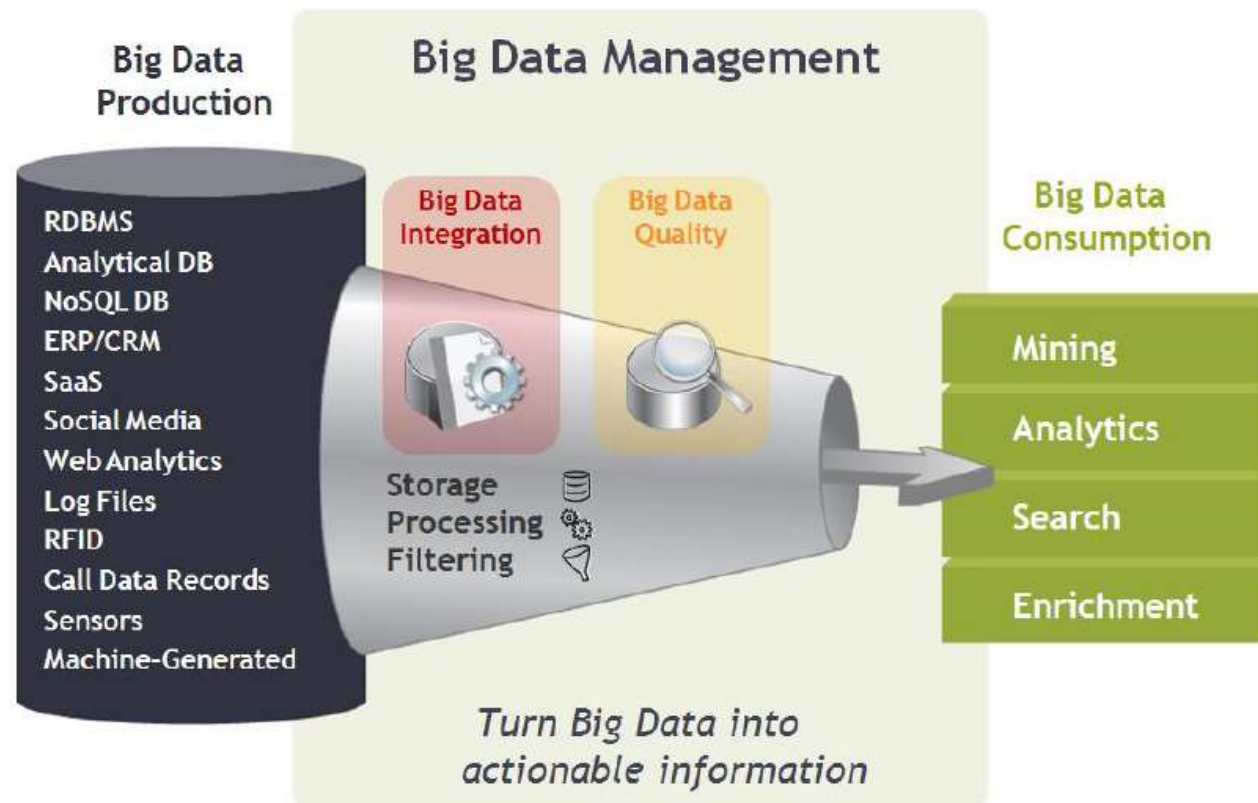# Module 6

## Introduction to Apache Hadoop and HDFS

Cristóbal Barba González– University of Málaga
Antonio J. Nebro – University of Málaga
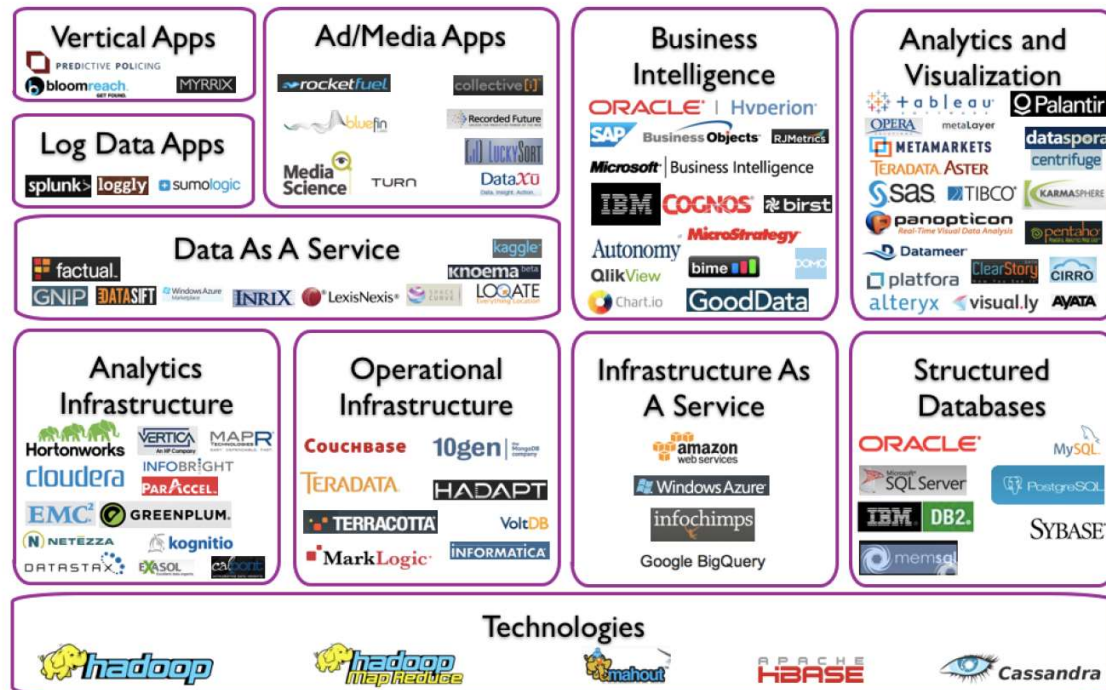
# Table of contents

- Big Data.

- Introduction to Apache Hadoop.

- Hadoop distribution.

- Hadoop cloud.

Máster de Formación Permanente
en **BIG DATA**
e Inteligencia Artificial

Master en
Big Data e Inteligencia Artificial

kha0s
R E S E A R C H

# Big Data Management

Máster de Formación Permanente
en **BIG DATA**
e Inteligencia Artificial

Master en
# Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# Big Data Landscape
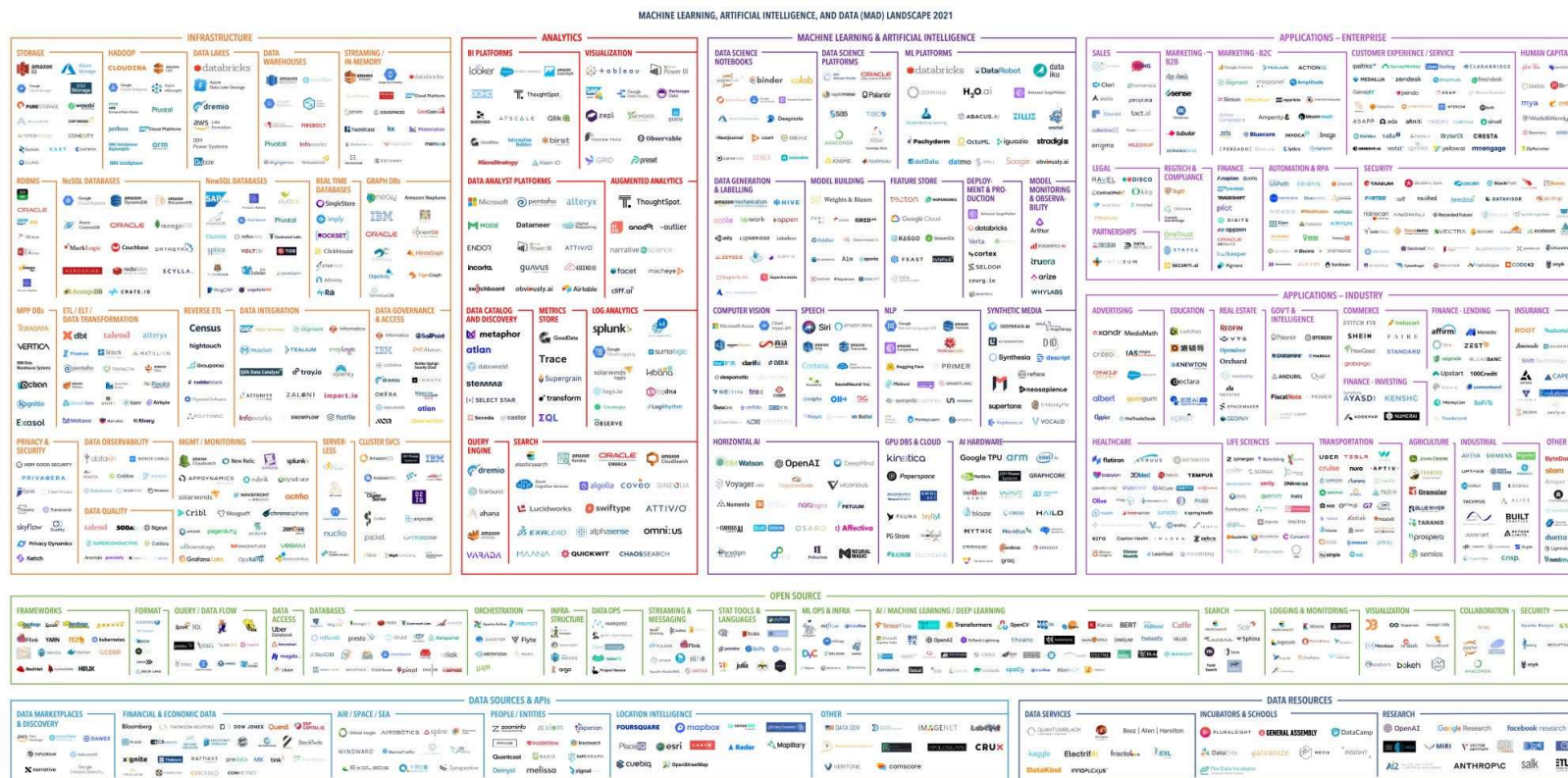
# Big Data Landscape



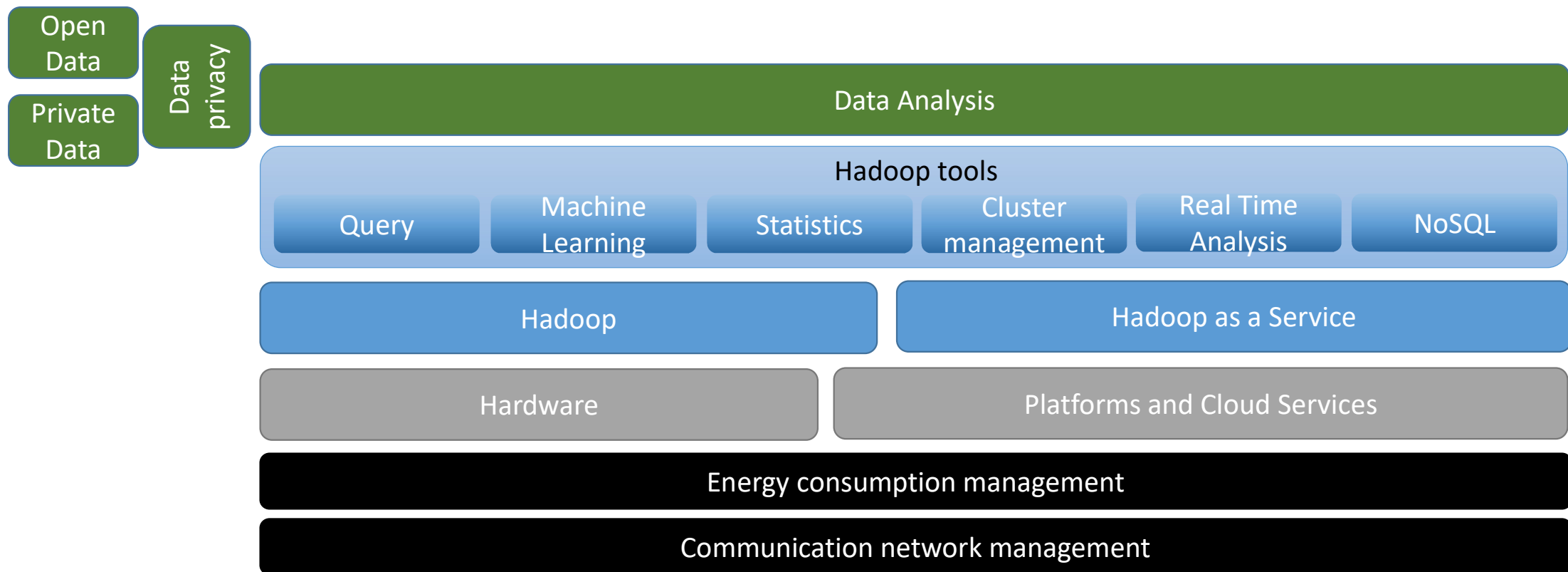© Matt Turck (@mattturck) and ShivonZilis (@shivonz) Bloomberg Ventures

# Big Data Landscape



MACHINE LEARNING, ARTIFICIAL INTELLIGENCE, AND DATA (MAD) LANDSCAPE 2021
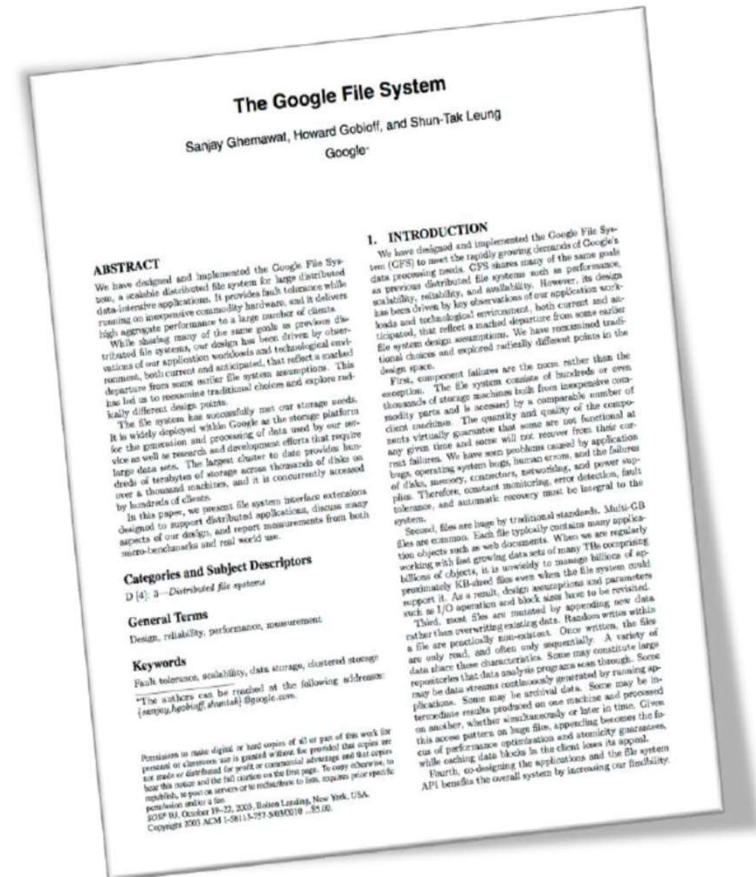
# Global Vision

# Big Data Business Intelligence

# Massive Data Storage

- On a PC, we store data on our hard drive.

- The data can be of any type:
  - Tables, text files, web pages…
  - …images, videos,

- BIG DATA: What happens when the data does not fit on a hard drive?

Máster de Formación Permanente
en **BIG DATA** e Inteligencia Artificial

Master en
# Big Data e Inteligencia Artificial
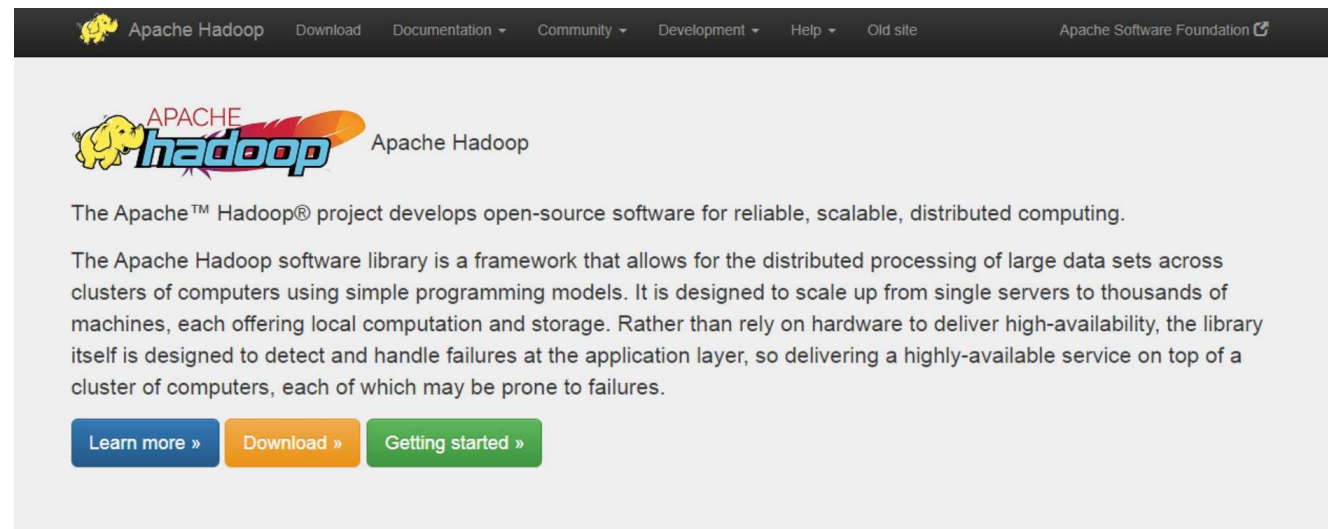
kha⊙s
R E S E A R C H

# The Google Approach

- In 2003, Google presents Google file system

- It is a distributed file system. It allows:
  - Store large amounts of data.
  - Distribute (partition) the data among several computers.
  - Allow more efficient file readings.
  - Prevent data loss if a computer crashes spoils.

Máster de Formación Permanente
en **BIG DATA**
e Inteligencia Artificial

Master en
Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# Apache Hadoop

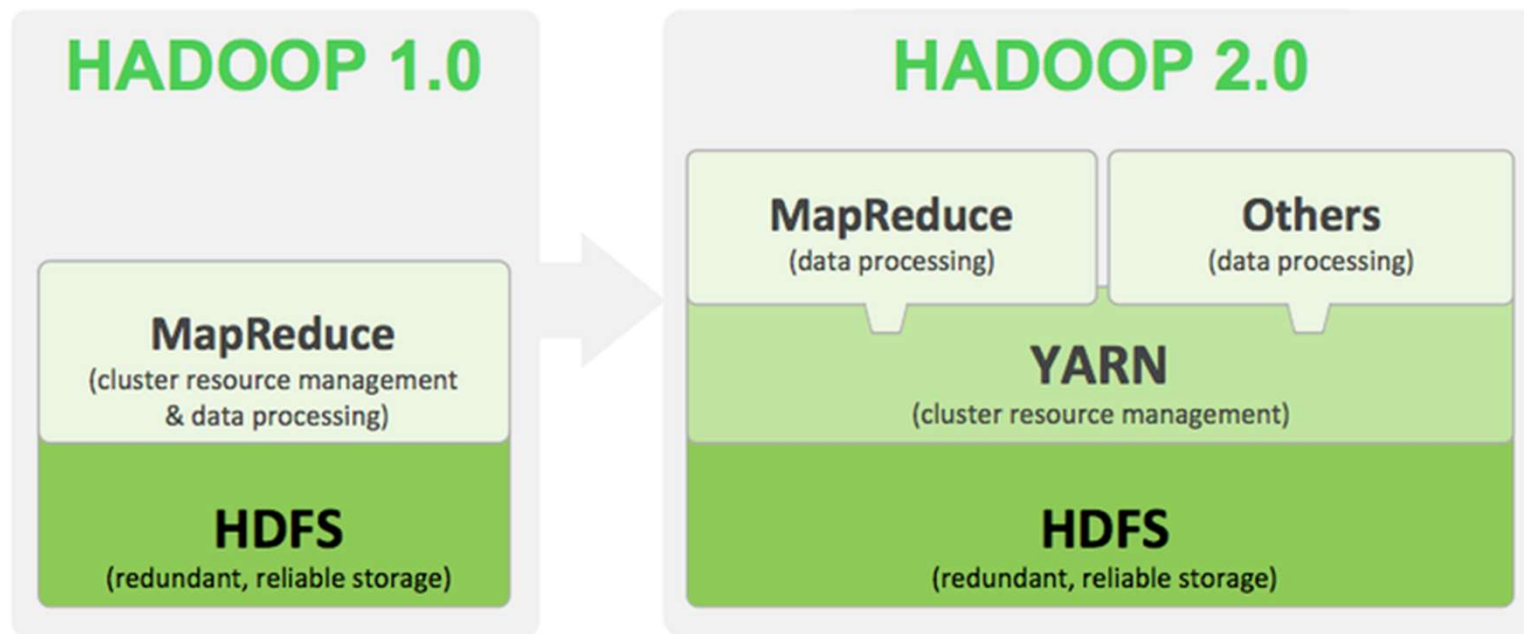- **Hadoop** is a framework that allows for the distributed processing of large data sets across clusters.

  - Performance.
  - Storage.
  - Scalability.
  - Fault tolerance.
  - Cost efficiency.



https://hadoop.apache.org/

# Apache Hadoop

Máster de Formación Permanente
en **BIG DATA**
e Inteligencia Artificial

Master en
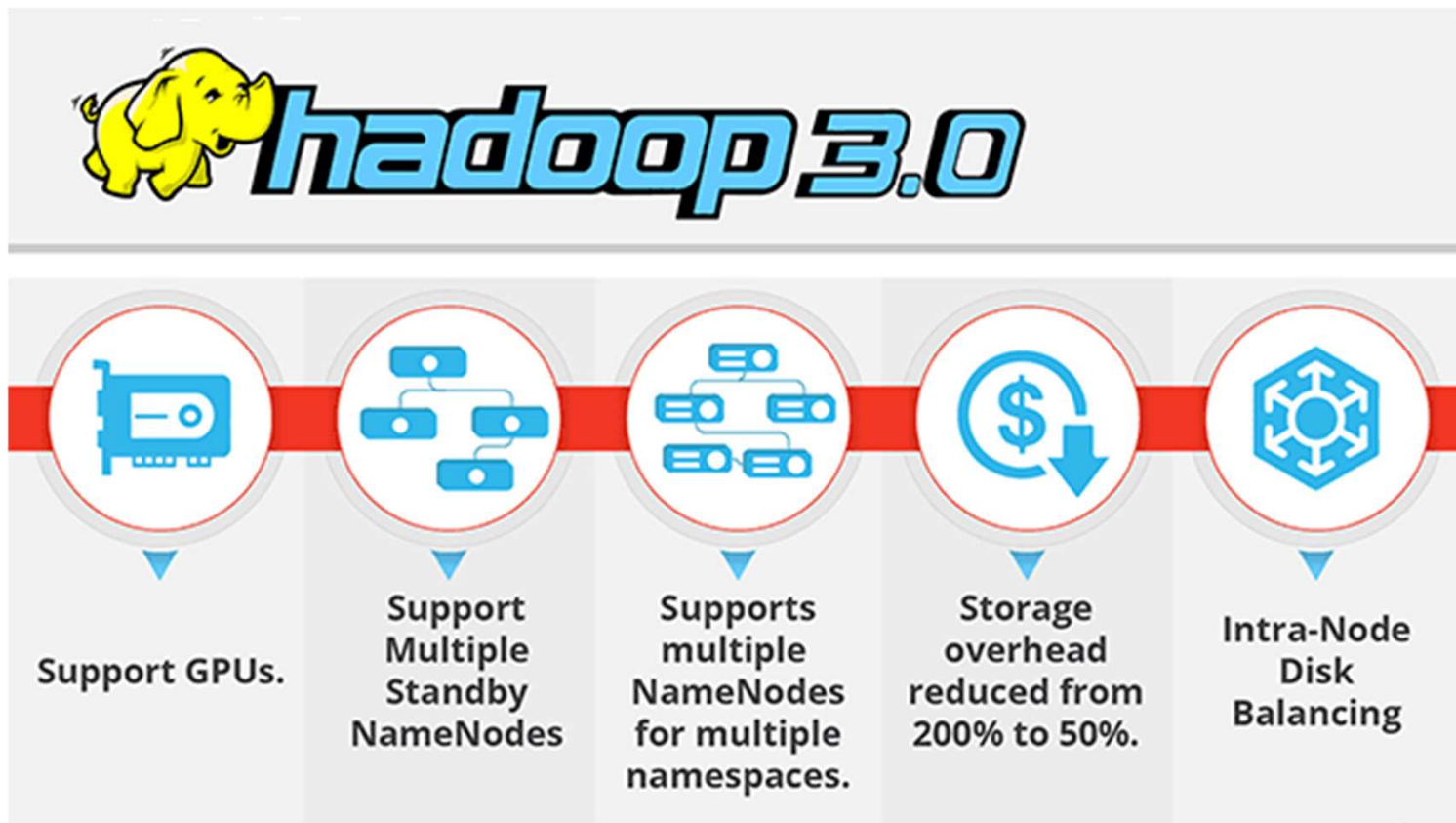Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# Apache Hadoop

# Apache Hadoop

- When to use Hadoop?
  - For processing really Big Data.
  - For storing a diverse set of data.
  - For parallel data processing.

- When not to use Hadoop?
  - For real-time data analysis.
  - For replacement relational database system.
  - For multiple smaller datasets.

# What is Apache Hadoop used for?

## Analytics and big data

A wide variety of companies and organizations use Hadoop for research, production data processing, and analytics that require processing terabytes or petabytes of big data, storing diverse datasets, and data parallel processing.

## Data storage and archiving

As Hadoop enables mass storage on commodity hardware, it is useful as a low-cost storage option for all kinds of data, such as transactions, click streams, or sensor and machine data.

## Data lakes

Since Hadoop can help store data without preprocessing, it can be used to complement to data lakes, where large amounts of unrefined data are stored.

## Marketing analytics

Marketing departments often use Hadoop to store and analyze customer relationship management (CRM) data.

## Risk management

Banks, insurance companies, and other financial services companies use Hadoop to build risk analysis and management models.
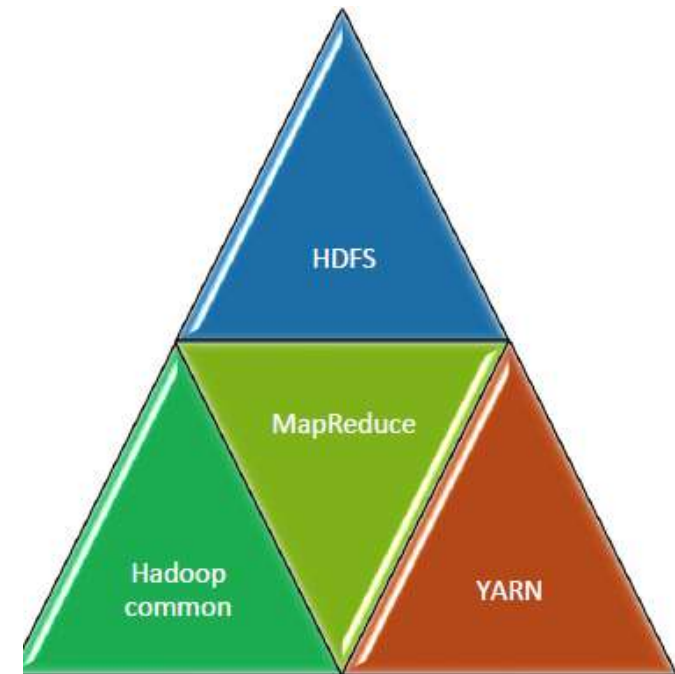
## AI and machine learning

Hadoop ecosystems help with the processing of data and model training operations for machine learning applications.

# Apache Hadoop

- Basic components of Hadoop
  - **Hadoop HDFS** - Hadoop Distributed File System (HDFS) is the storage unit of Hadoop.
  - **Hadoop MapReduce** - Hadoop MapReduce is the processing unit of Hadoop.
  - **Hadoop YARN** - Hadoop YARN is a resource management unit of Hadoop.
  - **Hadoop Common –** Haddop common is a collection of common libraries and utilities that work with different Hadoop modules.

Máster de Formación Permanente
en **BIG DATA**
e Inteligencia Artificial

Master en
Big Data e Inteligencia Artificial

kha○s
R E S E A R C H

# Hadoop HDFS

- HDFS has a master/worker(slave) architecture. An HDFS cluster consists of one or more **NameNodes**

  - NameNode is a master server that manages the file system namespace and regulates access to files by clients Simple Coherency Model.

  - HDFS applications need a write-once-read-many access model for files. A file once created, written, and closed need not be changed.

  - The NameNode executes file system namespace operations like opening, closing, and renaming files and directories. It also determines the mapping of blocks to DataNodes and their replications.

Máster de Formación Permanente
en **BIG DATA**
e **Inteligencia Artificial**

Master en
Big Data e Inteligencia Artificial
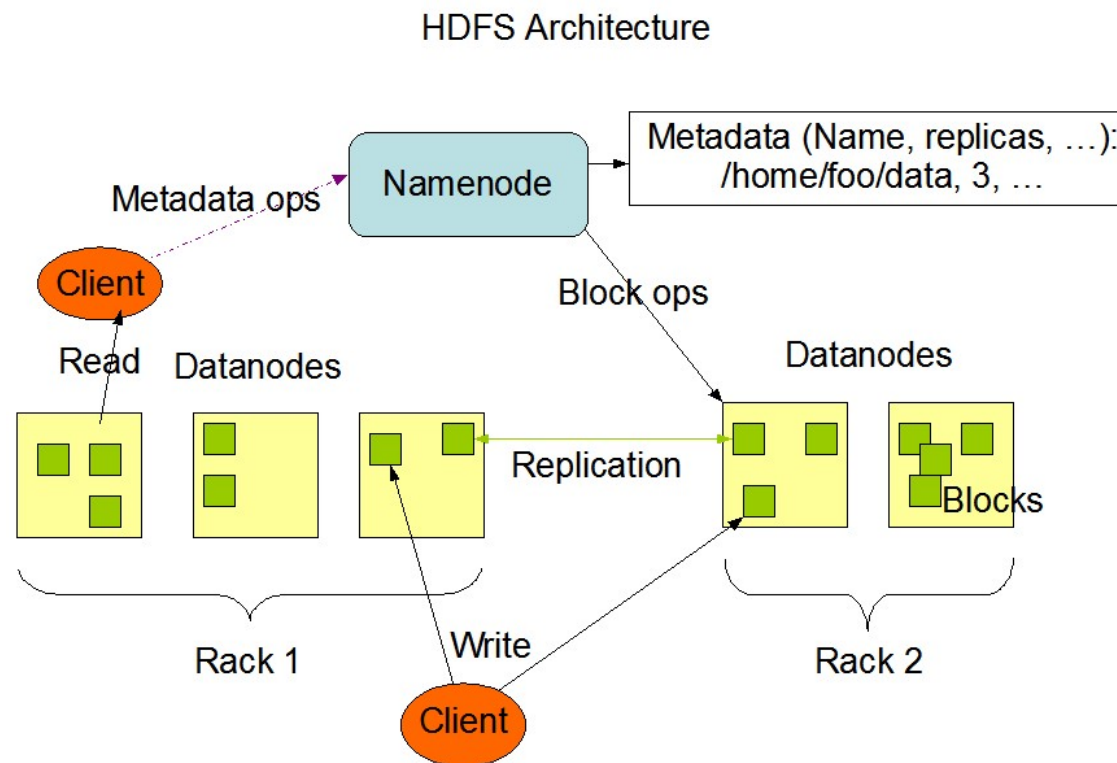
khaos
R E S E A R C H

# Hadoop HDFS

- There are a number of **DataNodes**, usually one per node in the cluster.
  - DataNodes manage storage attached to the nodes that they run on. HDFS exposes a file system namespace and allows user data to be stored in files.
  - The DataNodes are responsible for serving read and write requests from the file system's clients.
  - The DataNodes also perform block creation, deletion, and replication upon instruction from the NameNode.

# Hadoop HDFS

- Main features
  - Based on commodity hardware (not special hardware is required).
  - Fault tolerant and self healing.
  - Developed specifically for large scale data processing workloads where scalability, flexibility and throughput are critical.
  - Load Balancing: Place data intelligently for maximum efficiency and utilization.
  - Tunable Replication: Multiple copies of each file provide data protection and computational performance.
  - Based on "Moving Computation is Cheaper than Moving Data".
  - Master/ slave architecture.
  - Written in Java.
  - Implemented in Linux.
  - TCP communication protocol.

Master en
Big Data e Inteligencia Artificial

Máster de Formación Permanente
en BIG DATA
e Inteligencia Artificial

khaos
R E S E A R C H

# Hadoop HDFS

Máster de Formación Permanente
en BIG DATA
e Inteligencia Artificial

Master en
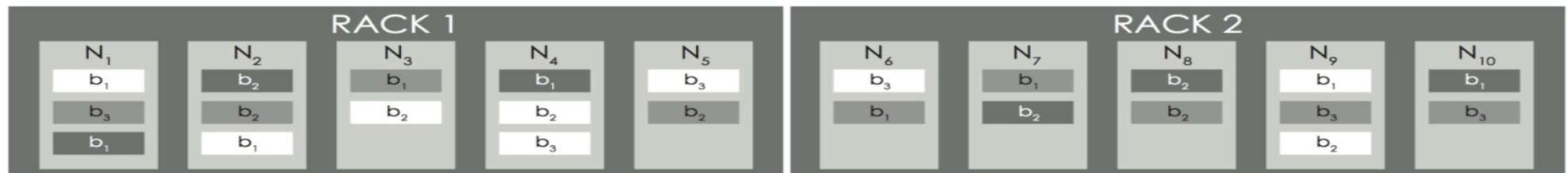Big Data e Inteligencia Artificial

khaos
RESEARCH

# Hadoop HDFS

- Each file is divided into blocks of 64 or 128 MB.
- The blocks are stored divided between several computers.
- Each computer is called a node.
- The set of computers is called a cluster.
- ADVANTAGE 1: If you double the number of nodes, you can fit twice as much data.
- ADVANTAGE 2: You can store files that do not fit on a computer.
- ADVANTAGE 3: Readings are faster: you can read a file from several computers at the same time.

# Hadoop HDFS

- PROBLEM: If there are many nodes, the probability that one fails is high.
- SOLUTION: Each block is stored in several nodes, by default 3.
- Each block stored on a node is called a replica.
- Replicas are topology-aware: better to store them in different racks* or data centers.



| RACK 1 | | | | | RACK 2 | | | | |
|---|---|---|---|---|---|---|---|---|---|
| $N_1$ | $N_2$ | $N_3$ | $N_4$ | $N_5$ | $N_6$ | $N_7$ | $N_8$ | $N_9$ | $N_{10}$ |
| $b_1$ | $b_2$ | $b_1$ | $b_1$ | $b_3$ | $b_3$ | $b_1$ | $b_2$ | $b_1$ | $b_1$ |
| $b_3$ | $b_2$ | $b_2$ | $b_2$ | $b_2$ | $b_1$ | $b_2$ | $b_2$ | $b_3$ | $b_3$ |
| $b_1$ | $b_1$ | | $b_3$ | | | | | $b_2$ | |

* A rack or cabinet is a container where several nodes are stored that share a power and network system

Máster de Formación Permanente
en **BIG DATA** e Inteligencia Artificial

Master en
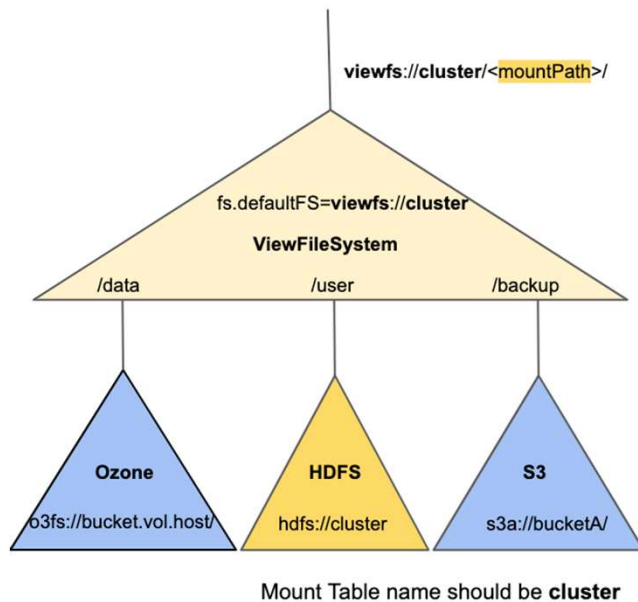Big Data e Inteligencia Artificial

khaos
RESEARCH

# Hadoop HDFS

- The File System (FS) shell includes various shell-like commands that directly interact with the Hadoop Distributed File System (HDFS) as well as other file systems that Hadoop supports, such as Local FS, WebHDFS, S3 FS, and others
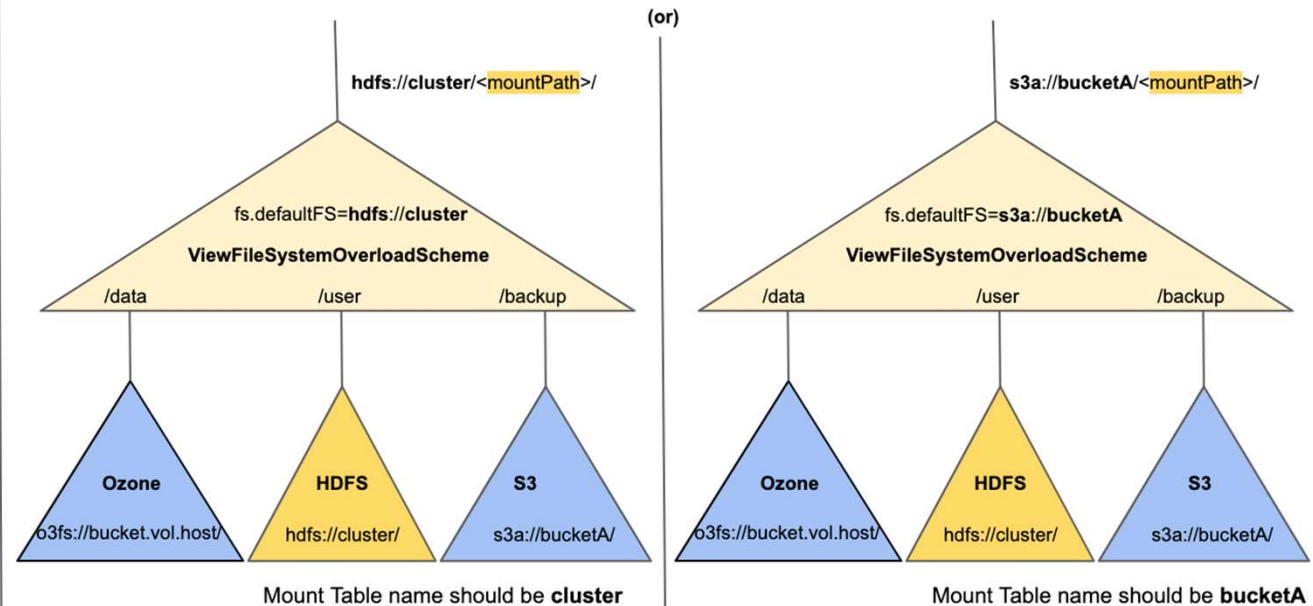
  https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/FileSystemShell.html

Máster de Formación Permanente
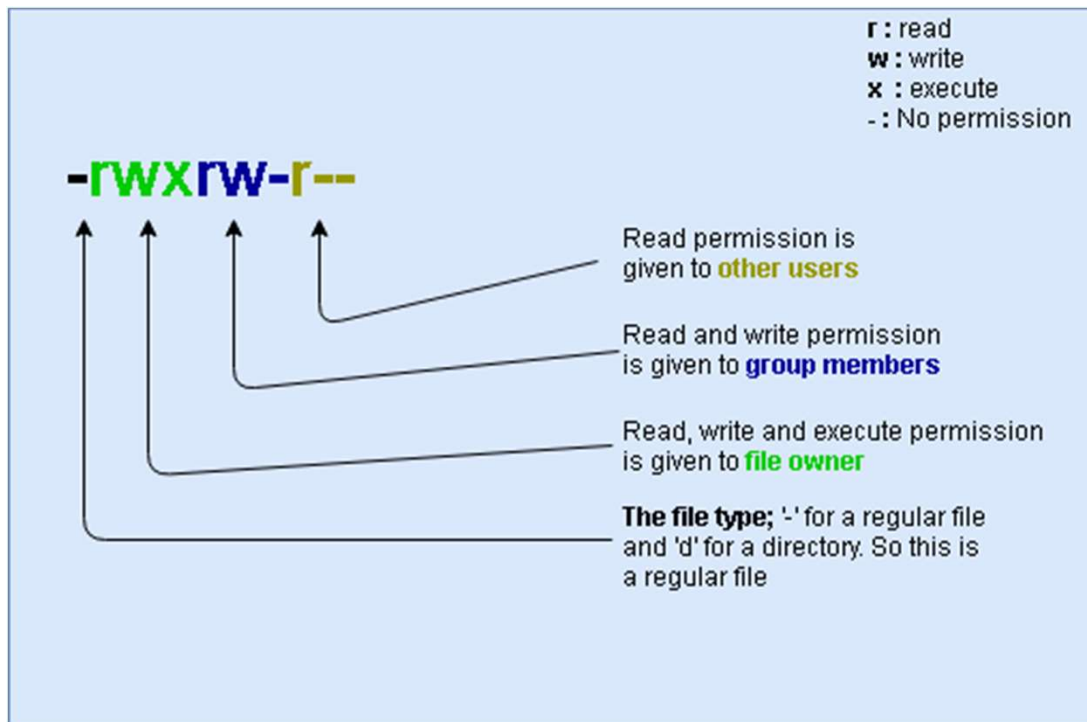en **BIG DATA**
e Inteligencia Artificial

Master en
Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# Hadoop with S3

Máster de Formación Permanente
en **BIG DATA**
e **Inteligencia Artificial**

Master en
Big Data e Inteligencia Artificial

kha⦿s
R E S E A R C H

# Linux: Commands



**r** : read
**w** : write
**x** : execute
**-** : No permission

-**rwx****rw**-**r**--

Read permission is given to **other users**

Read and write permission is given to **group members**

Read, write and execute permission is given to **file owner**

**The file type;** '-' for a regular file and 'd' for a directory. So this is a regular file

- ls (list directory)
- mkdir (make a directory)
- touch (touch file)
- mv (move file)
- cp (copy files)
- rm (remove file)

# Hadoop Distributed File System (HDFS)

- hadoop fs <args>
  - fs is used for generic file system, and it can point to any file system such as local file system, HDFS, WebHDFS, S3 FS, etc.

- hdfs dfs <args>
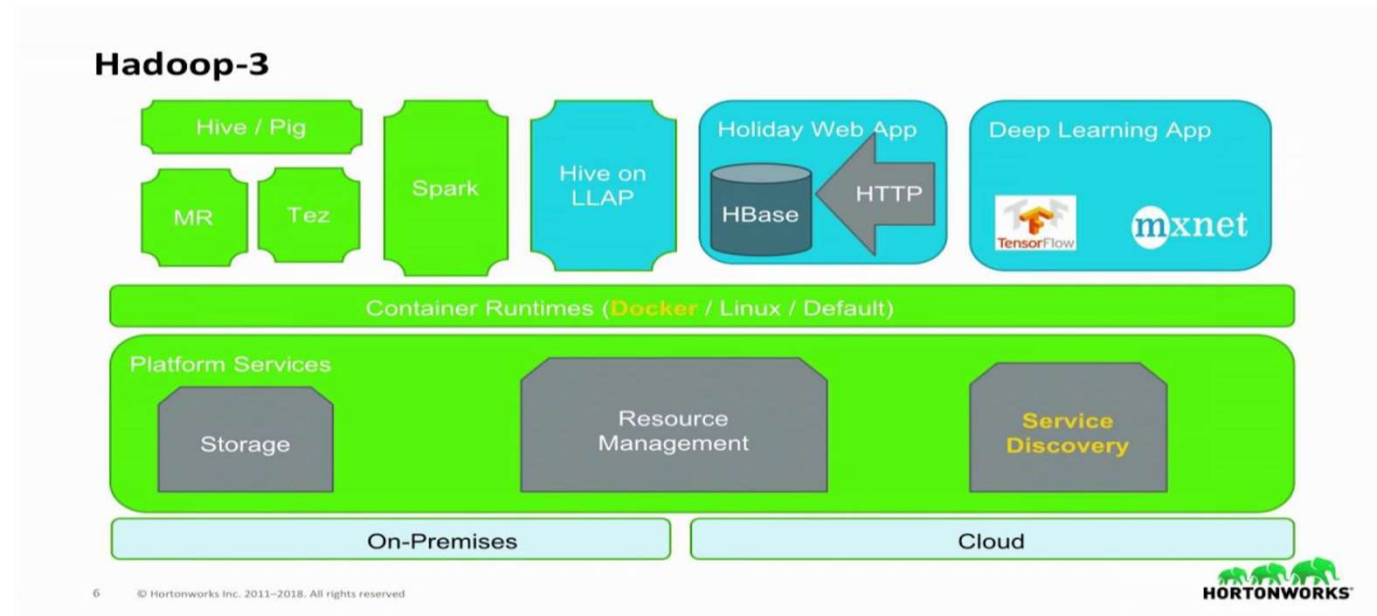  - dfs points to the Distributed File System and it is specific to HDFS

Máster de Formación Permanente
en **BIG DATA**
e Inteligencia Artificial

Master en
Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# Hadoop HDFS: Commands

| Command | Advantages |
|---|---|
| hdfs dfs –put localFile hdfsFile | Put a local file in hdfs |
| hdfs dfs –mkdir hdfsDir | Create a new directory in hdfs |
| hdfs  dfs –cat hdfsFile | Display the content of a file |
| hdfs dfs –rm hdfsFile | Delete a file |
| hdfs dfs | Help |

https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/HDFSCommands.html

* In old versions of Hadoop : hadoop dfs cat …

Máster de Formación Permanente
en BIG DATA
e Inteligencia Artificial

Master en
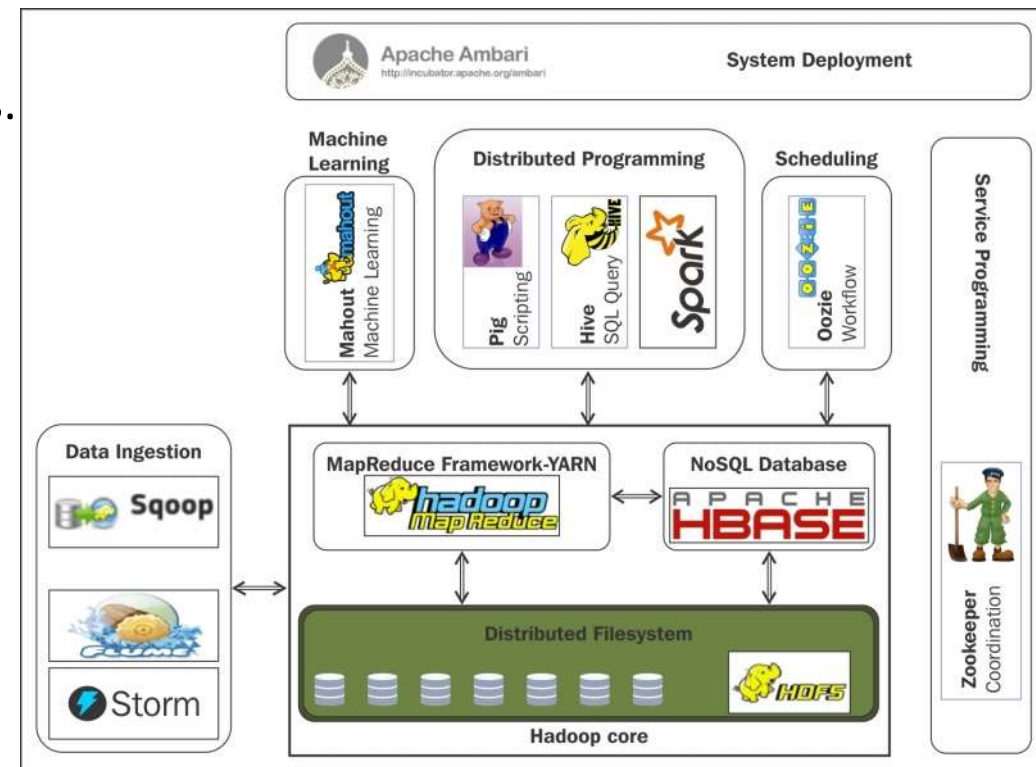Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# YARN: Yet another resource negotiator

- YARN — Yet Another Resource Negotiator, is a part of Hadoop, is one of the two major components of Apache Hadoop (with HDFS).
- YARN also allows different data processing engines like graph processing, interactive processing, stream processing as well as batch processing to run and process data stored in HDFS.

Máster de Formación Permanente
en **BIG DATA**
e Inteligencia Artificial

Master en
# Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# Apache Hadoop Ecosystem

- Hbase: A scalable, distributed database that supports structured data storage for large tables.
- Hive: A data warehouse infrastructure that provides data summarization and ad hoc querying.
- Pig: A high-level data flow language and execution framework for parallel computation.
- Zookeeper: A high-performance coordination service for distributed applications.
- Mahout: A Scalable machine learning and data mining library.
- …

Máster de Formación Permanente
en **BIG DATA**
e **Inteligencia Artificial**

Master en
Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# Big Data processing

- Two main paradigms:
  - Batch.
  - In real time (Streaming).

- Batch processing
  - Large volumes of historical data are analyzed.
  - Processing can take minutes, hours, days...
- Real-time processing (Streaming)
  - Large volumes of data are analyzed that reach high speed.
  - Processing must be immediated.

# Google's Solution: MapReduce

- In 2004, Google announces MapReduce.
- It is a framework for data processing:
  - Process data stored over GFS.
  - Do it in a distributed and parallel way.
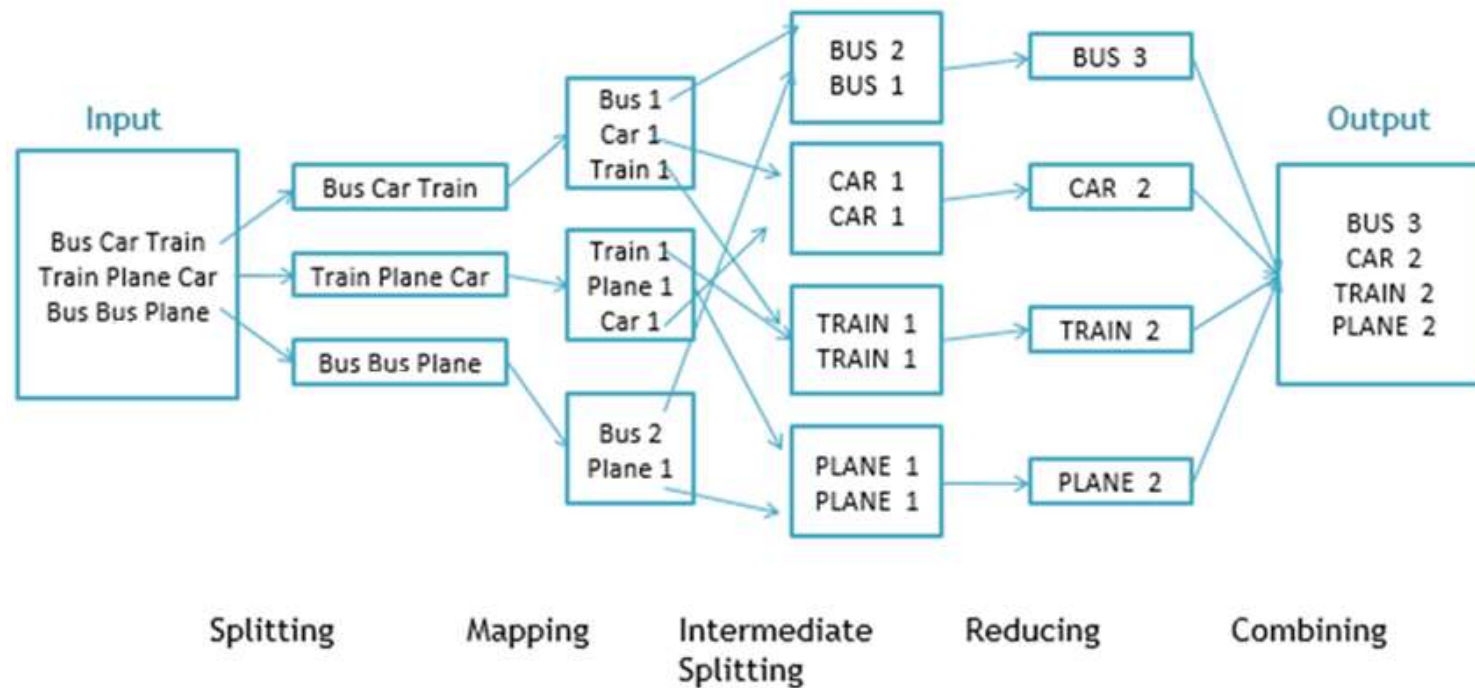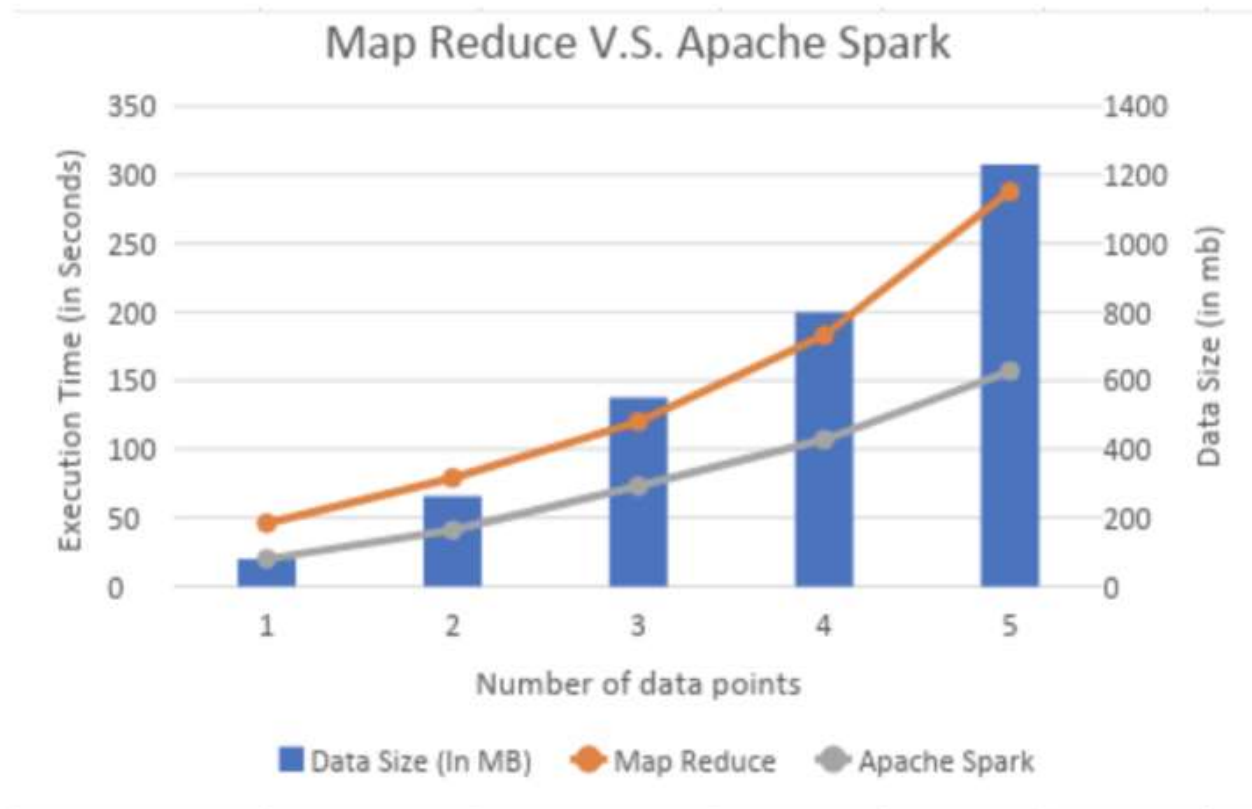  - Avoid data loss if a computer crashes.

Máster de Formación Permanente
en **BIG DATA**
e Inteligencia Artificial

Master en
# Big Data e Inteligencia Artificial

kha s
R E S E A R C H

# Hadoop MapReduce

- In 2005, an open-source version of MapReduce is included in HADOOP.

- It consists of two main functions.

- The **map** function performs a transformation of the data.
  - INPUT: a key value tuple.
  - OUTPUT: one or more key-value tuples.

- An intermediate phase called **shuffle** is responsible for grouping the output of the map by keys.

- The **reduce** function performs a grouping or aggregation of the data:
  - INPUT: A key, along with all its associated values.
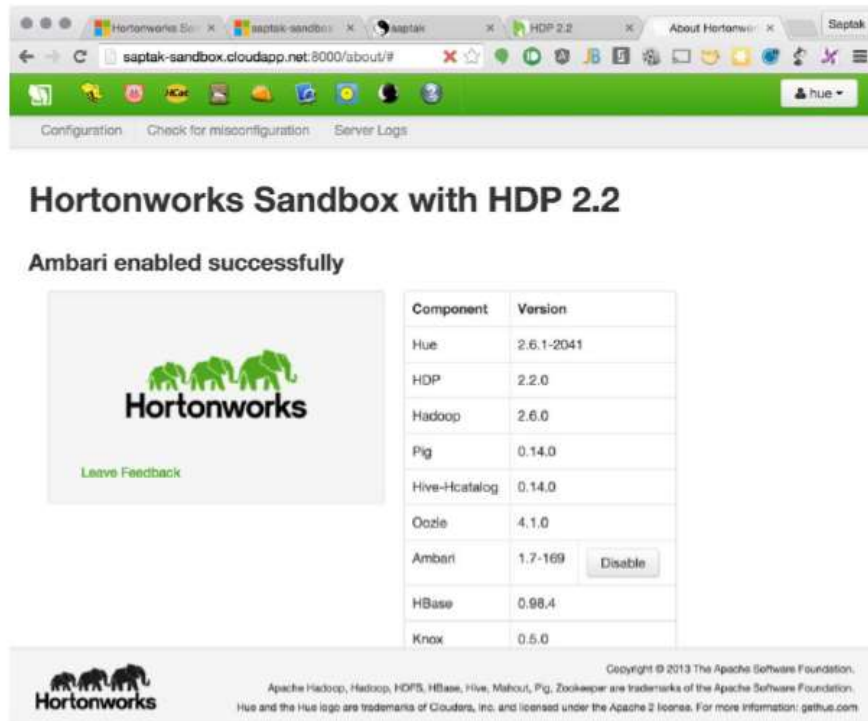  - OUTPUT: a value.

Máster de Formación Permanente
en BIG DATA
e Inteligencia Artificial

Master en
Big Data e Inteligencia Artificial

khaos
RESEARCH

# Hadoop MapReduce



Splitting    Mapping    Intermediate Splitting    Reducing    Combining

# Hadoop MapReduce



Map Reduce V.S. Apache Spark

# Hadoop distribution (Virtual Machines & Docker)

- Cloudera ( https://www.cloudera.com/)



- Hortonworks ( https://es.hortonworks.com/)



- MapR (https://mapr.com/)

Máster de Formación Permanente

en **BIG DATA**
e Inteligencia Artificial

Master en
# Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# Hadoop distribution (Virtual Machines & Docker)

Máster de Formación Permanente
en **BIG DATA**
e Inteligencia Artificial

Master en
Big Data e Inteligencia Artificial

khaos
R E S E A R C H

# Hadoop distribution (Hortonworks)

## DOCKER

Download HDP_3.0.1_docker-deploy-scripts.zip from
https://archive.cloudera.com/hwx-sandbox/hdp/hdp-3.0.1/HDP_3.0.1_docker-deploy-scripts_18120587fc7fb.zip
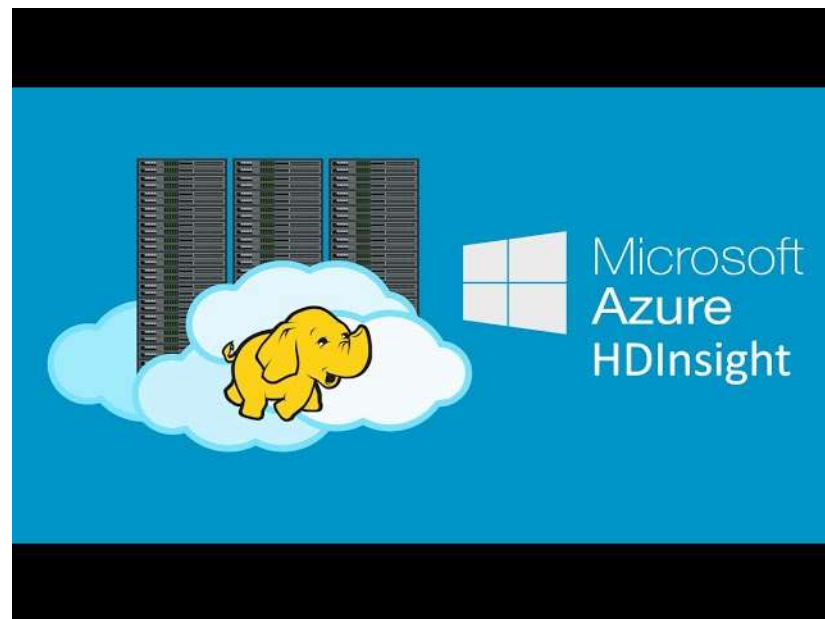bash docker-deploy-hdp30.sh

http://127.0.0.1:8080
admin/admin

docker stop sandbox-hdp
docker stop sandbox-proxy

docker start sandbox-hdp
docker start sandbox-proxy

## Máquina virtual VirtualBOX o VMWare

https://archive.cloudera.com/hwx-sandbox/hdp/hdp-3.0.1/HDP_3.0.1_virtualbox_181205.ova

# Hadoop cloud



https://azure.microsoft.com/es-es/products/hdinsight



https://aws.amazon.com/es/emr/features/hadoop/



https://cloud.google.com/architecture/hadoop?hl=es-419