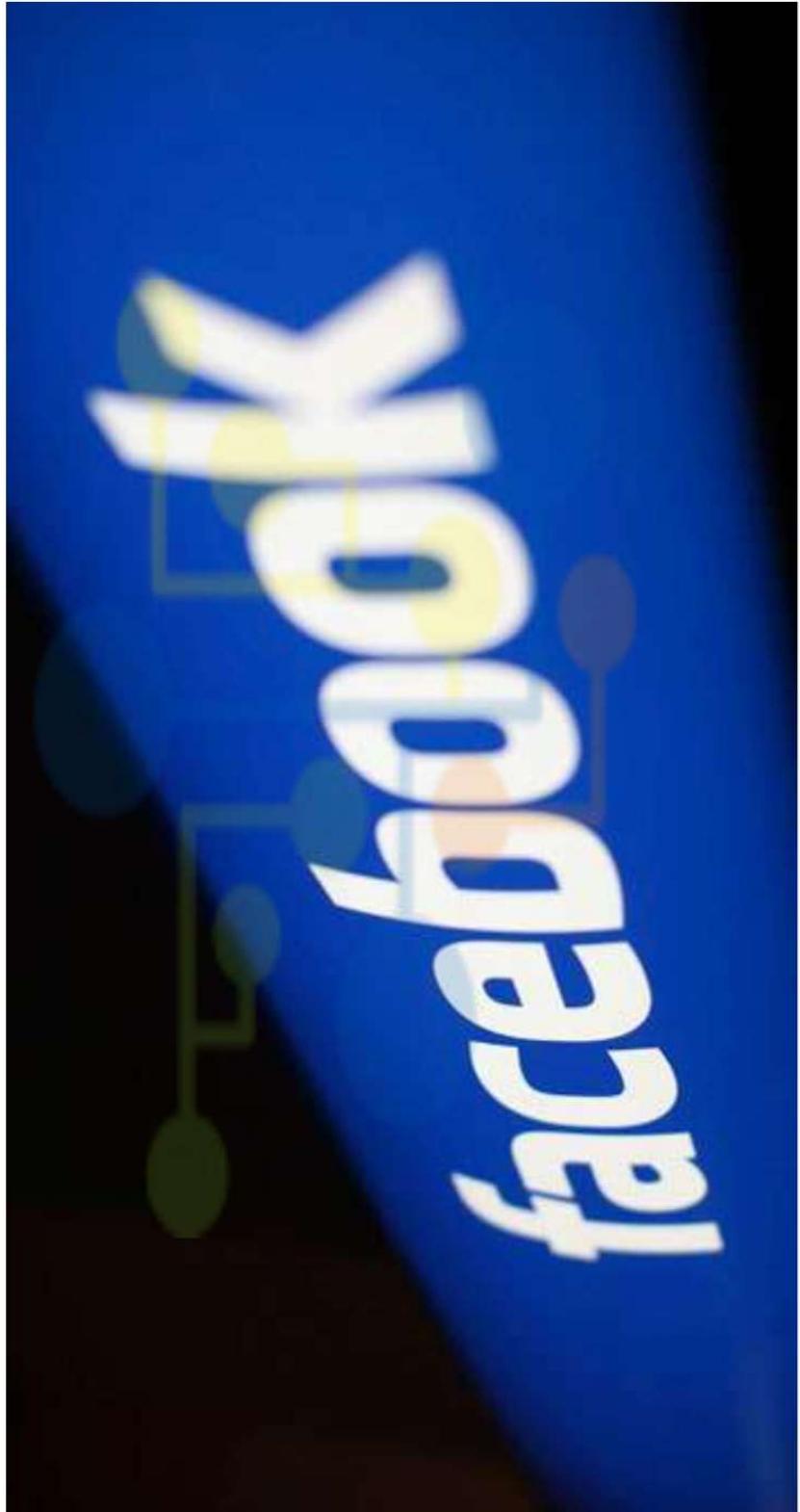
Sistemas corporativos, sistemas Web, mídias sociais, entre outros, produzem juntos um volume impressionante de dados, alcançando a dimensão de petabytes diários.



O que é o Apache Hadoop?

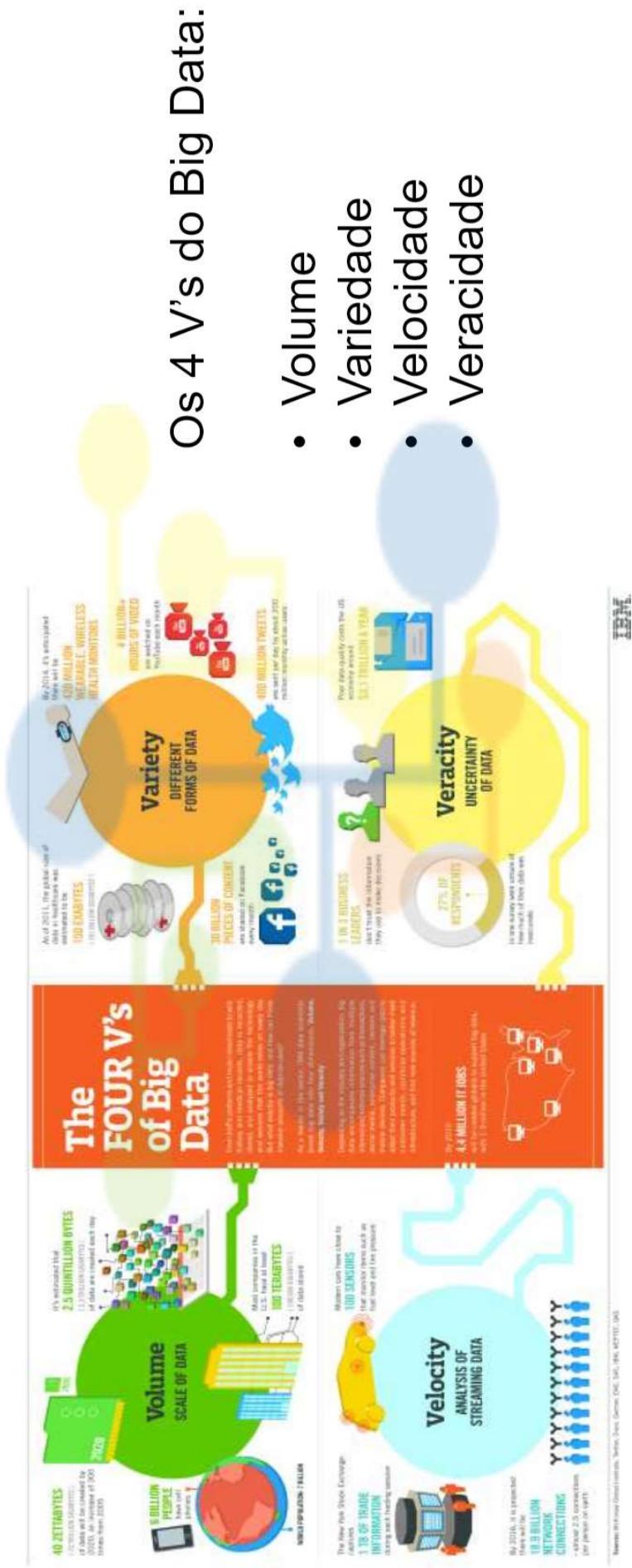


Data Science Academy Data.Science.Academy@helipr.ufsc.br emprebon@outlook.com 5c8a62005e4cd1acb8b45a3



O que é o Apache Hadoop?

Data Science Academy Data.Science.Academy@helipe.ufssemprbon@outlook.com 5c8a6/2005e4cd1acd8b45a3



O que é o Apache Hadoop?



Data Science Academy nhelene.ufssemprobony@outlook.com 5c8a62005e4cd1acb8b45a3



O que é o Apache Hadoop?



Data Science Academy nhelipe.ufsemprebon@gmail.com 5c8a62005e4cd1acb8b45a3

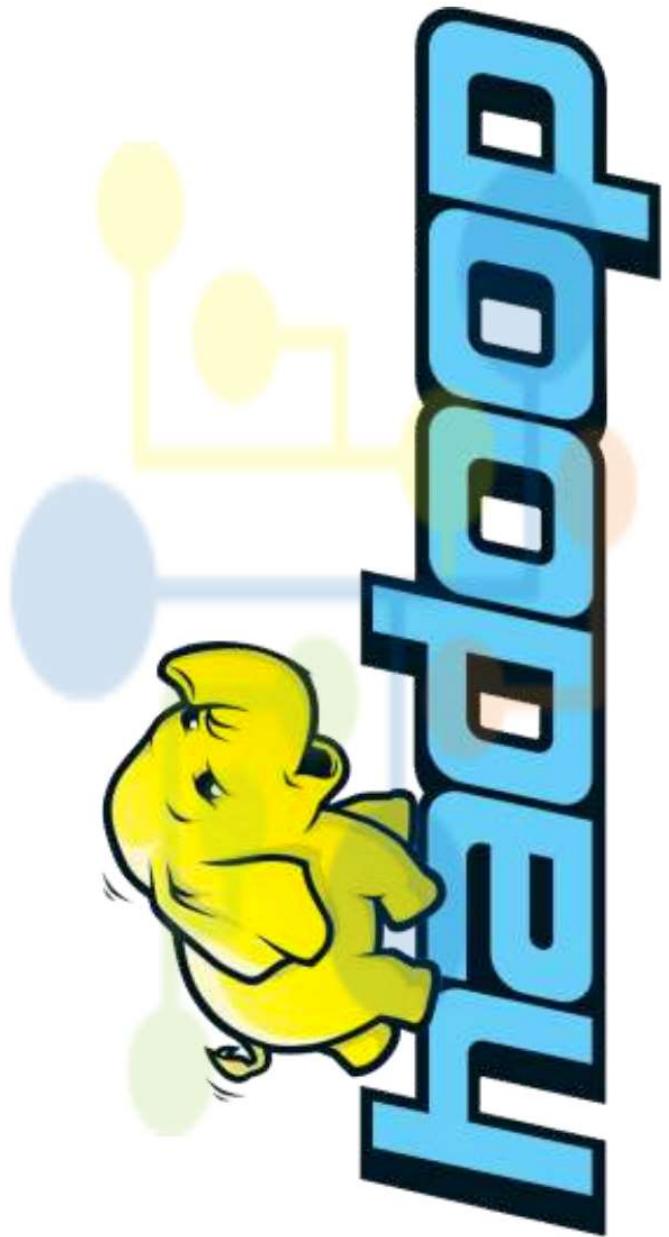


Computação Paralela

O que é o Apache Hadoop?



Data Science Academy nheine.ufsemprebon@gmail.com 5c8a62005e4cd1acb8b45a3



Uma Breve História do Apache Hadoop

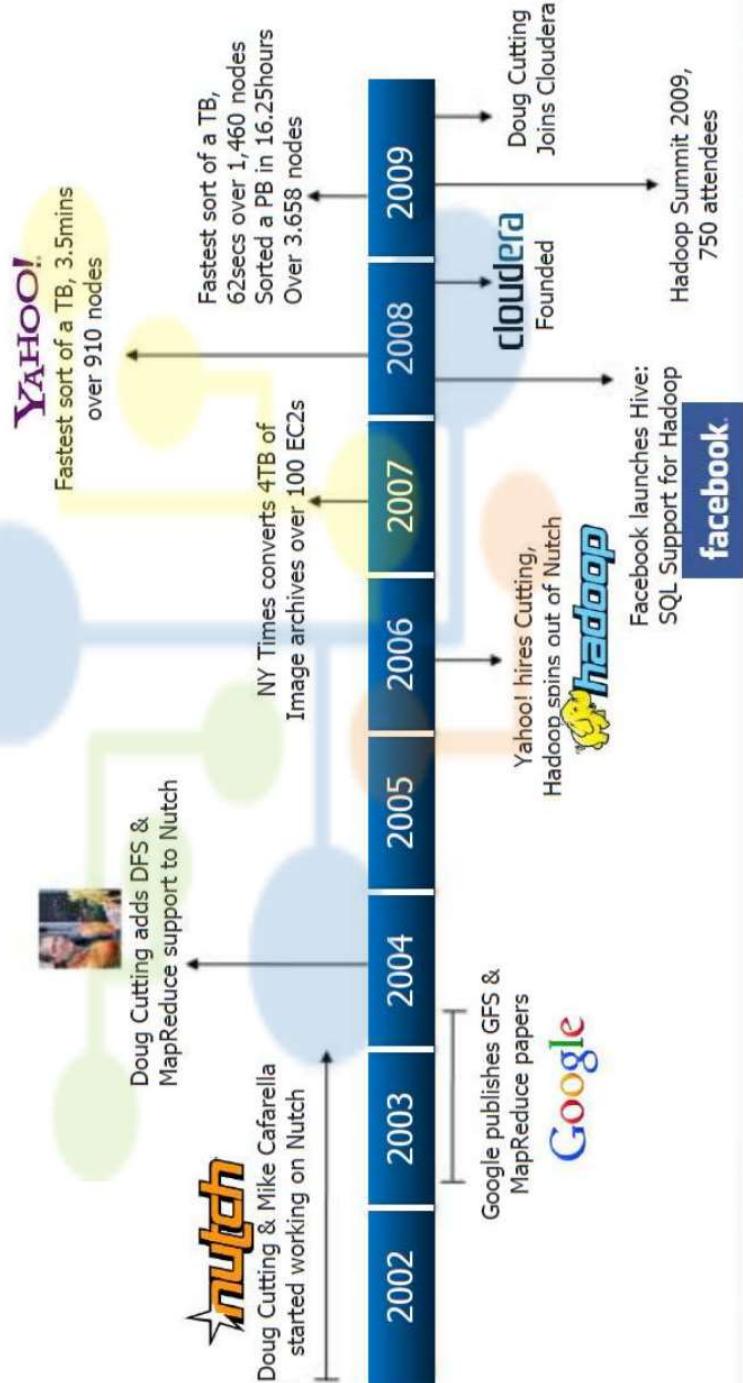


Data Science Academy

Academy

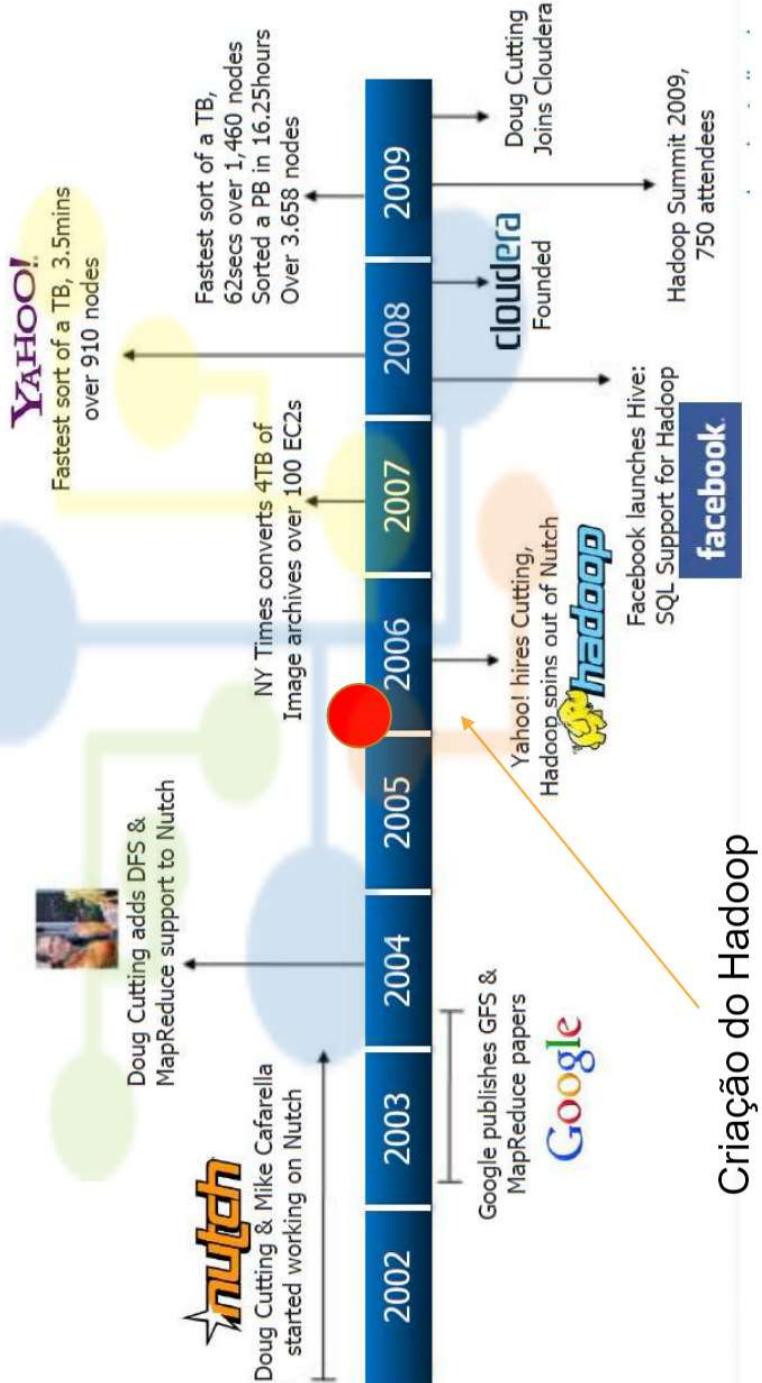
Academy

Hadoop History



Uma Breve História do Apache Hadoop

Hadoop History

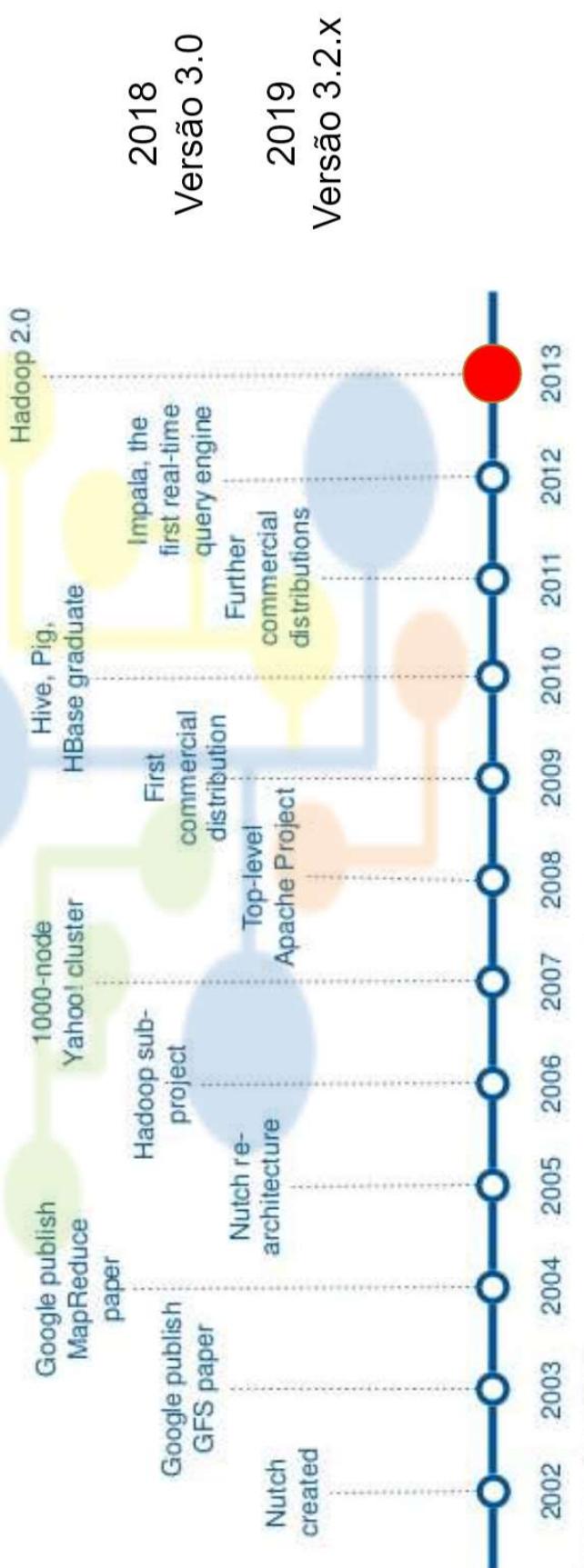


Uma Breve História do Apache Hadoop



Data Science Academy pheilipe.jttempreboni@outlook.com 5c8362005e4cded1aach8b45a3

A Brief History of Hadoop



O que é o Hadoop?

Data Science Academy Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cde1acb8b45a3



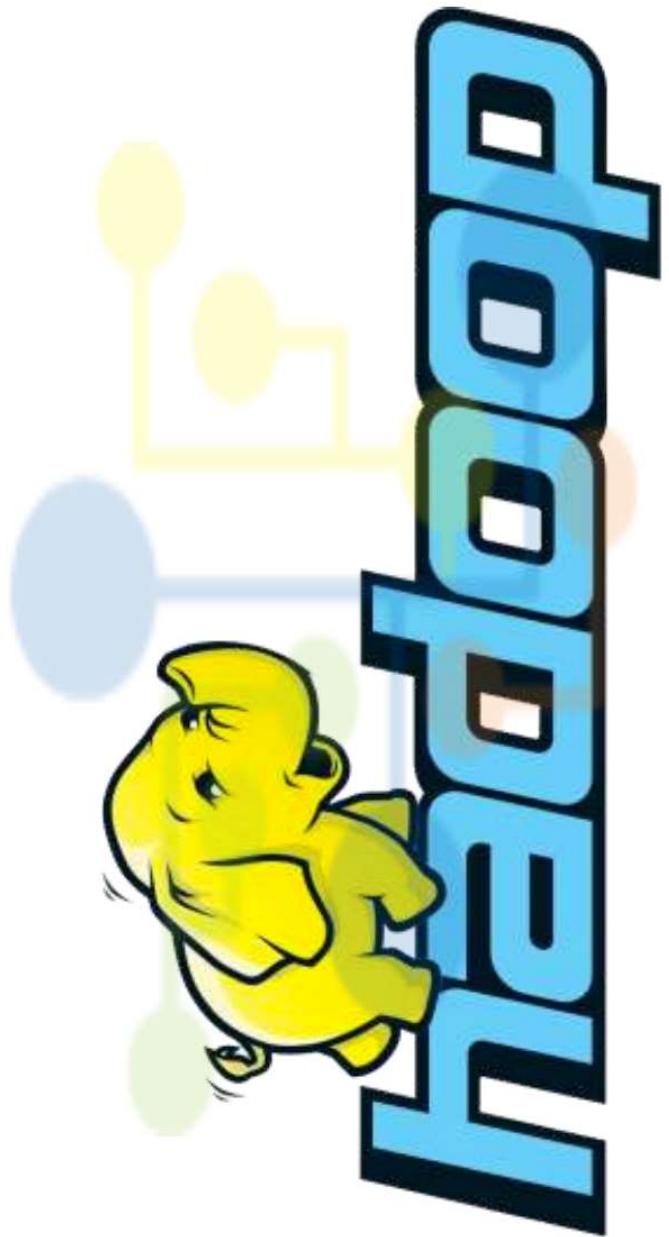


Quais os benefícios para as Empresas ao utilizar o Hadoop?



Benefícios do Hadoop

Data Science Academy, neliipe,utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3
Data Science Academy



Benefícios do Hadoop

Data Science Academy, phelipe.utsempreboni@outlook.com 5c8a62005e4cded1acb8b45a3



Open Source



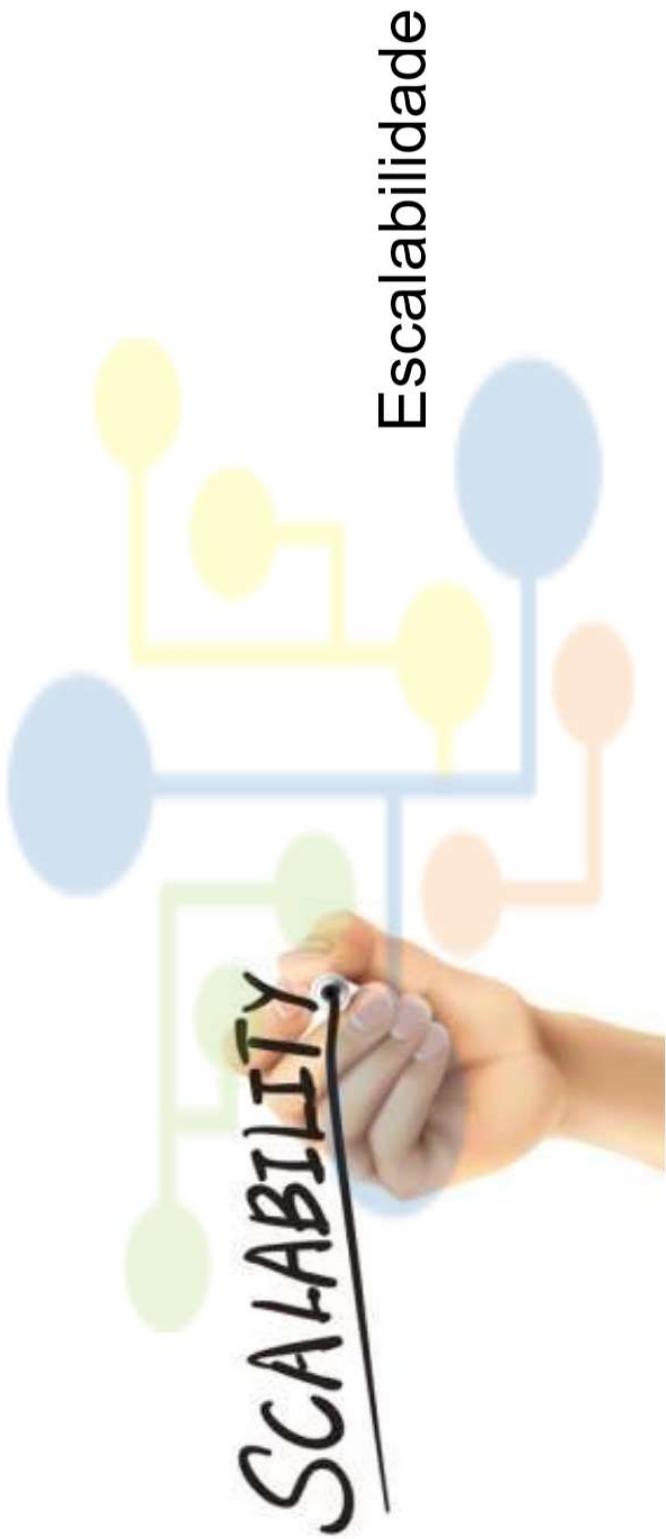
Benefícios do Hadoop

Data Science Academy, neliipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3



Benefícios do Hadoop

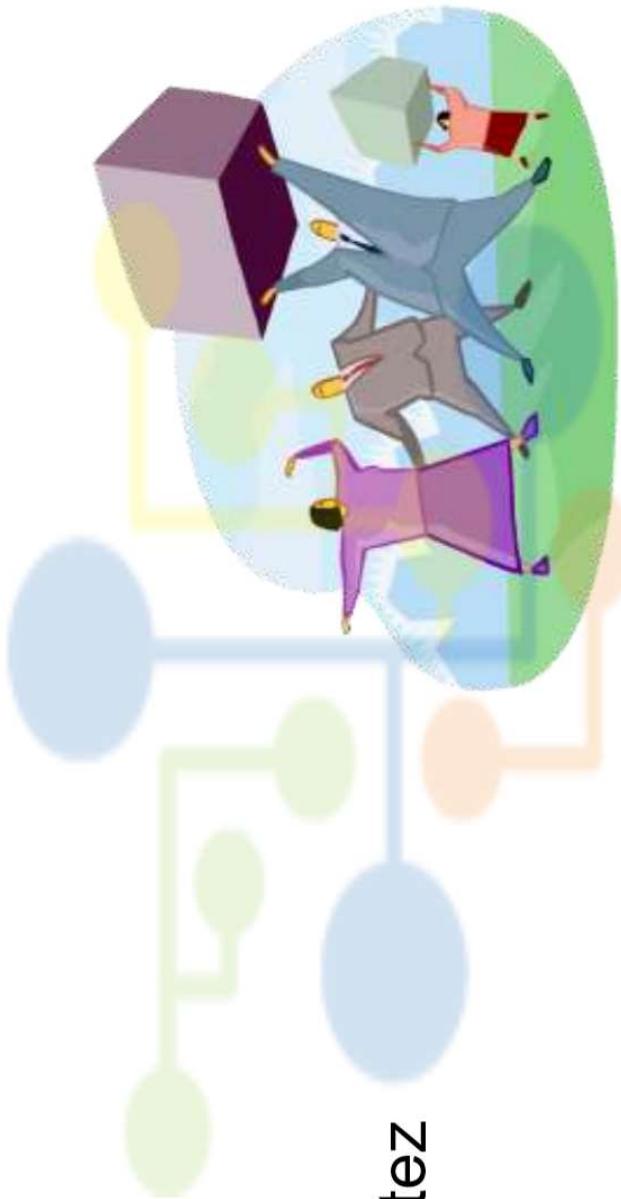
Data Science Academy Data Science Academy, neliipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3



Escalabilidade

Benefícios do Hadoop

Data Science Academy, neliipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

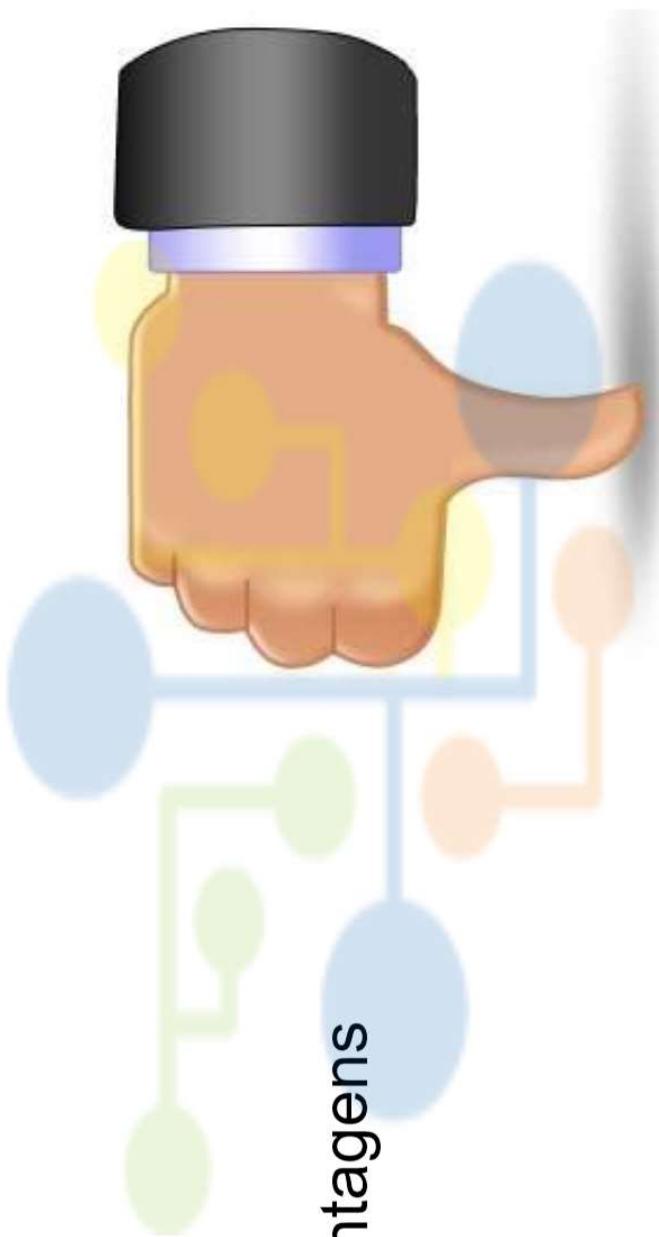


Robustez

Desvantagens do Hadoop



Data Science Academy, pheipe, ursempreboni@outlook.com 5c8a62005e4cd1acb8b45a3



Desvantagens

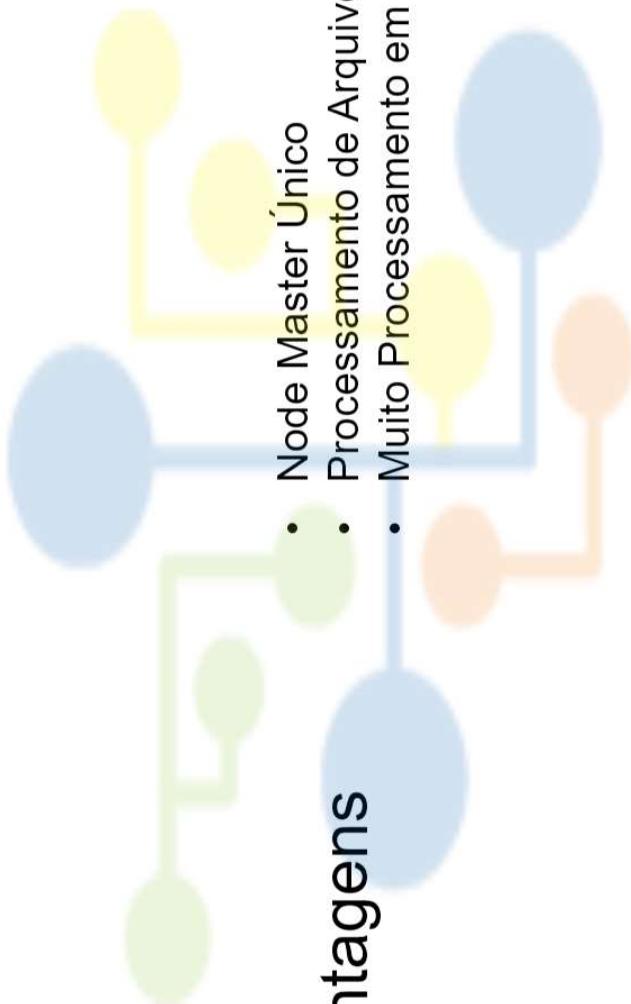
Desvantagens do Hadoop



Data Science Academy, pheilipe.ufsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

Desvantagens

- Node Master Único
- Processamento de Arquivos Pequenos
- Muito Processamento em Poucos Dados





Data Science
Academy

Data Science Academy phelipe.utsemprboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

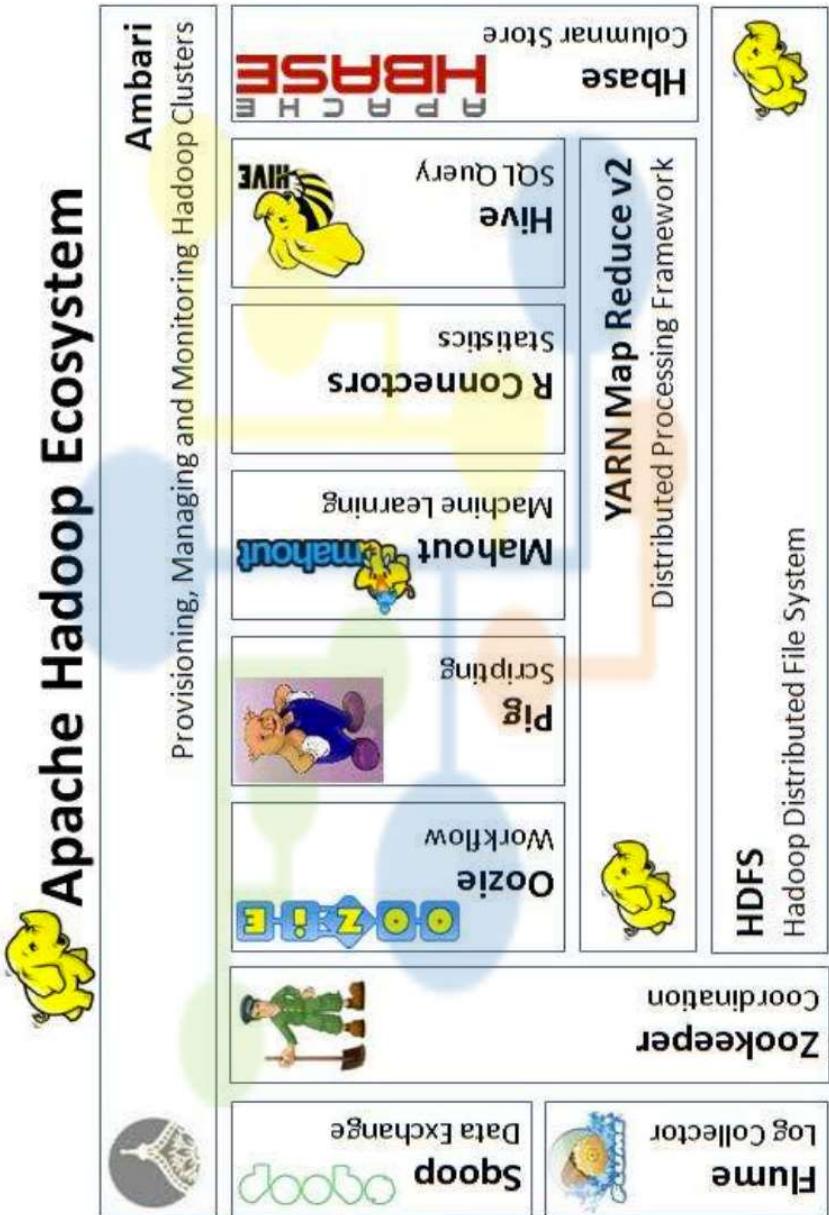


EcoSistema Hadoop

Ecosistema Hadoop

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

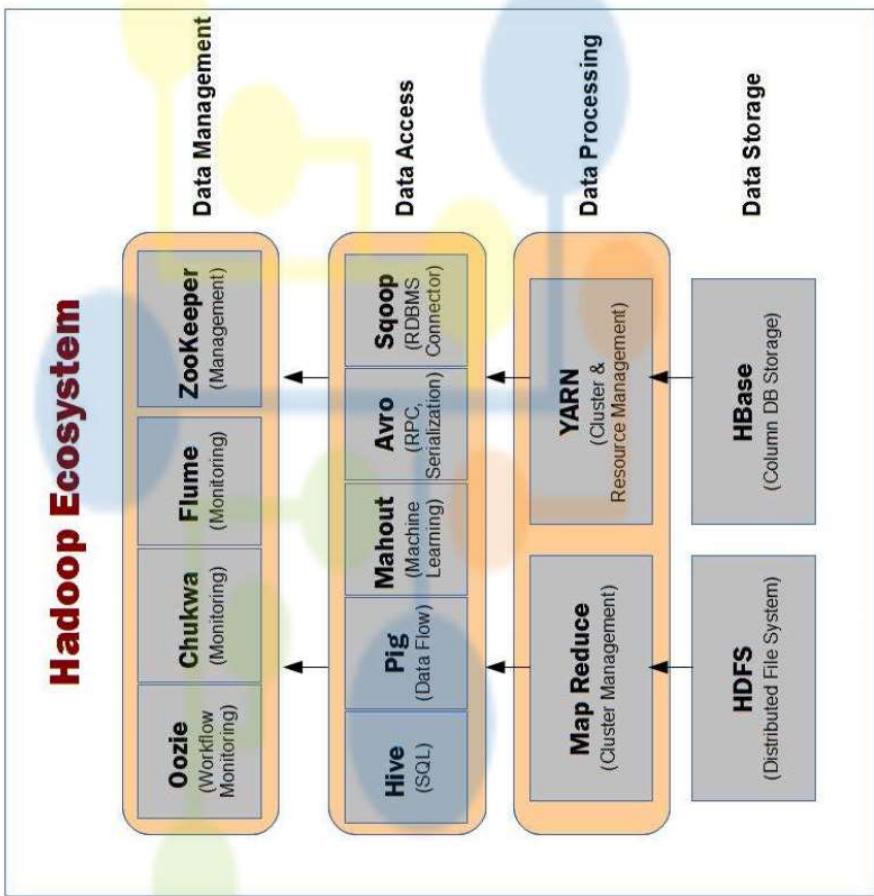
Apache Hadoop Ecosystem



Ecosistema Hadoop

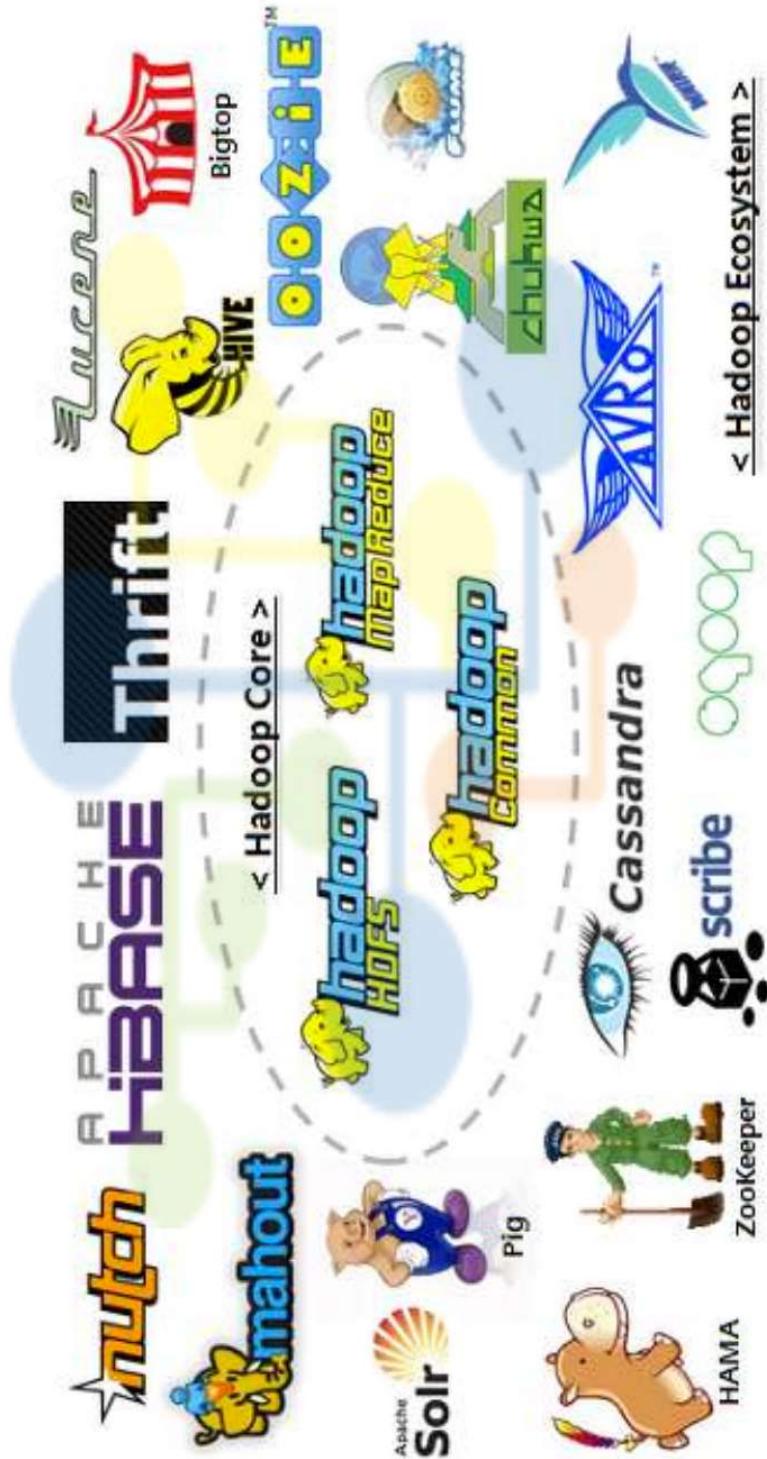
Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

Hadoop Ecosystem



Ecossistema Hadoop

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3



Ecossistema Hadoop

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cde1acb8b45a3

Projetos Principais do Ecossistema Hadoop



Ecosistema Hadoop

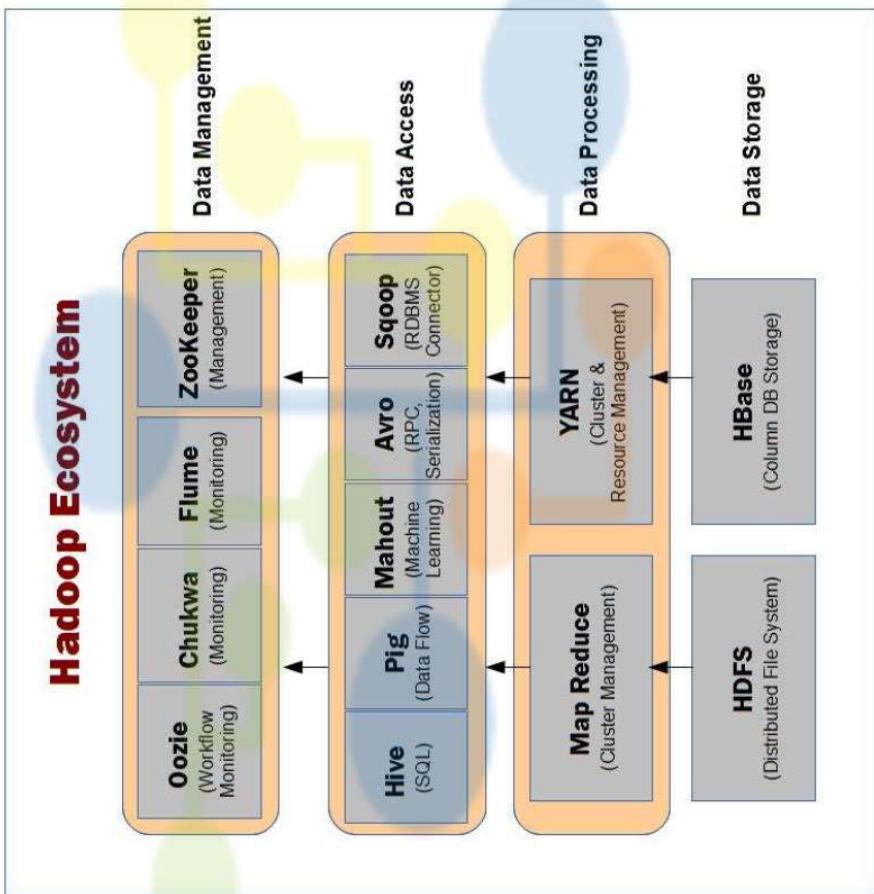
Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

Outros Projetos Importantes



Ecosistema Hadoop

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3





Data Science
Academy

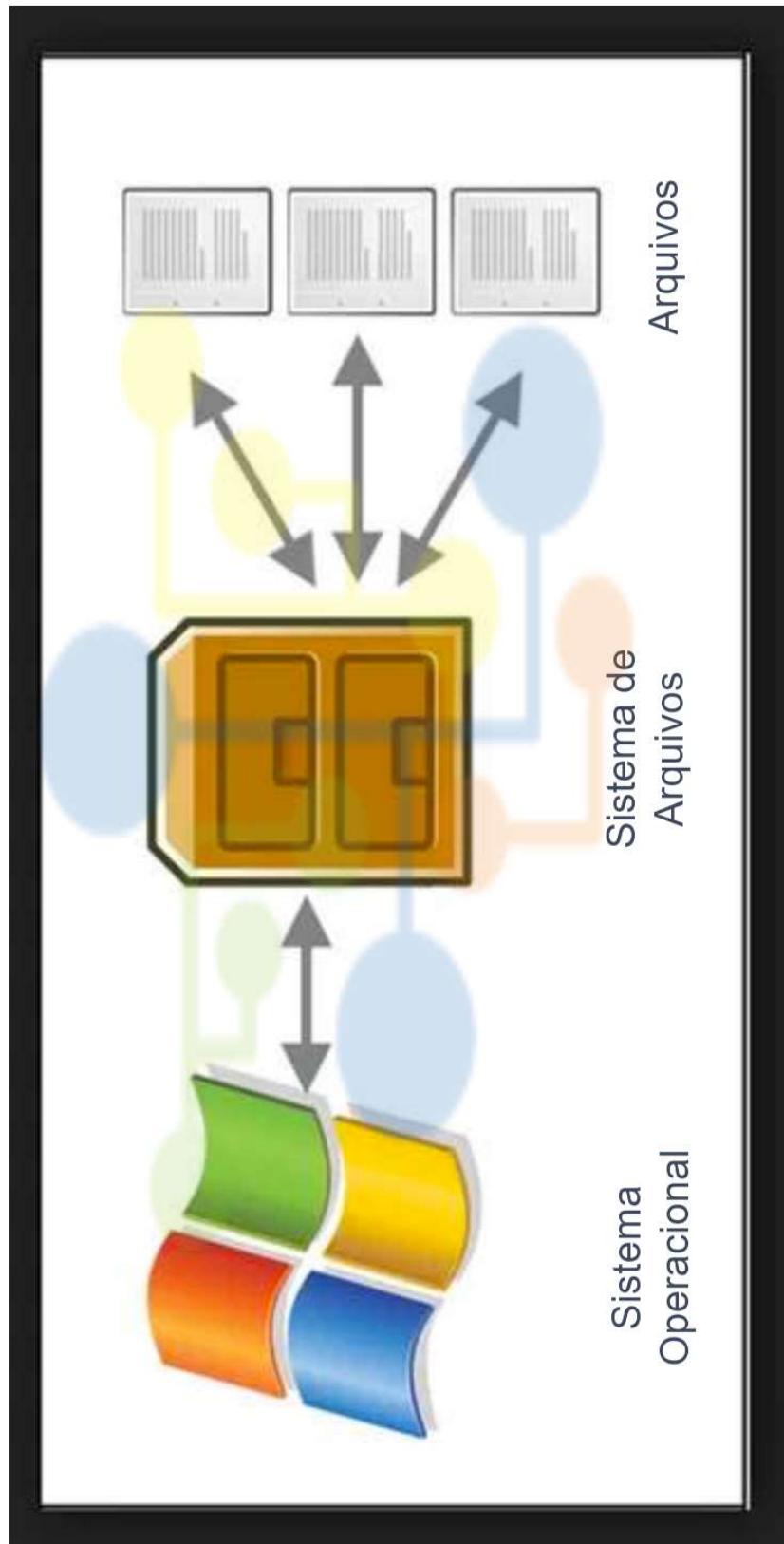
Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

HDFS (Hadoop Distributed File System)

Conceito e Importância

HDFS – Conceito e Importância

Data Science Academy phelipe.ufsemprebon@outlook.com 5c8a62005e4cd1acb8b45a3



HDFS – Conceito e Importância

Data Science Academy Data Science Academy phelege.ufsemprebon@outlook.com 5c8a6/2005e4cd1acb8b45a3

Os tipos de Sistemas de Arquivos são:

Tipo	Descrição
ext2	Sistema de arquivos padrão do Linux
ext3	Sistema de arquivos ext2 melhorado
reiserfs	Sistema de arquivos do tipo Journaling
msdos	Sistema de arquivos FAT da Microsoft DOS
vfat	Sistema de arquivos FAT-32 do Microsoft Windows
iso9660	Sistema de arquivos do CD-ROM
nfs	Network File System. Usado para montar dispositivos em computadores remotos.
swap	Sistema de arquivos de troca utilizando para memória virtual.
proc	Uma janela especial dentro do Kernel do Linux. Utilizada pelos usuários, programas e utilitários para escrever ou ler parâmetros do Kernel. Geralmente montado no diretório /proc.

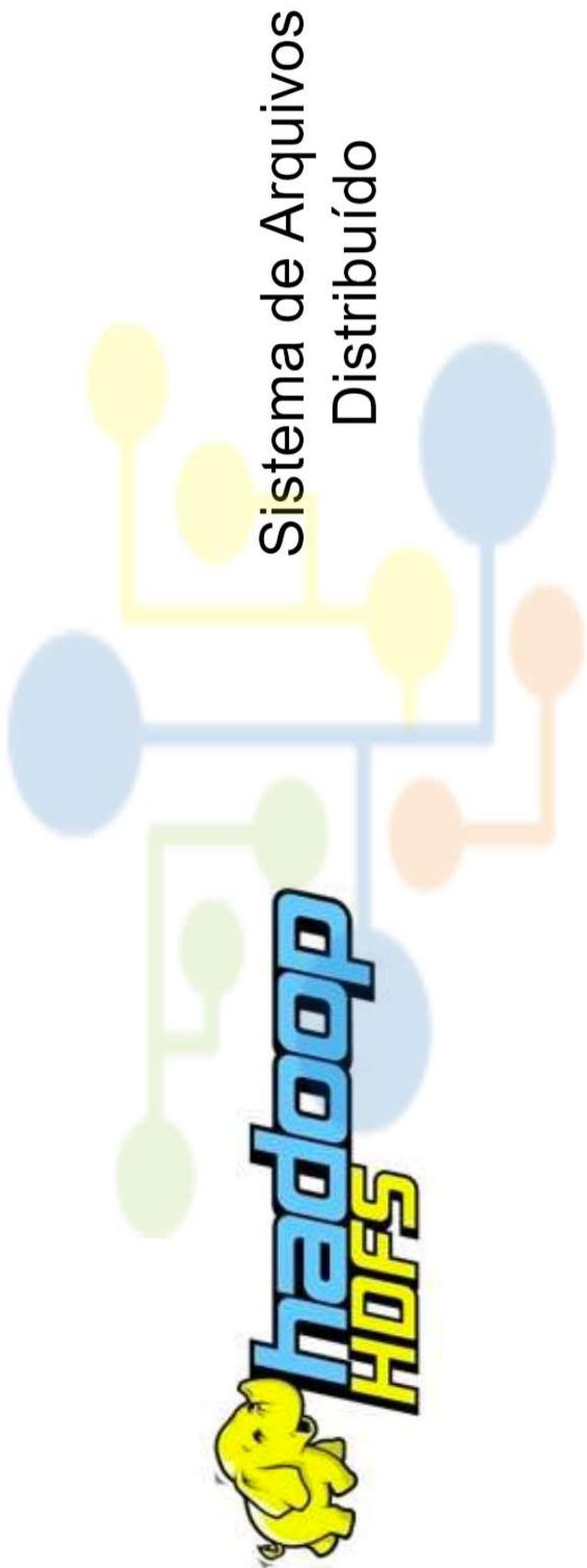
HDFS – Conceito e Importância

Data Science Academy phelege.ufsemprebon@outlook.com 5c8a62005e4cd1acb8b45a3



HDFS – Conceito e Importância

Data Science Academy  Data Science Academy phelipe.ufsemprebon@outlook.com 5c8a62005e4cd1acb8b45a3



HDFS – Conceito e Importância

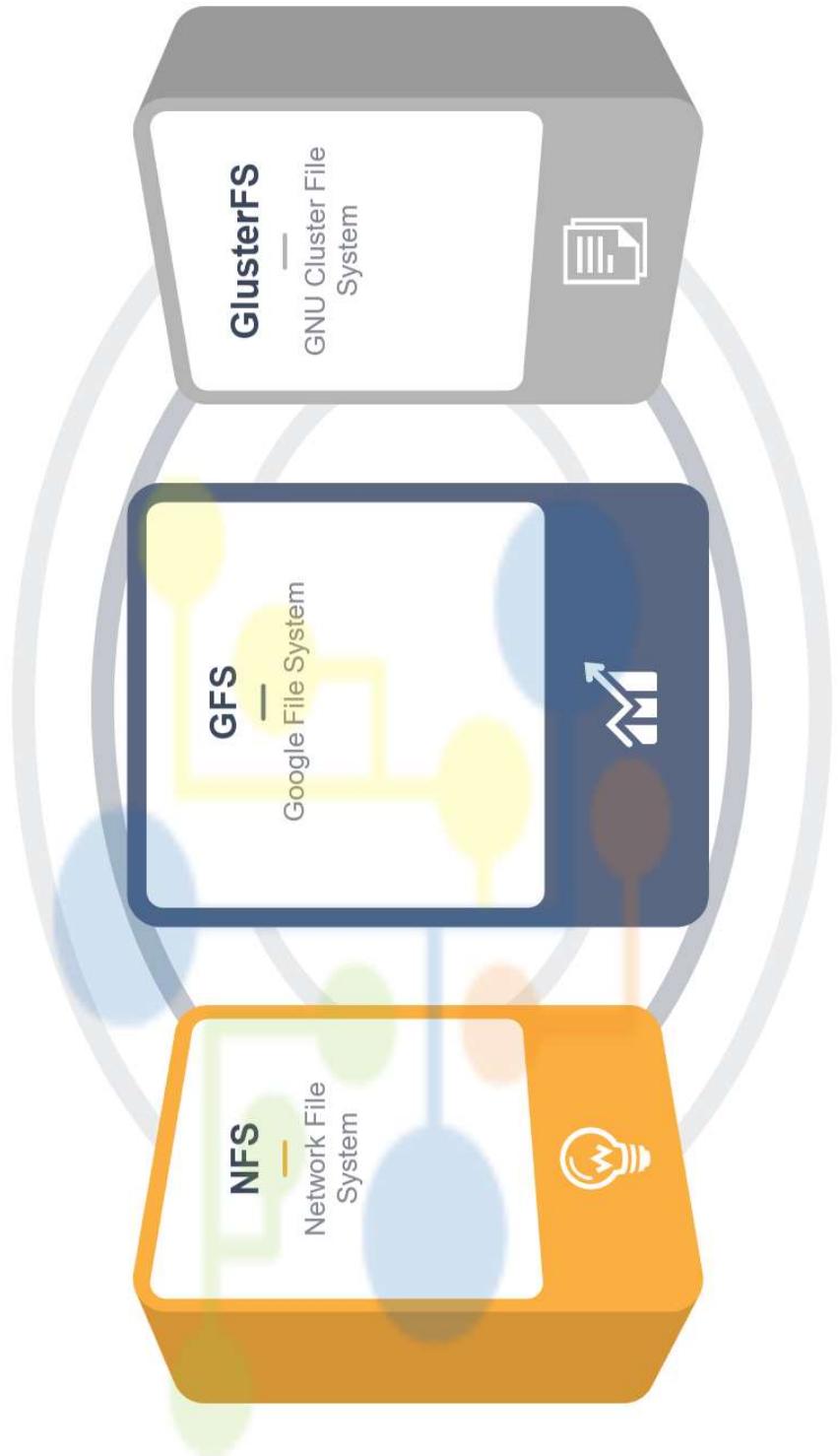
Data Science Academy phelipe.ufsemprebon@outlook.com 5c8a62005e4cd1acb8b45a3

- Tolerância a Falhas
- Integridade
- Segurança
- Desempenho
- Consistência



HDFS – Conceito e Importância

Data Science Academy pheipe.ufsemprebon@outlook.com 5c8a62005e4cd1acb8b45a3

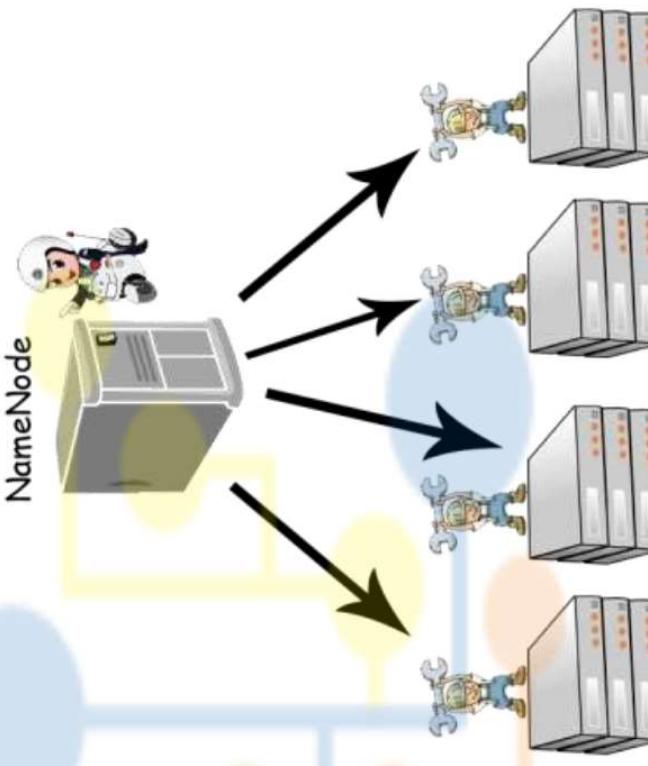


Outros Sistemas
de Arquivos
Distribuídos

HDFS – Conceito e Importância

Data Science Academy phelipe.ufssemprbon@outlook.com 5c8a62005e4cd1acb8b45a3

Hadoop
Distributed File
System

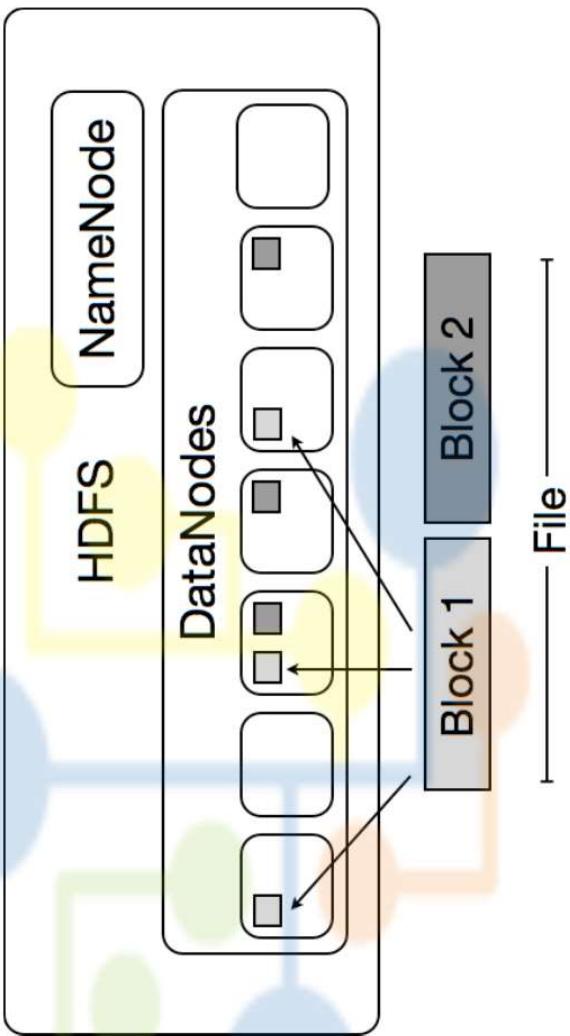


Data Nodes (Commodity Hardware)

HDFS – Conceito e Importância

Data Science Academy phelege.ufsemprebon@outlook.com 5c8a62005e4cd1acb8b45a3

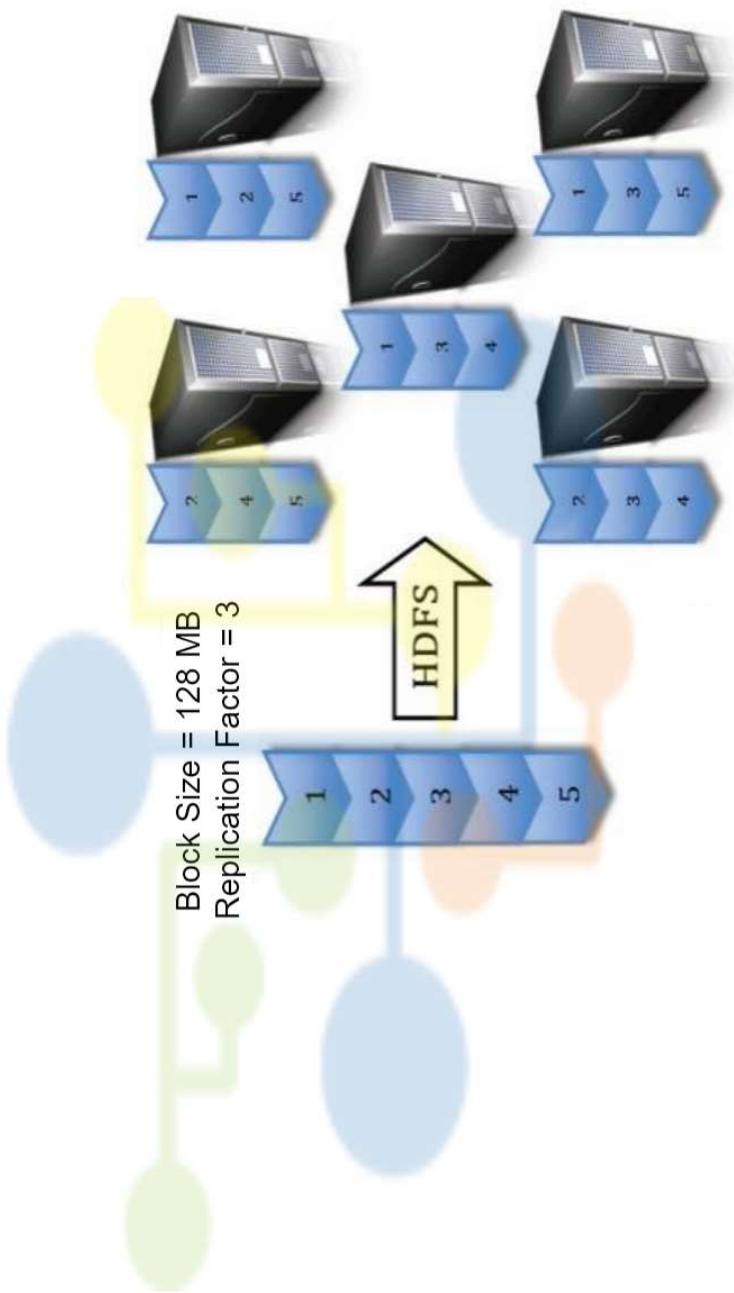
Hadoop
Distributed File
System



HDFS – Conceito e Importância

Data Science Academy pheilipe.ufsemprebon@outlook.com 5c8a6/2005e4cd1acb8b45a3

Hadoop
Distributed File
System



HDFS – Conceito e Importância



Data Science Academy phelipe.ufsmprebon@outlook.com 5c8a62005e4cd1acb8b45a3

O HDFS foi criado para resolver “Big Problems” e por isso seu funcionamento e arquitetura são próprios para se trabalhar com grandes arquivos de dados e distribuir esses arquivos em blocos ao longo de um cluster de computadores, para que possam ser processados em paralelo.



Data Science
Academy

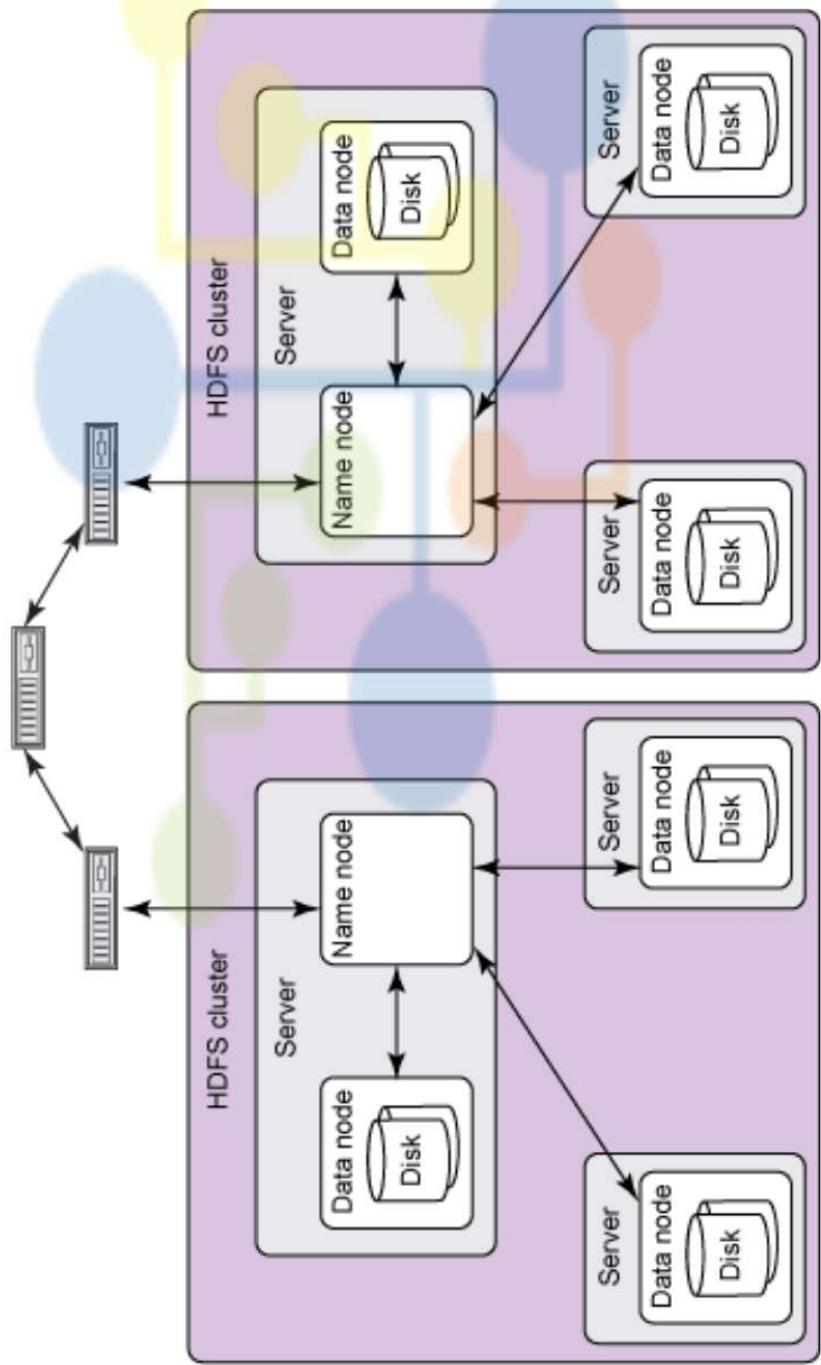
Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

HDFS (Hadoop Distributed File System) Arquitetura

HDFS – Arquitetura

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

Arquitetura Master/Worker



HDFS – Arquitetura



Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

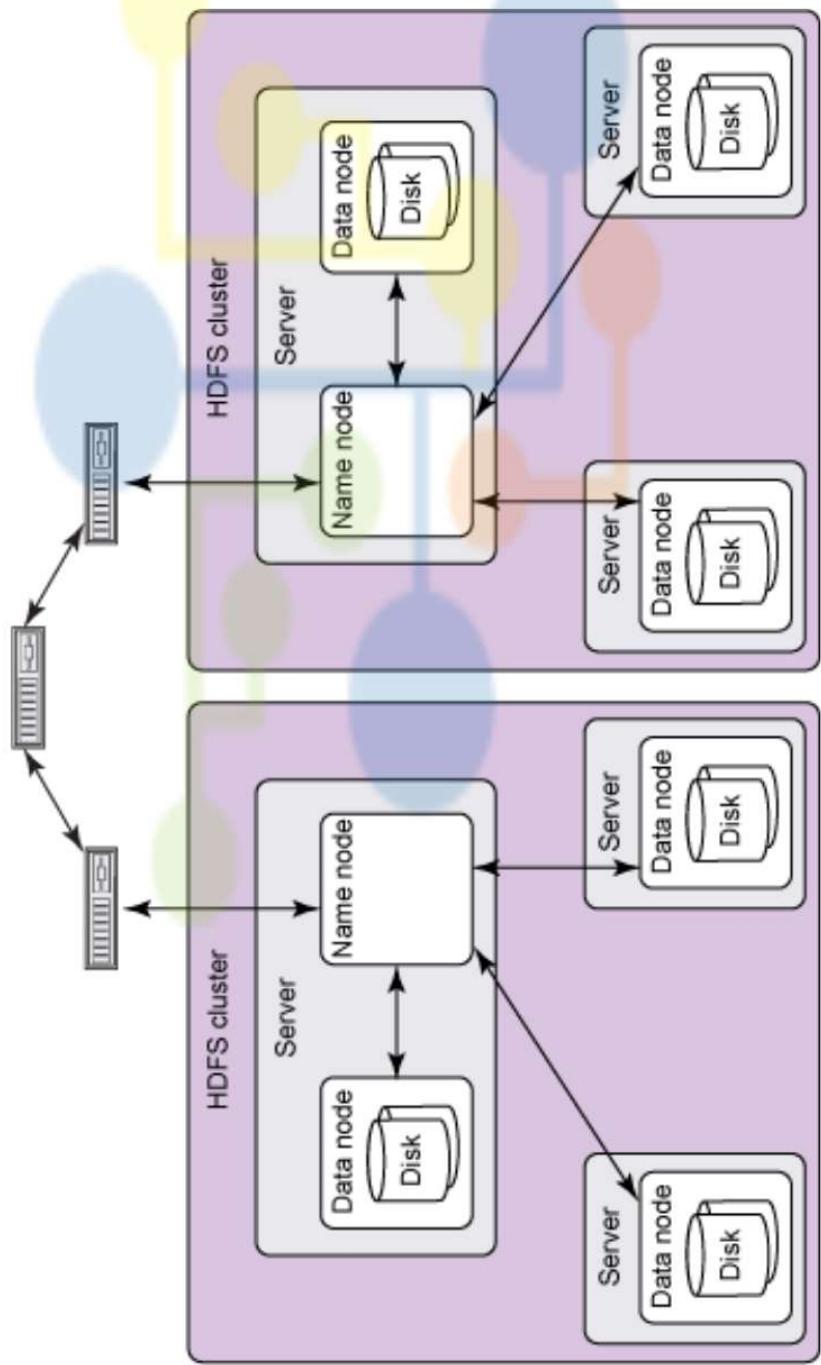


Arquitetura
Master/Worker

HDFS – Arquitetura

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

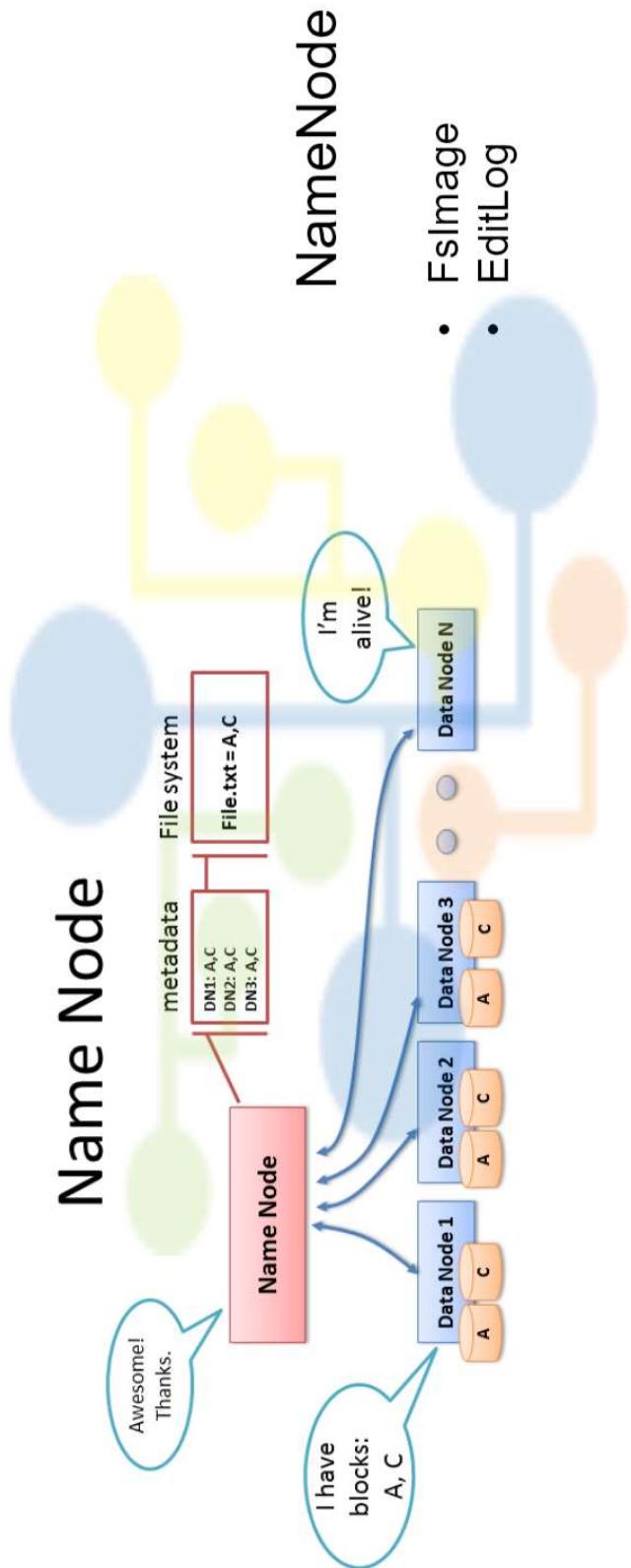
Arquitetura Master/Worker



HDFS – Arquitetura



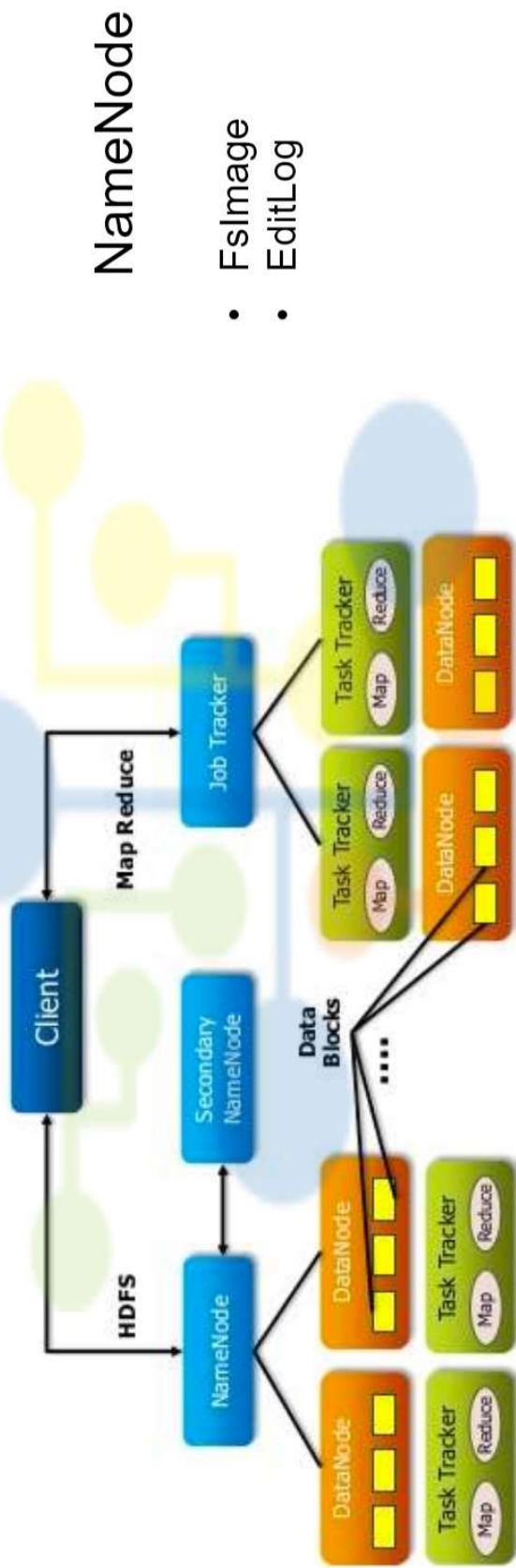
Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1ac8bb45a3



HDFS – Arquitetura

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

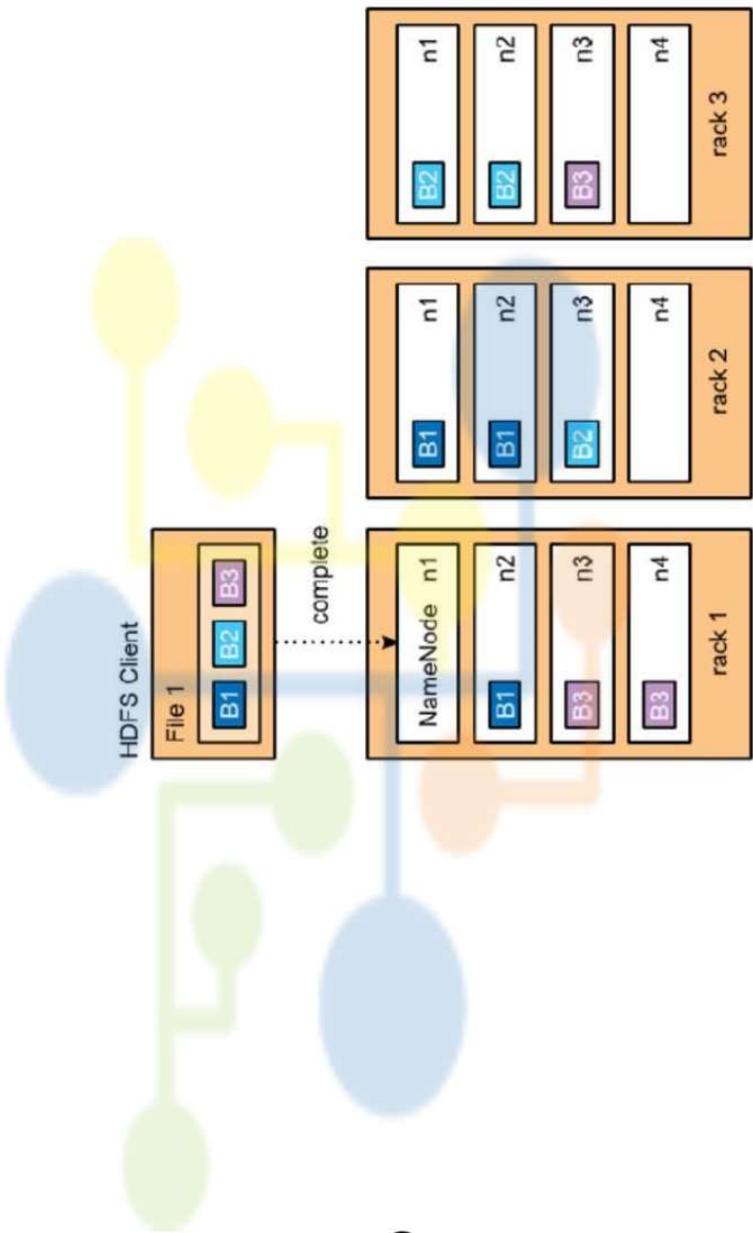
Data Science
Academy



HDFS – Arquitetura

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

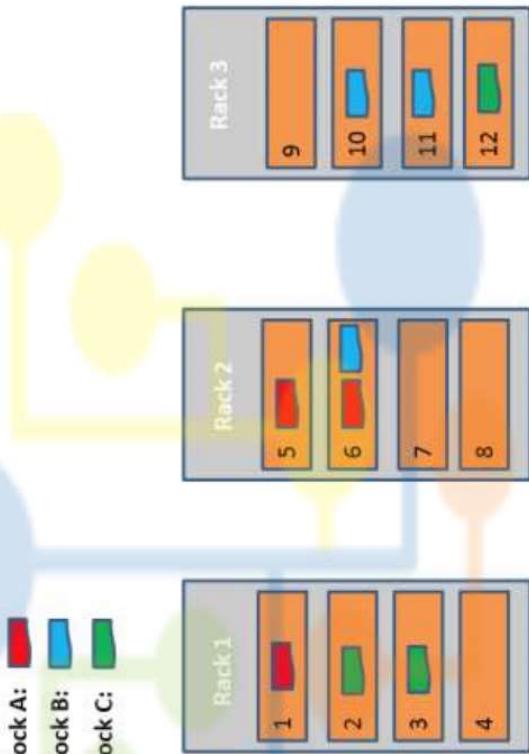
Replicação



HDFS – Arquitetura

Data Science Academy phelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3

Replicação





Data Science
Academy

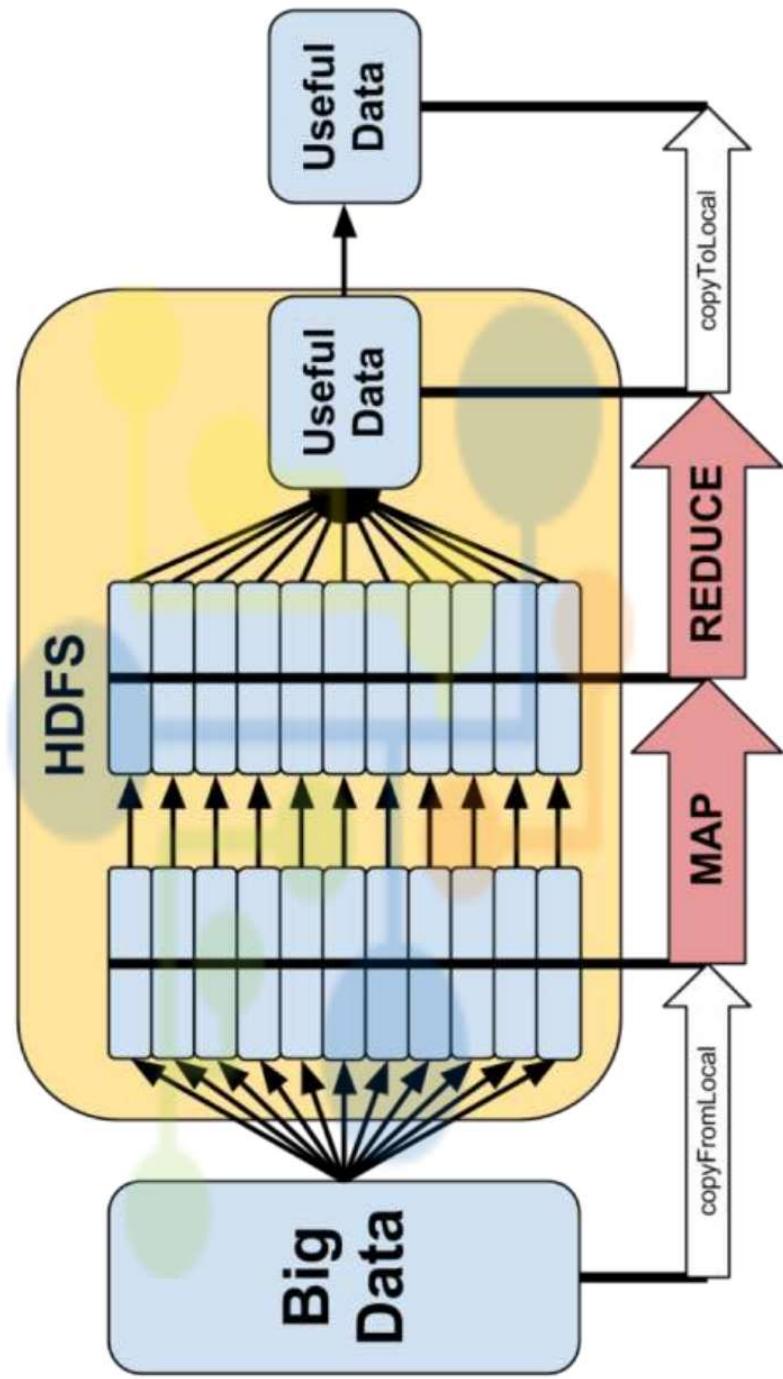
Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

Definindo MapReduce



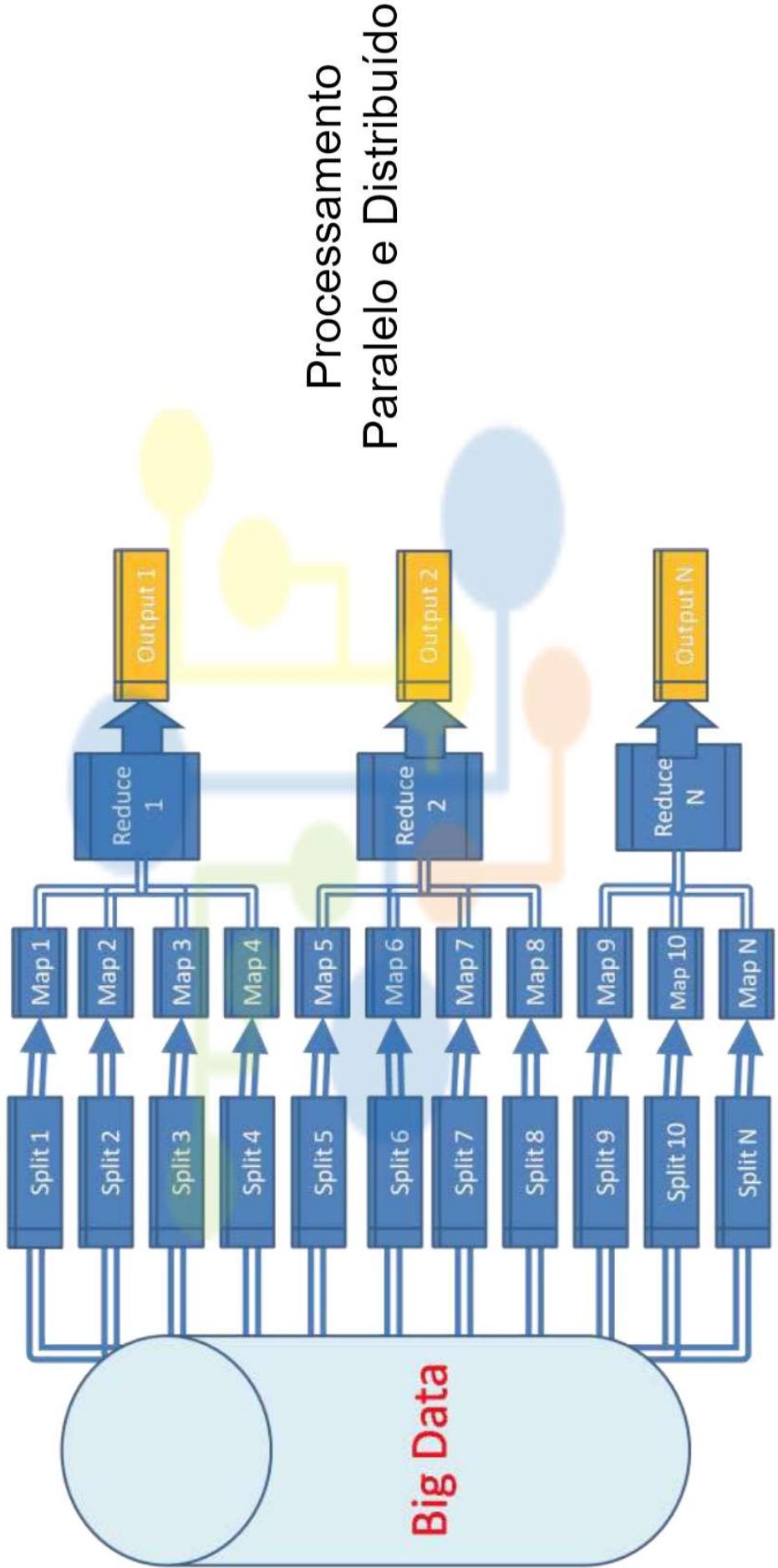
Definindo MapReduce

Data Science Academy Data Science Academy nhelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3



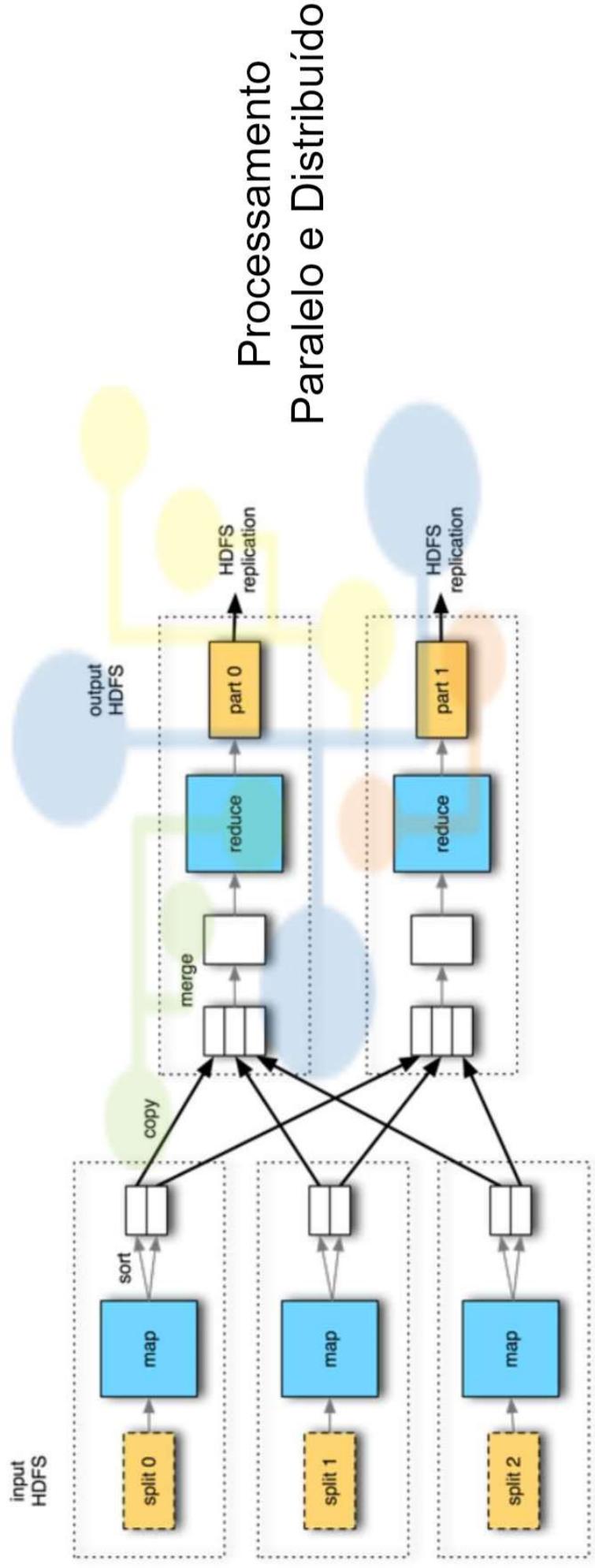
Definindo MapReduce

Data Science Academy Data Science Academy nhelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3



Definindo MapReduce

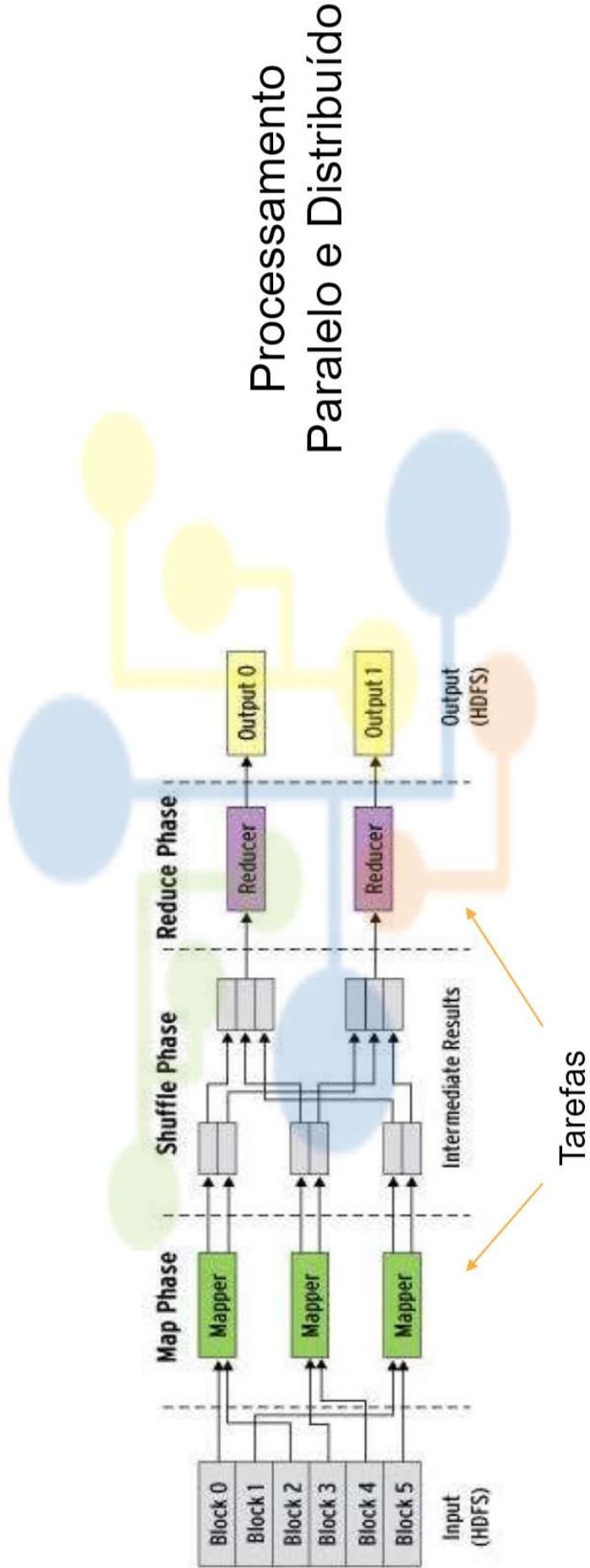
Data Science Academy Data Science Academy nhelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3



Processamento Paralelo e Distribuído

Definindo MapReduce

Data Science Academy Data Science Academy nhelipe.utsempreboni@outlook.com 5c8a62005e4cd1acb8b45a3





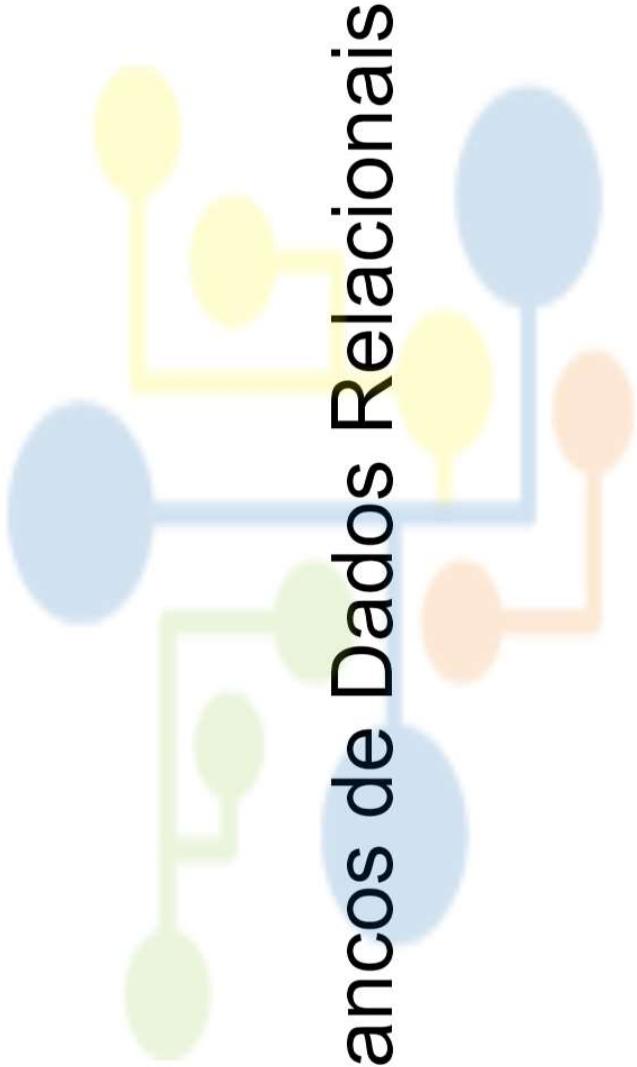
Data Science Academy phelipe.utsempreboni@outlook.com 5c8a6/2005e4cd1acb8b45a3

Hadoop x Bancos de Dados Relacionais

Hadoop x Bancos de Dados Relacionais



Data Science Academy philipe.utsempreboni@outlook.com 5c8a62005e4cd1acd8b45a3



Bancos de Dados Relacionais

Hadoop x Bancos de Dados Relacionais



Data Science Academy philipe.utsempreboni@outlook.com 5c8a62005e4cded1acb8b45a3

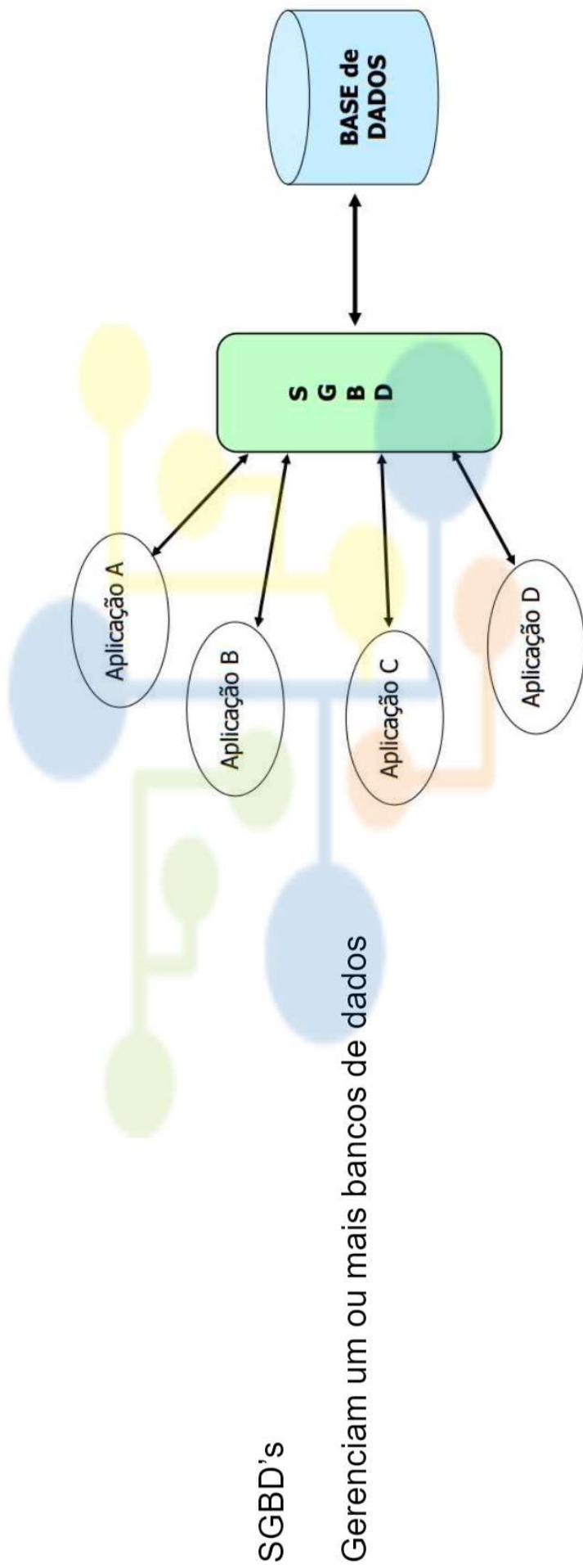


RDBMS
Relational Database
Management Systems

Hadoop x Bancos de Dados Relacionais



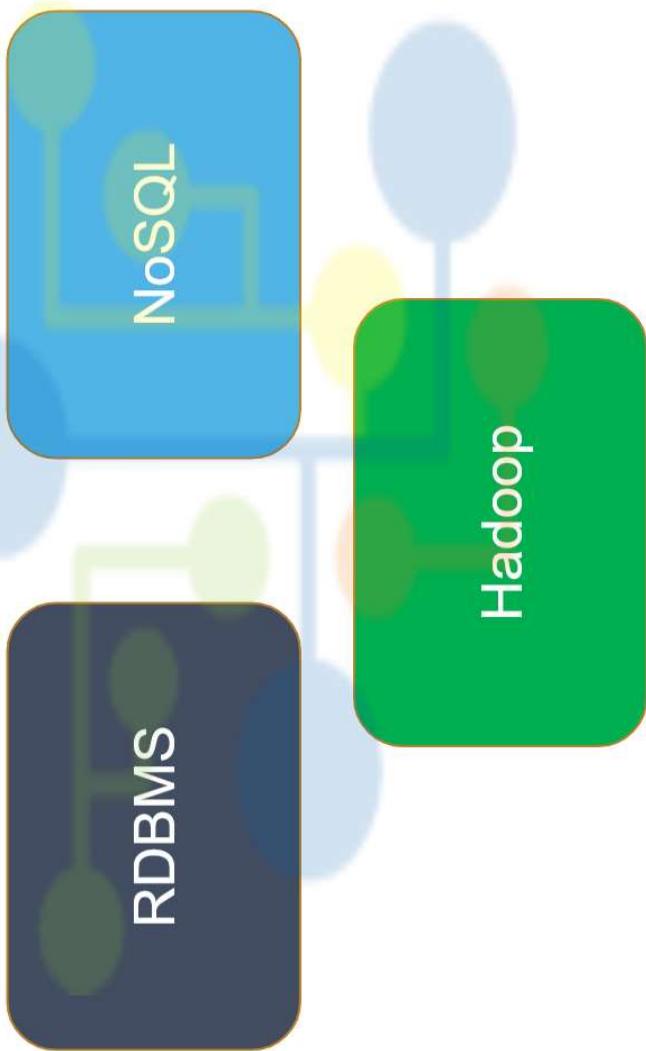
Data Science Academy philipe.utsempreboni@outlook.com 5c8a62005e4cd1acd8b45a3



Hadoop x Bancos de Dados Relacionais



Data Science Academy philipe.utsempreboni@outlook.com 5c8aa62005e4cd1acd8b45a3



Hadoop x Bancos de Dados Relacionais

Data Science Academy philipe.utsempreboni@outlook.com 5c8a62005e4cded1acb8b45a3

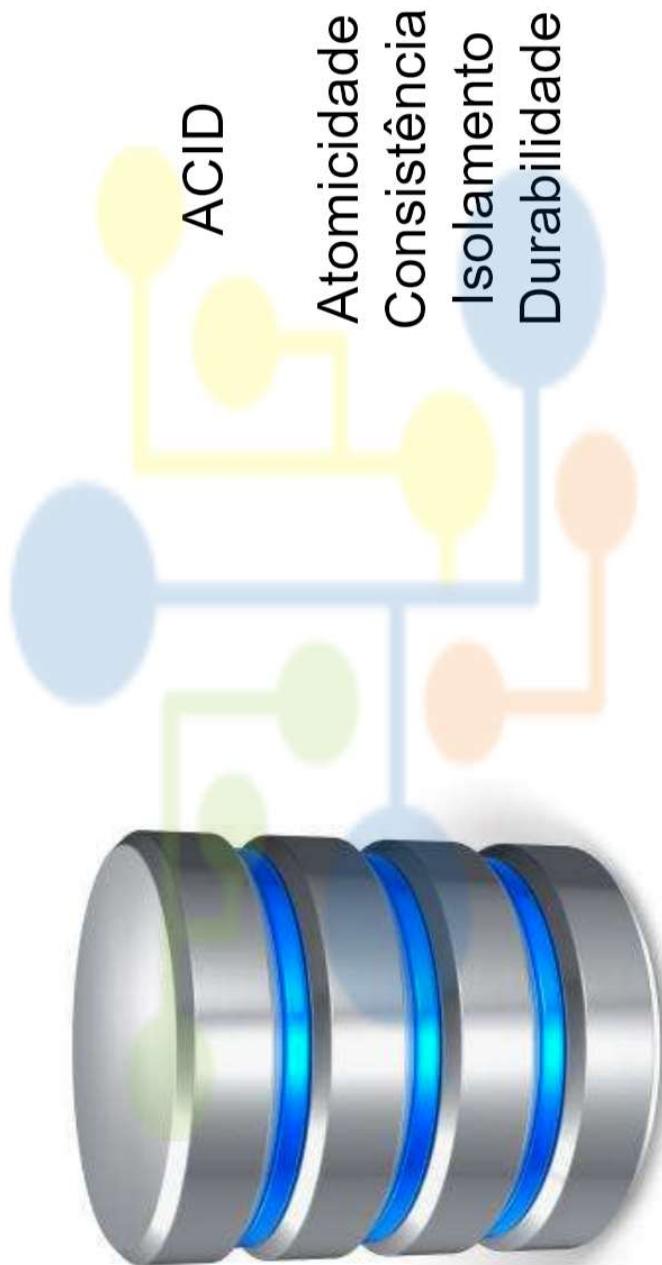
Data Science Academy



Hadoop x Bancos de Dados Relacionais



Data Science Academy philipe.utsempreboni@outlook.com 5c8aa62005e4cded1acb8b45a3³



Hadoop x Bancos de Dados Relacionais



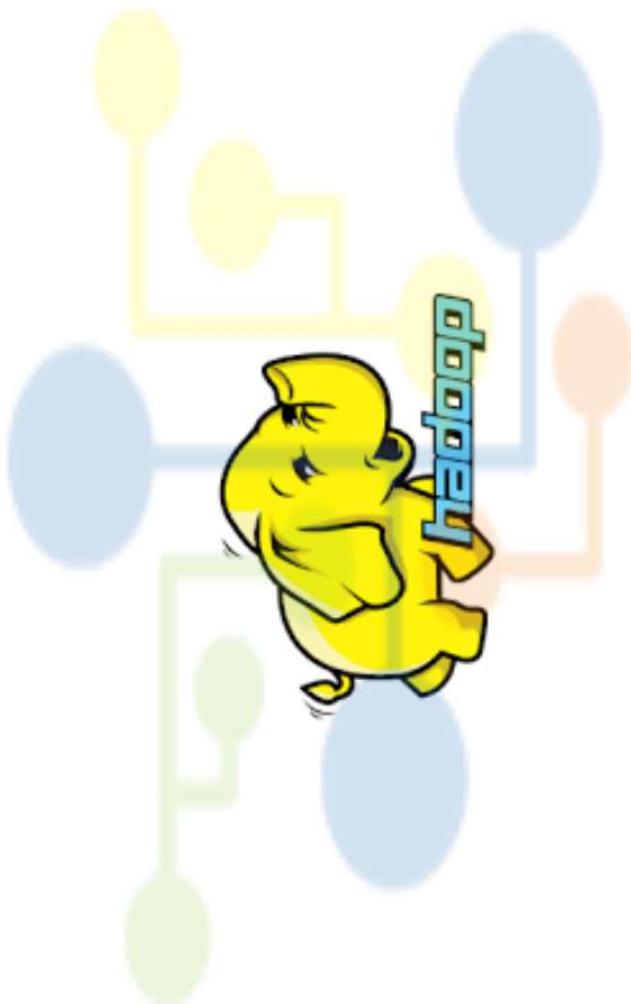
Data Science Academy philipe.utsempreboni@outlook.com 5c8aa62005e4cded1acdb8b45a3³



Hadoop x Bancos de Dados Relacionais



Data Science Academy philipe.utsempreboni@outlook.com 5c8aa62005e4cd1acd8b45a3

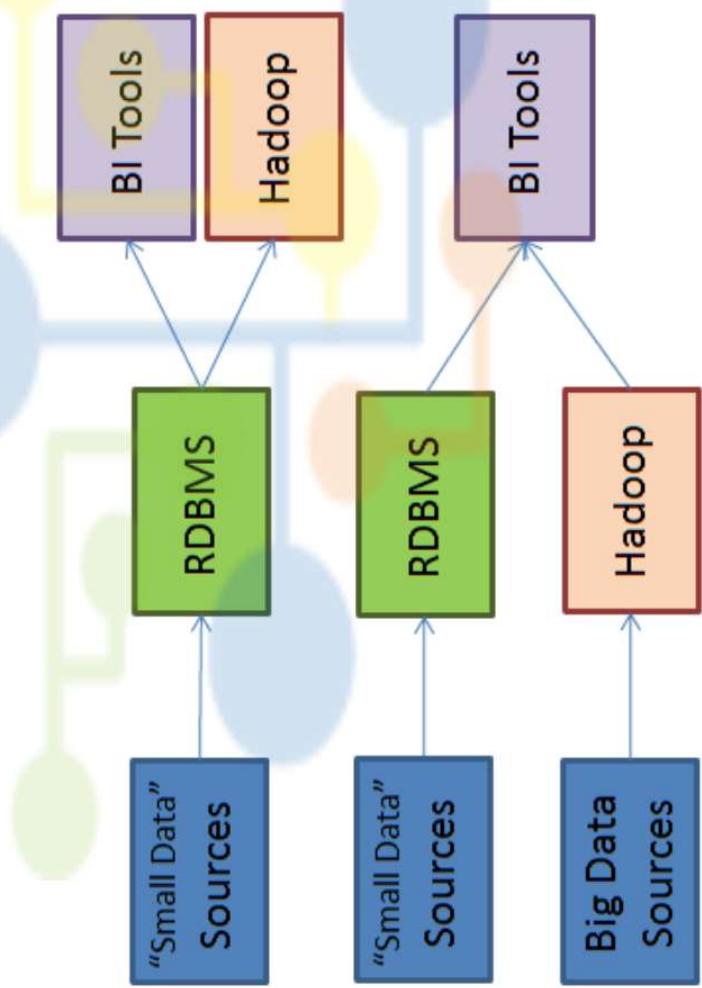


Hadoop x Bancos de Dados Relacionais



Data Science Academy philipe.utsempreboni@outlook.com 5c8aa62005e4cd1acdb8b45a3

Hadoop → Grandes volumes de dados (estruturados ou não estruturados)
RDBMS → Dados transacionais (dados estruturados)



Hadoop x Bancos de Dados Relacionais



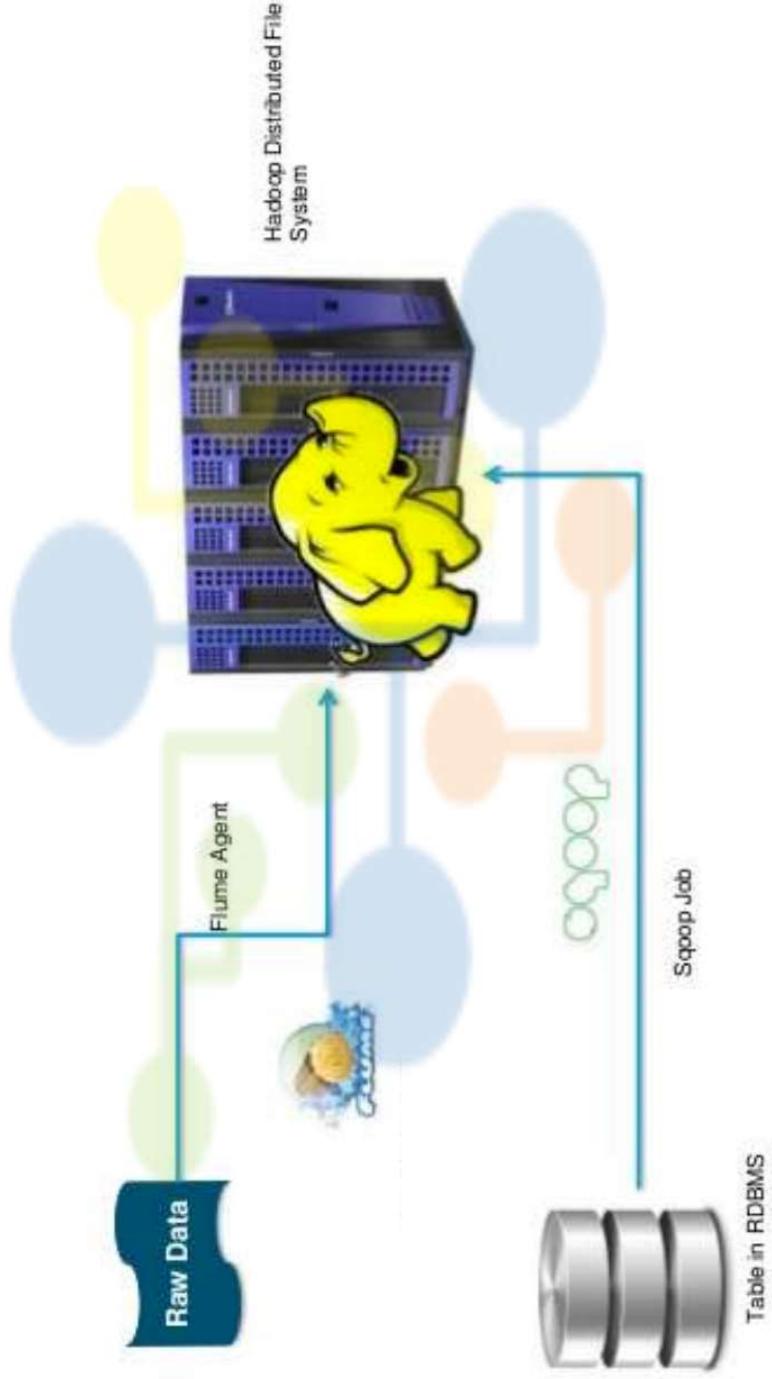
Data Science Academy philipe.utsempreboni@outlook.com 5c8aa62005e4cd1acdb8b45a3



Hadoop x Bancos de Dados Relacionais



Data Science Academy philipe.utsempreboni@outlook.com 5c8a62005e4cd1acd8b45a3



Hadoop x Bancos de Dados Relacionais



Data Science Academy philipe.utsempreboni@outlook.com 5c8a62005e4cd1acd8b45a3

Hadoop processa dados em batch. Consequentemente, ele não deve ser usado para processar dados transacionais. Mas o Hadoop pode resolver muitos outros tipos de problemas relacionados ao Big Data.



Obrigado
