

TI TOTAL

ÁREA FISCAL E CONTROLE



Professor
Ramon Souza

Tecnologia da Informação

RESUMO

Big Data

CONCEITOS DE BIG DATA

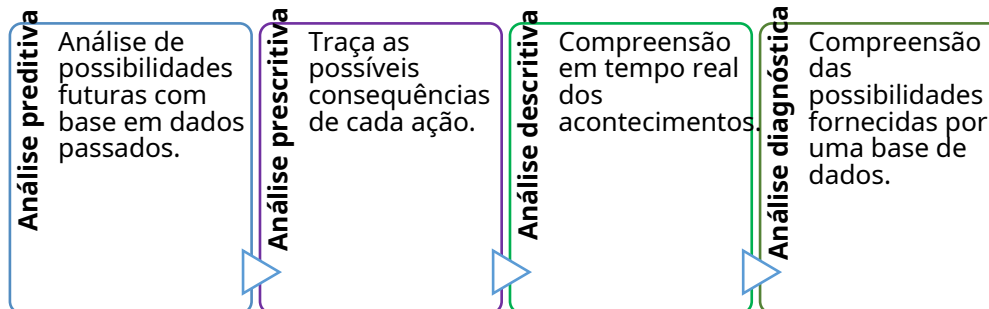
Big Data é um termo amplamente utilizado na atualidade para nomear **conjuntos de dados muito grandes ou complexos**, que os **aplicativos de processamento de dados tradicionais ainda não conseguem lidar**.

O **Big Data** pode ser definido genericamente como a **captura, gerenciamento e a análise de dados que vão além dos dados tipicamente estruturados**, que podem ser consultados e pesquisados através de bancos de dados relacionais.

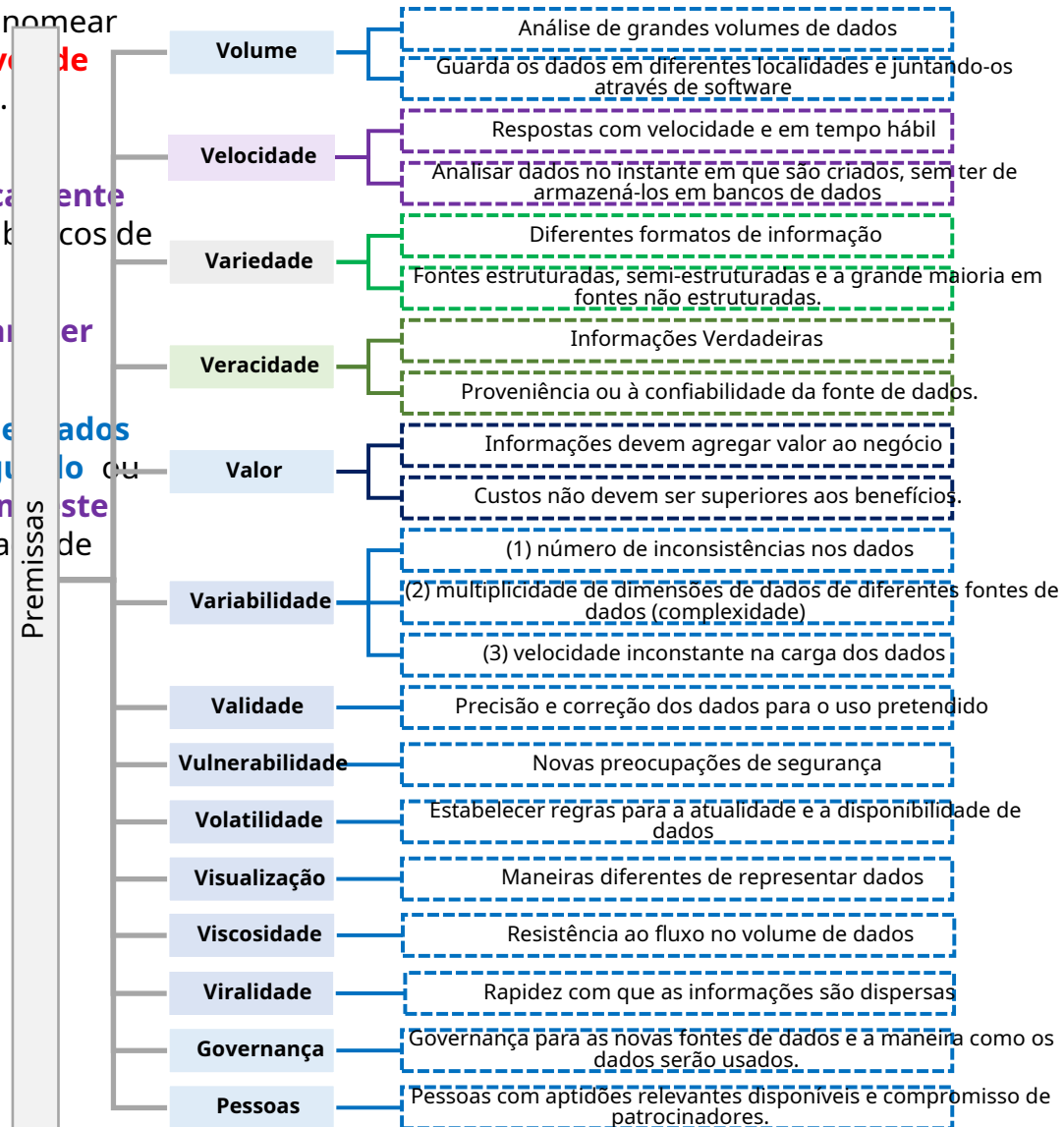
O **objetivo do Big Data** é **propiciar dados e informações que possam ser analisados visando subsidiar tomadas de decisão**.

O **Big Data** tanto pode ser encarado como o **grande volume de dados estruturados e não estruturados que são gerados a cada segundo** ou também como as **tecnologias que são utilizadas para lidar com este grande volume de dados**. Esta segunda acepção é, por vezes, chamada de **Big Data Analytics**.

TIPOS DE ANÁLISE

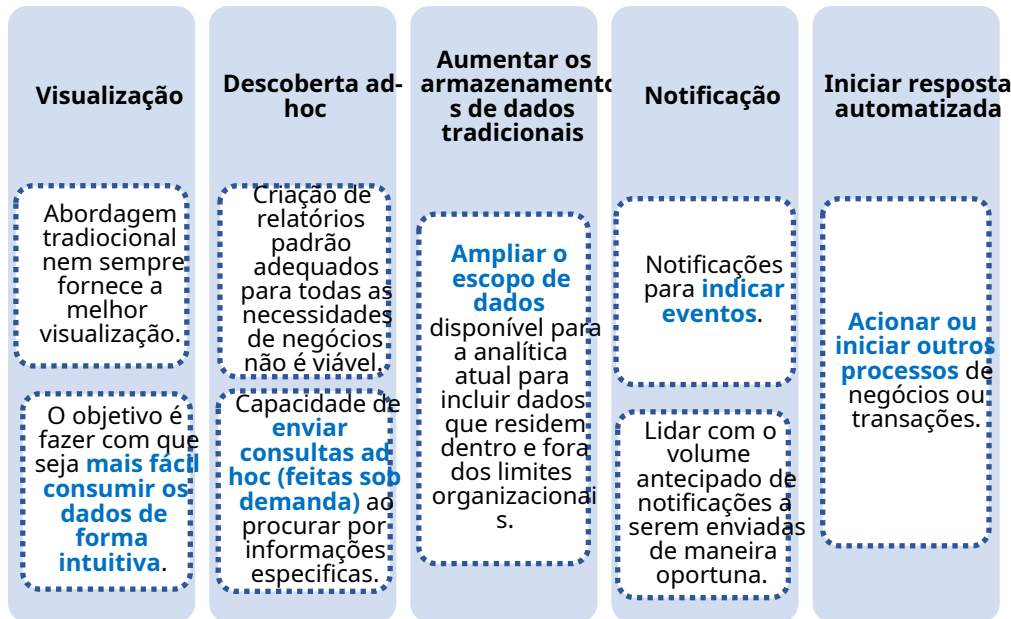


PREMISSAS DE BIG DATA

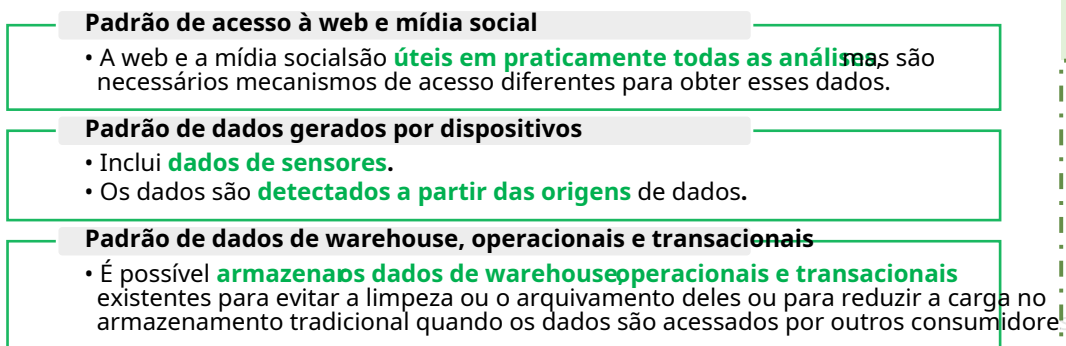


PADRÕES DE BIG DATA

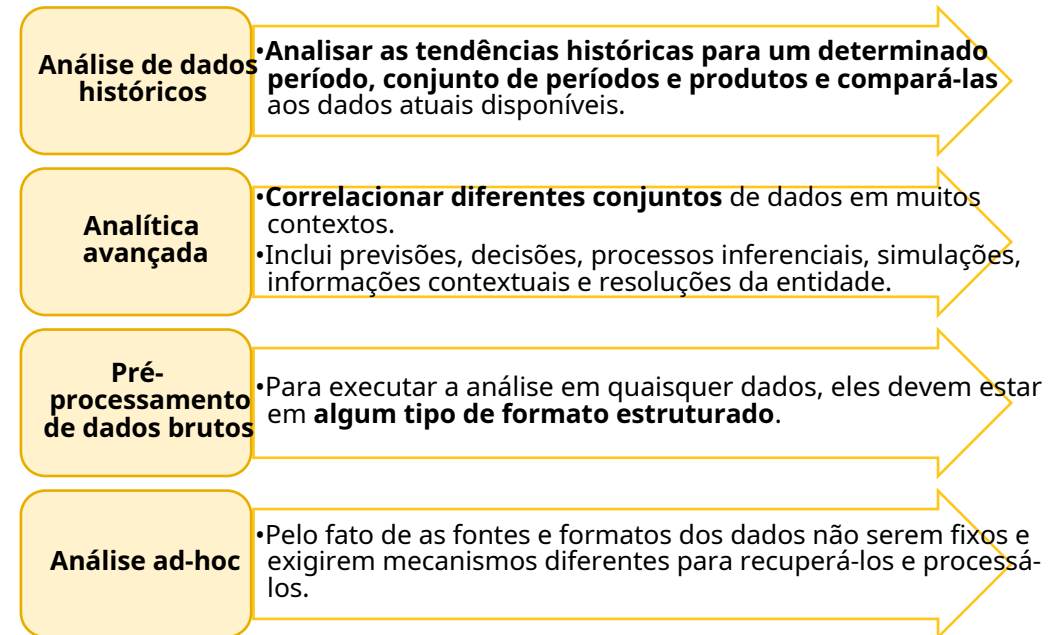
Padrões de consumo:



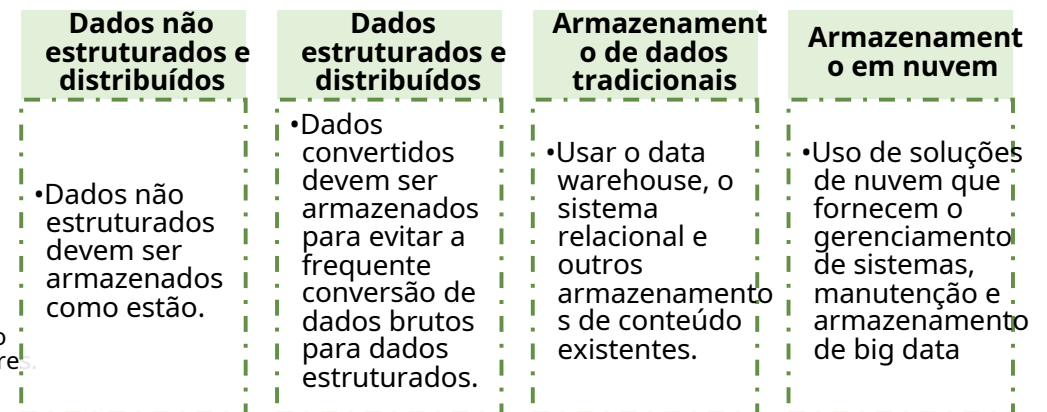
Padrões de acesso:



Padrões de processamento:



Padrões de armazenamento:

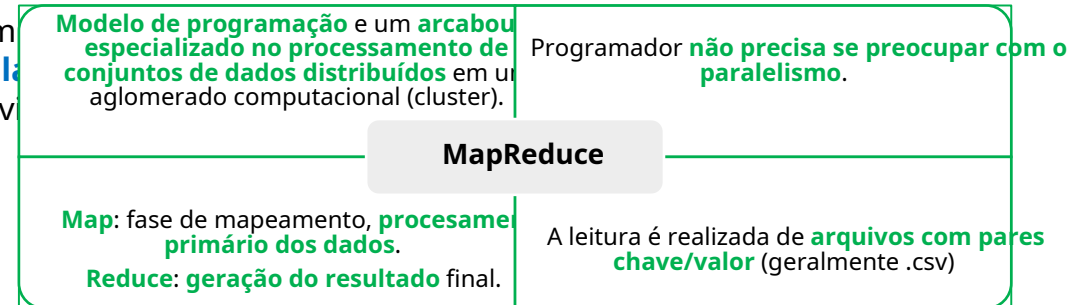


HADOOP

O **Hadoop** é um framework de código aberto, implementado em **utilizado para o processamento e armazenamento em larga escala** **alta demanda de dados**, **utilizando máquinas comuns**. Os servidores Hadoop fornecem armazenamento, processamento, acesso, governança, segurança e operações de dados.

O Hadoop é formado por vários subprojetos.

Subprojetos para processamento de dados



Subprojetos para gerenciamento

O **Chukwa** é o **sistema especialista em coleta e análise de logs em sistemas de larga escala**.

O **ZooKeeper** é o arcabouço criado pelo Yahoo! em 2007 com o objetivo de fornecer um **serviço de coordenação para aplicações distribuídas de alto desempenho**.

Subprojetos para acesso aos dados

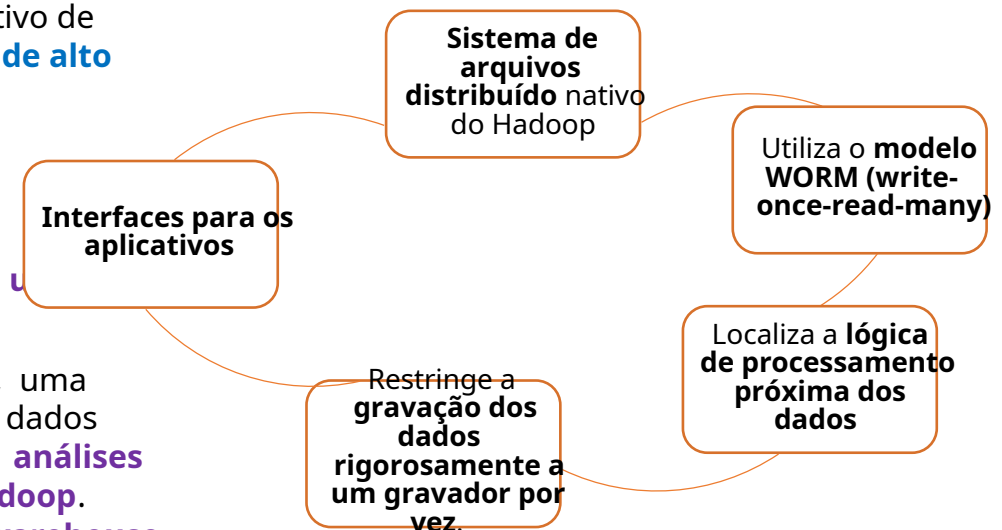
O **Pig** é uma **linguagem de alto nível orientada a fluxo de dados e um arcabouço de execução para computação paralela**.

O **Hive** fornece uma infraestrutura que permita utilizar Hive QL, uma linguagem de consulta similar a SQL bem como demais conceitos de dados relacionais tais como tabelas, colunas e linhas, para **facilitar as análises complexas feitas nos dados não relacionais de uma aplicação Hadoop**. Existe também uma definição que trata o **Hive** como um **datawarehouse distribuído** que facilita o uso de grandes conjuntos de dados. Nesse caso, ele seria enquadrado como um subprojeto para armazenamento dos dados.

O **Avro** é o **sistema de serialização de dados baseado em schemas**.

Subprojetos para armazenamento de dados

O **HDFS**:



O Hadoop é um **banco de dados distribuído e escalável que dá suporte a armazenamento estruturado e otimizado para grandes tabelas**.

Sqoop (SQL to Hadoop)

O **Sqoop (SQL to Hadoop)** é um aplicativo de interface de linha de comando para **transferência de dados entre bancos de dados relacionais e Hadoop**.

Flume

O **Flume** é um software distribuído, confiável e disponível para **coletar, agregar e mover com eficiência grandes quantidades de dados de log**.

Componentes do Hadoop

NameNode

- Gerenciar os arquivos armazenados no HDFS.

SecondaryNameNode

- Auxiliar o NameNode a manter seu serviço.
- Ser uma **alternativa de recuperação** no caso de uma falha do NameNode.

DataNode

- Efetivamente realizam o **armazenamento dos dados**.

JobTracker

- Função de **gerenciamento sobre o plano de execução das tarefas** a serem processadas pelo MapReduce.

TaskTracker

- Responsável pela **execução de tarefas MapReduce**.